

Foundations of Machine Learning

E. Le Pennec



MAP553 - Foundations of Machine Learning - Spring 2023

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

1 Statistical Learning: Introduction, Setting and Risk Estimation

- Introduction
- Machine Learning
- Supervised Learning
- Risk Estimation and Model Selection
- Cross Validation and Test
- References

2 ML Methods: Probabilistic Point of View

- Motivation
- Supervised Learning
- A Probabilistic Point of View
- Parametric Conditional Density Modeling
- Non Parametric Conditional Density Modeling
- Generative Modeling
- Model Selection
- Penalization

3 ML Methods: Optimization Point of View

- Supervised Learning
- Optimization Point of View
- SVM
- Penalization
- Cross Validation and Weights

4 Optimization: Gradient Descent Algorithms

- Introduction
- Gradient Descent
- Proximal Descent
- Coordinate Descent
- Gradient Descent Acceleration
- Stochastic Gradient Descent

• Gradient Descent Step

- Non-Convex Setting
- References

5 ML Methods: Neural Networks and Deep Learning

- Introduction
- From Logistic Regression to NN
- NN Optimization
- NN Regularization
- Image and CNN
- Text, Recurrent Neural Networks and Transformers
- NN Architecture
- References

6 ML Methods: Trees and Ensemble Methods

- Trees
- Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
- Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
- Ensemble Methods
- References

7 Unsupervised Learning: Dimension Reduction and Clustering

- Unsupervised Learning?
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling

• Dimension Reduction

- Simplification
- Reconstruction Error
- Relationship Preservation
- Comparing Methods?

• Clustering

- Prototype Approaches
- Contiguity Approaches
- Agglomerative Approaches
- Other Approaches
- Scalability

• Generative Modeling

- (Plain) Parametric Density Estimation
- Latent Variables
- Approximate Simulation
- Diffusion Model
- Generative Adversarial Network

• Applications to Text

- References

8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory

- Supervised Learning
- Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization

- References

9 References

1 Statistical Learning: Introduction, Setting and Risk Estimation

• Introduction

- Machine Learning
- Supervised Learning
- Risk Estimation and Model Selection
- Cross Validation and Test
- References

2 ML Methods: Probabilistic Point of View

- Motivation
- Supervised Learning
- A Probabilistic Point of View
- Parametric Conditional Density Modeling
- Non Parametric Conditional Density Modeling
- Generative Modeling
- Model Selection
- Penalization

3 ML Methods: Optimization Point of View

- Supervised Learning
- Optimization Point of View
- SVM
- Penalization
- Cross Validation and Weights

4 Optimization: Gradient Descent Algorithms

- Introduction
- Gradient Descent
- Proximal Descent
- Coordinate Descent
- Gradient Descent Acceleration
- Stochastic Gradient Descent

• Gradient Descent Step

- Non-Convex Setting
- References

5 ML Methods: Neural Networks and Deep Learning

- Introduction
- From Logistic Regression to NN
- NN Optimization
- NN Regularization
- Image and CNN
- Text, Recurrent Neural Networks and Transformers
- NN Architecture
- References

6 ML Methods: Trees and Ensemble Methods

- Trees
- Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
- Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
- Ensemble Methods
- References

7 Unsupervised Learning: Dimension Reduction and Clustering

- Unsupervised Learning?
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling

• Dimension Reduction

- Simplification
- Reconstruction Error
- Relationship Preservation
- Comparing Methods?

• Clustering

- Prototype Approaches
- Contiguity Approaches
- Agglomerative Approaches
- Other Approaches
- Scalability

• Generative Modeling

- (Plain) Parametric Density Estimation
- Latent Variables
- Approximate Simulation
- Diffusion Model
- Generative Adversarial Network

• Applications to Text

- References

8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory

- Supervised Learning
- Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization

- References

9 References

Major Influences

Four major influences act today:

- The formal theories of statistics
- Accelerating developments in computers and display devices
- The challenge, in many fields, of more and ever larger bodies of data
- The emphasis on quantification in an ever wider variety of disciplines

Major Influences - Tukey (1962)

Four major influences act today:

- The formal theories of statistics
 - Accelerating developments in computers and display devices
 - The challenge, in many fields, of more and ever larger bodies of data
 - The emphasis on quantification in an ever wider variety of disciplines
-
- He was talking of Data Analysis.
 - Data Mining, Machine Learning, Big Data, Data Science, Artificial Intelligence. . .

Large Scale ML Is (Quite) Easy

```
def run(params: Params) {
  val conf = new SparkConf()
  .setAppName(s"BinaryClassification with $params")
  val sc = new SparkContext(conf)

  Logger.getRootLogger.setLevel(Level.WARN)

  val examples = MLUtils.loadLibSVMFile(sc, params.input).cache()

  val splits = examples.randomSplit(Array(0.8, 0.2))
  val training = splits(0).cache()
  val test = splits(1).cache()
  val numTraining = training.count()
  val numTest = test.count()
  println(s"Trainings: $numTraining, Test: $numTest.")
  examples.unpersist(blocking = false)

  val updater = params.regType match {
    case L1 => new L1Updater()
    case L2 => new SquaredUpdater()
  }

  val algorithm = new LogisticRegressionWithSGD()
  .setNumIterations(params.numIterations)
  .setStepSize(params.stepSize)
  .setUpdater(updater)
  .setRegParam(params.regParam)
  val model = algorithm.run(training).clearThreshold()

  val prediction = model.predict(test.map(_.features))
  val predictionAndLabel = prediction.zip(test.map(_.label))

  val metrics = new BinaryClassificationMetrics(predictionAndLabel)
  val myMetrics = new MyBinaryClassificationMetrics(predictionAndLabel)

  println(s"Empirical CrossEntropy = ${myMetrics.crossEntropy().}")
  println(s"Test areaUnderPR = ${metrics.areaUnderPR().}")
  println(s"Test areaUnderROC = ${metrics.areaUnderROC().}")

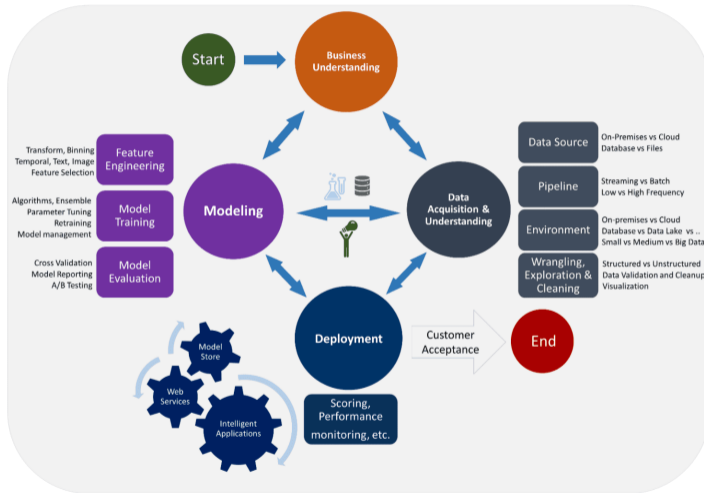
  sc.stop()
}
```



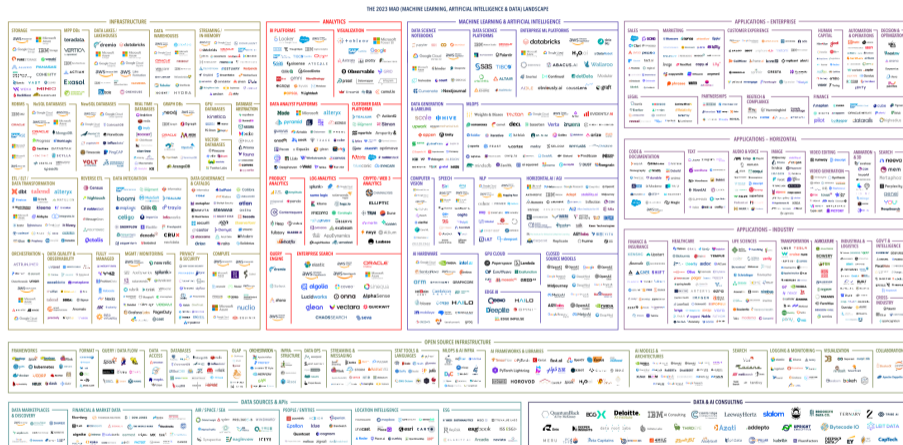
Example of **off the shelves** solution

- Algorithm implementation + copy/paste + cloud computing.
- Machine learning on an arbitrary large dataset!

Data Science Is (Quite) Complex!



Data Ecosystem Is (Quite) Complex!



Version 1.0 - Feb 2023

© Matt Turck (@matturck), Kevin Zhang (@kylevzhang) & FirstMark (@firstmarkcap)

Blog post: matturck.com/MAD2023

Interactive version: MAD.firstmarkcap.com

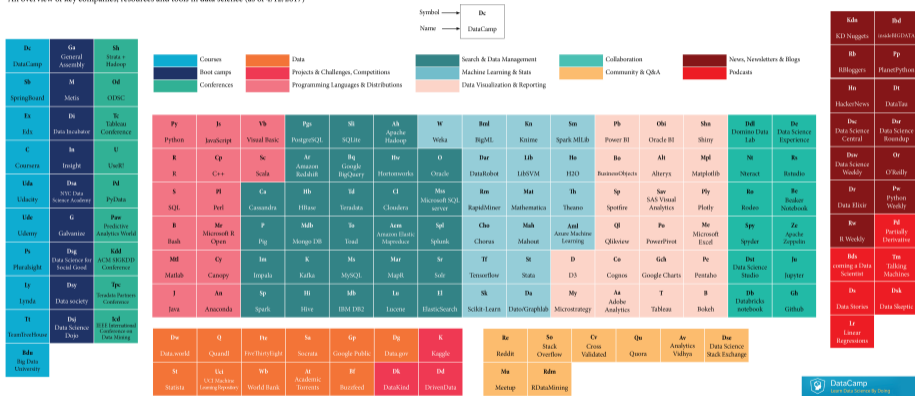
Comments? Email MAD2023@firstmarkcap.com



Data Ecosystem Is (Quite) Complex!

The Periodic Table of Data Science

An overview of key companies, resources and tools in data science (as of 4/12/2017)



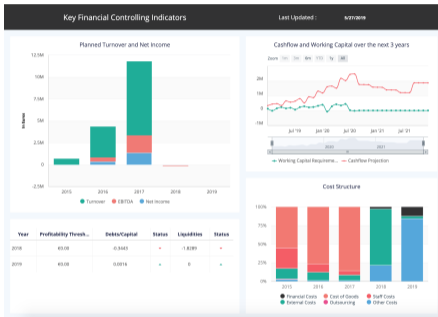
Dc	Ga	Sh
DataCamp	General Assembly	Strata + Hadoop
Sb	M	Od
Springboard	Metis	ODSC
Es	Ds	Tc
Eds	Data Incubate	Tableau Conference
C	In	UwR
Coursera	Insight	UwR
Lda	Du	Pl
Udacity	NYC Data Science Academy	PyData
Lda	G	Paw
Udemy	Gabmatix	Predictive Analytics World
Ps	Dug	KdM
Pharablight	Data Science for Social Good	ACM SIGKDD Conference
Is	Duy	Tpc
Lytnda	Data society	Twinkl Parisian Conference
Ti	Dsj	Kd
TeamFoodHouse	Data Science Days	IEEE International Conference on Data Mining
Edx	Big Data University	

Py	Js	Vb	Pgs	Sl	Ah	W	Brd	Kn	Sm	Pb	Obi	Sha	Ddb	Ds
Python	JavaScript	Visual Basic	PostgreSQL	SQLite	Apache Hadoop	Weka	BigML	Kntime	Spark MLlib	Power BI	Oracle BI	Shiny	Dominio Data Lab	Data Science Experience
R	Cp	Sc	Ar	Bq	Hr	O	Dar	Lib	Hs	Bo	Alt	Mpl	Ni	Ko
R	C++	Scala	Amazon Redshift	Google BigQuery	Hortonworks	Oracle	DataRobot	LibSVM	H2O	BusinessObjects	Alteryx	Maplotlib	Stratoc	Refradio
S	Pl	Ca	Hb	td	Cl	Mis	Rm	Mat	Th	Sp	Sar	Ply	Ro	Re
SQL	Perl	Cassandra	HBase	YaraData	Cloudera	Microsoft SQL server	RapidMiner	Mathematica	Theano	Spotify	SAS Visual Analytics	Plotly	Rodao	ReeBooker Notebook
B	Mr	P	Mlb	To	Arm	Spl	Cho	Mah	AmI	Ql	Po	Me	Spr	Zz
Bash	Microsoft R Open	Pig	Mongo DB	Toad	Amazon EMR/Hadoop/Azure	Spark	Chorus	Mahout	Azure Machine Learning	Qlikview	PowerPivot	Microsoft Excel	Spider	Agave Zippin
Md	Cy	In	K	Ms	Mar	Sr	Tf	Se	D	Co	Gch	Pe	Dst	Jo
Matlab	Canopy	Impala	Kafka	MySQL	Magi	So	TensorFlow	Stata	D3	Cognos	Google Charts	Pentaho	Data Science Studio	Jupyter
J	An	Sp	H	Mb	Lo	El	Sk	Du	My	Aa	T	B	Hb	Gh
Java	Anaconda	Spark	Hive	IBM DB2	Lucene	ElasticSearch	Scikit-Learn	DansGraphLab	Microstrategy	Adobe Analytics	Tableau	Bokeh	Databricks notebook	GitHub
Dw	Q	Hr	Sa	Gp	Dg	K	Rz	So	Cv	Qu	Ar	Du	Lr	
Data.world	Quandl	FiveThirtyEight	Secranta	Google Public	Data.gov	Kaggle	Reddit	Stack Overflow	Cross Validated	Quora	Ar Analytics Vidhya	Data Science Stack Exchange	Linear Regressors	
St	Uct	Wb	M	Bf	Bk	Ed	Ma	Rdm						
Statista	U.S. Machine Learning Database	World Bank	Academic Torrents	Buzzfeed	DataKind	DrivenData	Mextup	RDataMining						

Kd	Bd
KD Nuggets	100k-100DATA
Bb	Pp
BBloggers	PlanetPython
Hn	Dt
HackerNews	DataTau
Dsc	Dur
Data Science Central	Data Science Roundup
Dow	Or
Data Science Weekly	O'Reilly
Dr	Pw
Data Elixir	Python Weekly
Bw	Pd
B Weekly	Partially Derivative
Bd	Tp
Building a Data Scientist	Talking Machines
Ds	Dsk
Data Stories	Data Sleights
Lr	
Linear Regressors	



Monthly KPI Dashboard

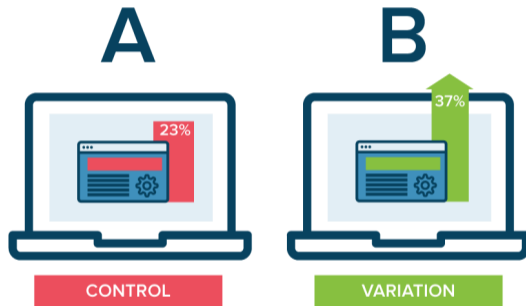


Realtime Log Dashboard

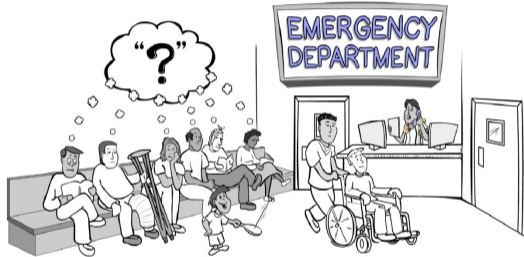


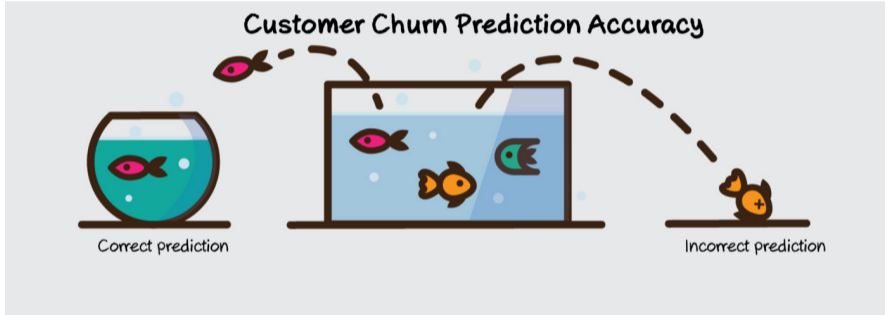
On-demand Legal Document Generation



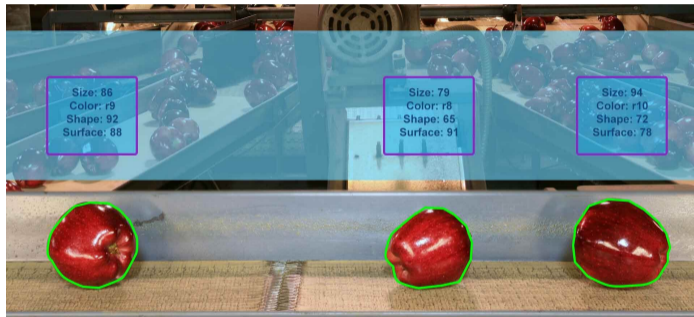


ER Waiting Time Prediction

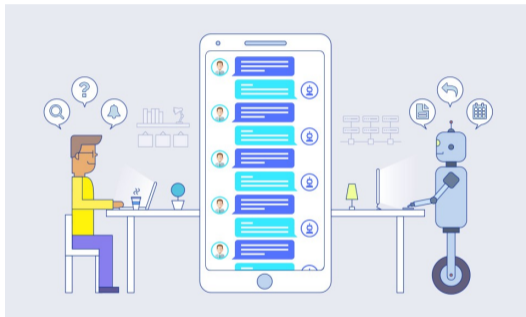


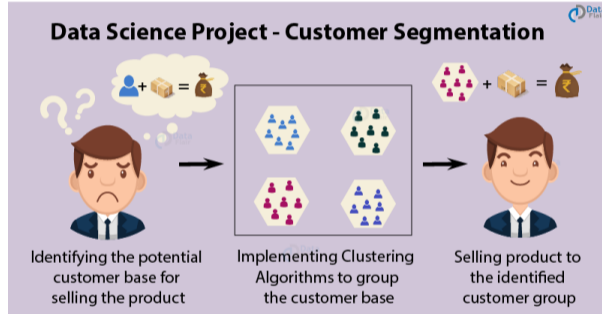


Realtime Automatic Fruit Sorting

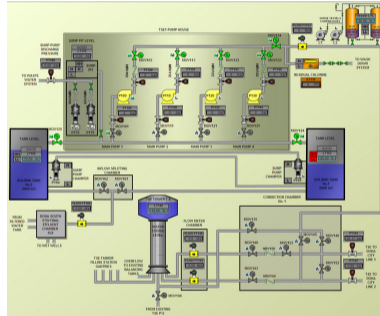


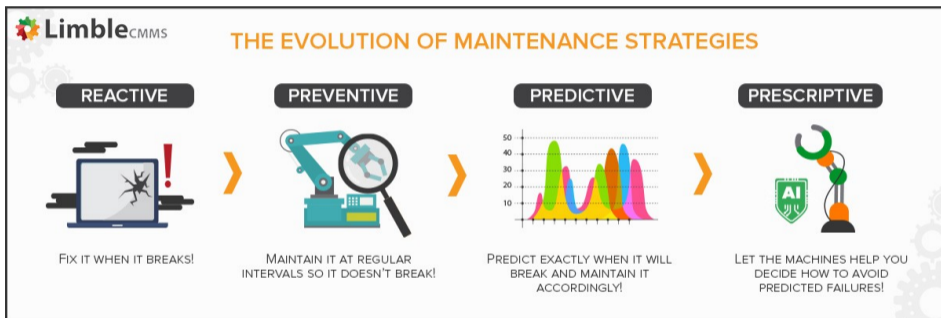
Realtime Chatbot





Realtime Anomaly Detection





1 Statistical Learning: Introduction, Setting and Risk Estimation

- Introduction
- **Machine Learning**
- Supervised Learning
- Risk Estimation and Model Selection
- Cross Validation and Test
- References

2 ML Methods: Probabilistic Point of View

- Motivation
- Supervised Learning
- A Probabilistic Point of View
- Parametric Conditional Density Modeling
- Non Parametric Conditional Density Modeling
- Generative Modeling
- Model Selection
- Penalization

3 ML Methods: Optimization Point of View

- Supervised Learning
- Optimization Point of View
- SVM
- Penalization
- Cross Validation and Weights

4 Optimization: Gradient Descent Algorithms

- Introduction
- Gradient Descent
- Proximal Descent
- Coordinate Descent
- Gradient Descent Acceleration
- Stochastic Gradient Descent

- Gradient Descent Step
- Non-Convex Setting
- References

5 ML Methods: Neural Networks and Deep Learning

- Introduction
- From Logistic Regression to NN
- NN Optimization
- NN Regularization
- Image and CNN
- Text, Recurrent Neural Networks and Transformers
- NN Architecture
- References

6 ML Methods: Trees and Ensemble Methods

- Trees
- Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
- Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
- Ensemble Methods
- References

7 Unsupervised Learning: Dimension Reduction and Clustering

- Unsupervised Learning?
- A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling

● Dimension Reduction

- Simplification
- Reconstruction Error
- Relationship Preservation
- Comparing Methods?

● Clustering

- Prototype Approaches
- Contiguity Approaches
- Agglomerative Approaches
- Other Approaches
- Scalability

● Generative Modeling

- (Plain) Parametric Density Estimation
- Latent Variables
- Approximate Simulation
- Diffusion Model
- Generative Adversarial Network

● Applications to Text

- References

8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory

- Supervised Learning
- Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization

- References

9 References



Google News

Search

Headlines Local For You U.S.

SETTINGS

- Top Stories
- World
- U.S.
- Business
- Technology
- Entertainment
- Sports
- Science
- Health
- Manage sections

Top Stories

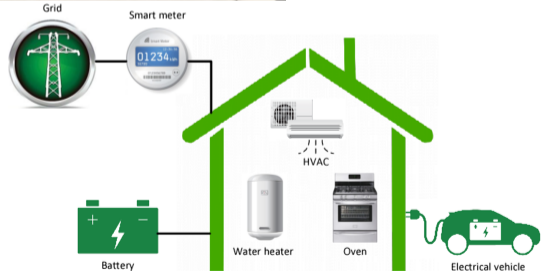
Sarah Huckabee Sanders rips CNN, media at heated briefing
Fox News - 10 ago
RELATED COVERAGE
These journalists leaving CNN after retracted article
Highly rated - CNNMoney - Jun 26, 2017
Reporter accuses White House of 'inflaming' media tensions in heated exchange
The Hill (blog) - 10 ago
Opinion: CNN journalists screwed up, then quit — should that be the standard at White House, on Wall Street and in ...
Opinion - MarketWatch - 10 ago
Reporter interjects as Sanders denounces media
Washington Post

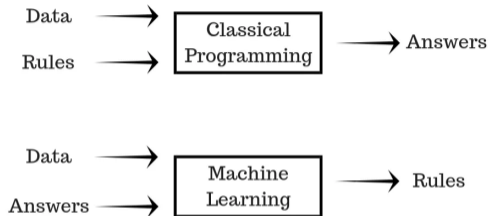
A Time Magazine with Trump on the cover hangs in his golf clubs. It's fake.
Washington Post - 40 ago
RELATED COVERAGE
Donald Trump: Time Magazine Covers: See Them All! Time.com
Mail Reference - Time Magazine - 10 ago

Google N
We've made with more o more covers
Read our list

In the Ne

- Executive1
- Sereno WI
- Dyia
- Google
- John McE
- Barbar al





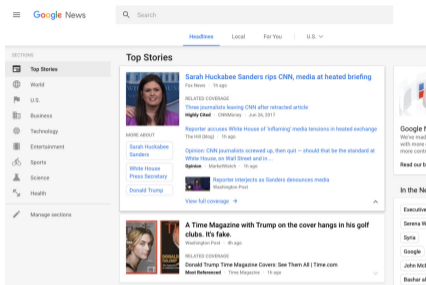
A definition by Tom Mitchell

A computer program is said to learn from **experience E** with respect to some **class of tasks T** and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E.



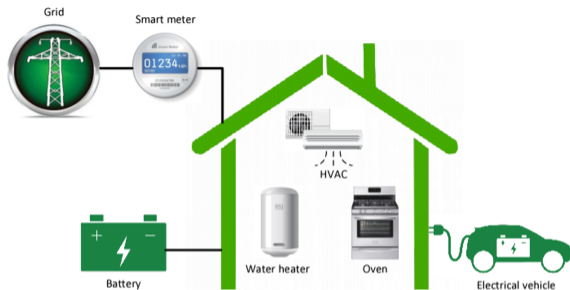
A detection algorithm:

- **Task:** say if an object is present or not in the image
- **Performance:** number of errors
- **Experience:** set of previously seen labeled images



An article clustering algorithm:

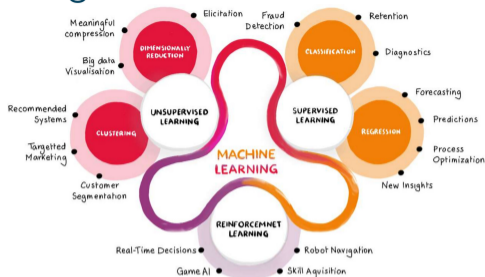
- **Task:** group articles corresponding to the same news
- **Performance:** quality of the clusters
- **Experience:** set of articles



A controller in its sensors in a home smart grid:

- **Task:** control the devices
- **Performance:** energy costs
- **Experience:**
 - previous days
 - current environment and performed actions

Three Kinds of Learning



Unsupervised Learning

- **Task:** Clustering/DR/Generative
- **Performance:** Quality
- **Experience:** Raw dataset (No Ground Truth)

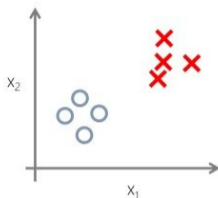
Supervised Learning

- **Task:** Prediction/Classification
- **Performance:** Average error
- **Experience:** Good Predictions (Ground Truth)

Reinforcement Learning

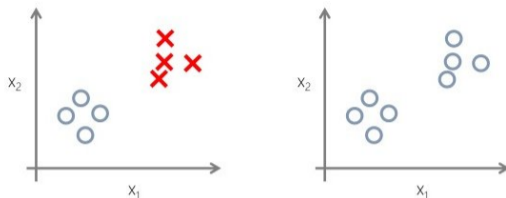
- **Task:** Actions
- **Performance:** Total reward
- **Experience:** Reward from env. (Interact. with env.)

- **Timing:** Offline/Batch (learning from past data) vs Online (continuous learning)



Supervised Learning (Imitation)

- **Goal:** Learn a function f predicting a variable Y from an individual \underline{X} .
- **Data:** Learning set with labeled examples (\underline{X}_i, Y_i)
- **Assumption:** Future data behaves as past data!
- **Predicting is not explaining!**



Supervised Learning (Imitation)

- **Goal:** Learn a function f predicting a variable Y from an individual \underline{X} .
- **Data:** Learning set with labeled examples (\underline{X}_i, Y_i)
- **Assumption:** Future data behaves as past data!
- **Predicting is not explaining!**

Unsupervised Learning (Structure Discovery)

- **Goal:** Discover a structure within a set of individuals (\underline{X}_i) .
- **Data:** Learning set with unlabeled examples (\underline{X}_i)
- Unsupervised learning is not a well-posed setting...



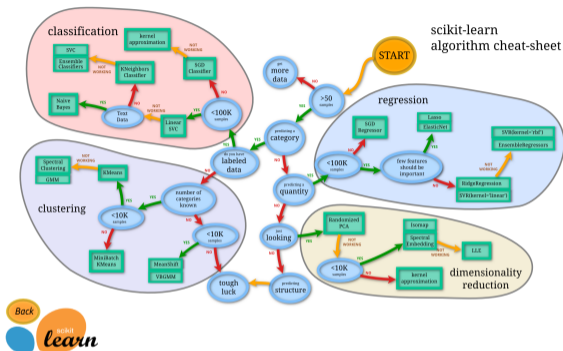
Machine Can

- Forecast (Prediction using the past)
- Detect expected changes
- Memorize/Reproduce
- Take a decision very quickly
- Learn from huge dataset
- Optimize a single task
- Replace/Help some humans

Machine Cannot

- Predict something never seen before
- Detect any new behaviour
- Create something brand new
- Understand the world
- Get smart really fast
- Go beyond their task
- Kill all humans

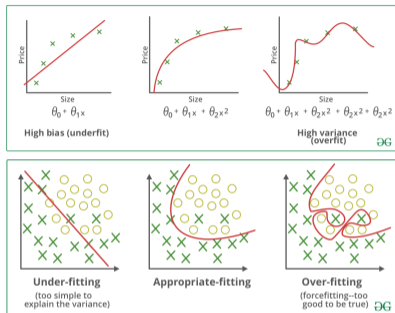
- Some progresses but still very far from the *singularity*...



Machine Learning Methods

- Huge catalog of methods,
- Need to define the performance,
- Numerous tricks: feature design, hyperparameter selection. . .

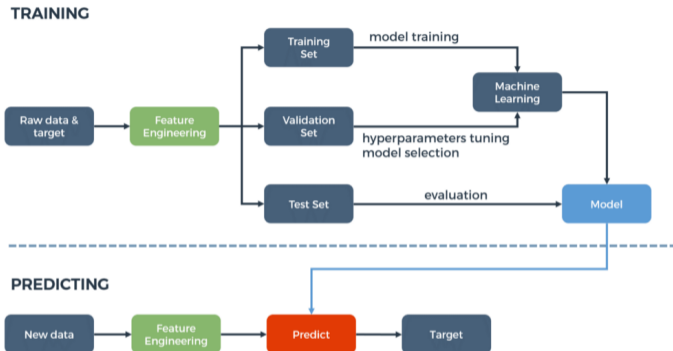
Under and Over Fitting



Finding the Right Complexity

- What is best?
 - A simple model that is stable but false? (*oversimplification*)
 - A very complex model that could be correct but is unstable? (*conspiracy theory*)
- Neither of them: tradeoff that depends on the dataset.

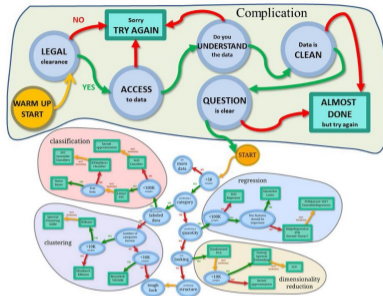
Machine Learning Pipeline



Learning pipeline

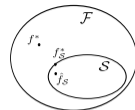
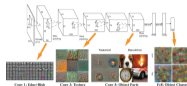
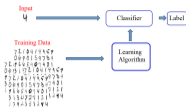
- Test and compare models.
- Deployment pipeline is different!

Data Science \neq Machine Learning



Main DS difficulties

- Figuring out the problem,
- Formalizing it,
- Storing and accessing the data,
- Deploying the solution,
- Not (always) the Machine Learning part!



Goal

- Master the **statistical learning** framework and its challenges.
- Know the inner machinery of the most classical supervised and unsupervised **ML methods** in order to understand their strengths, limitations and connections.
- Understand some **optimization** tools used in ML as well as some **theoretical aspects** of ML.
- Not a course on practical tricks to use machine learning in a data product!

Evaluation

- A practical lab (5 pt)
- A final exam or a project (15 pt)

- Erwan Le Pennec



`Erwan.Le-Pennec@polytechnique.edu`

- Randal Douc



`randal.douc@it-sudparis.eu`

- El Mahdi El Mhamdi



`elmahdi@elmhamdi.com`

- Thierry Klein



`thierry.klein@math.univ-toulouse.fr`

- Edouard Oyallon



`edouard.oyallon@lip6.fr`

9 Lectures (8h30-10h30)

- Mon. 19/09: Statistical Learning: Introduction, Setting and Risk Estimation
- Mon. 26/09: ML Methods: Probabilistic Point of View
- Mon. 03/10: ML Methods: Optimization Point of View
- Mon. 10/10: Optimization: Gradient Descent Algorithms
- Mon. 17/10: ML Methods: Neural Networks and Deep Learning
- Mon. 24/10: ML Methods: Trees and Ensemble Methods
- Mon. 07/11: Unsupervised Learning: Dimension Reduction
- Mon. 14/11: Unsupervised Learning: Clustering
- Mon. 21/11: Statistical Learning: PAC-Bayesian Approach and Complexity Theory

- Mon. ??/12: **Exam** (?)

References



T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning.
Springer Series in Statistics, 2009



M. Mohri, A. Rostamizadeh, and A. Talwalkar.
Foundations of Machine Learning (2nd ed.)
MIT Press, 2018



A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



Ch. Giraud.
Introduction to High-Dimensional Statistics.
CRC Press, 2014



S. Bubeck.
Convex Optimization: Algorithms and Complexity.
Now Publisher, 2015

1 Statistical Learning: Introduction, Setting and Risk Estimation

- Introduction
- Machine Learning
- **Supervised Learning**
- Risk Estimation and Model Selection
- Cross Validation and Test
- References

2 ML Methods: Probabilistic Point of View

- Motivation
- Supervised Learning
- A Probabilistic Point of View
- Parametric Conditional Density Modeling
- Non Parametric Conditional Density Modeling
- Generative Modeling
- Model Selection
- Penalization

3 ML Methods: Optimization Point of View

- Supervised Learning
- Optimization Point of View
- SVM
- Penalization
- Cross Validation and Weights

4 Optimization: Gradient Descent Algorithms

- Introduction
- Gradient Descent
- Proximal Descent
- Coordinate Descent
- Gradient Descent Acceleration
- Stochastic Gradient Descent

● Gradient Descent Step

- Non-Convex Setting
- References

5 ML Methods: Neural Networks and Deep Learning

- Introduction
- From Logistic Regression to NN
- NN Optimization
- NN Regularization
- Image and CNN
- Text, Recurrent Neural Networks and Transformers
- NN Architecture
- References

6 ML Methods: Trees and Ensemble Methods

- Trees
- Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
- Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
- Ensemble Methods
- References

7 Unsupervised Learning: Dimension Reduction and Clustering

- Unsupervised Learning?
- A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling

● Dimension Reduction

- Simplification
- Reconstruction Error
- Relationship Preservation
- Comparing Methods?

● Clustering

- Prototype Approaches
- Contiguity Approaches
- Agglomerative Approaches
- Other Approaches
- Scalability

● Generative Modeling

- (Plain) Parametric Density Estimation
- Latent Variables
- Approximate Simulation
- Diffusion Model
- Generative Adversarial Network

● Applications to Text

- References

8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory

- Supervised Learning
- Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization

- References

9 References

Supervised Learning Framework

- Input measurement $\underline{X} \in \mathcal{X}$
- Output measurement $Y \in \mathcal{Y}$.
- $(\underline{X}, Y) \sim \mathbb{P}$ with \mathbb{P} unknown.
- **Training data** : $\mathcal{D}_n = \{(\underline{X}_1, Y_1), \dots, (\underline{X}_n, Y_n)\}$ (i.i.d. $\sim \mathbb{P}$)
- Often
 - $\underline{X} \in \mathbb{R}^d$ and $Y \in \{-1, 1\}$ (classification)
 - or $\underline{X} \in \mathbb{R}^d$ and $Y \in \mathbb{R}$ (regression).
- A **predictor** is a function in $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \text{ meas.}\}$

Goal

- Construct a **good** predictor \hat{f} from the training data.
- Need to specify the meaning of good.
- Classification and regression are almost the **same** problem!

Loss function for a generic predictor

- **Loss function:** $\ell(Y, f(\underline{X}))$ measures the goodness of the prediction of Y by $f(\underline{X})$
- Examples:
 - 0/1 loss: $\ell(Y, f(\underline{X})) = \mathbf{1}_{Y \neq f(\underline{X})}$
 - Quadratic loss: $\ell(Y, f(\underline{X})) = |Y - f(\underline{X})|^2$

Risk function

- Risk measured as the average loss for a new couple:

$$\mathcal{R}(f) = \mathbb{E}_{(X, Y) \sim \mathbb{P}}[\ell(Y, f(\underline{X}))]$$

- Examples:
 - 0/1 loss: $\mathbb{E}[\ell(Y, f(\underline{X}))] = \mathbb{P}(Y \neq f(\underline{X}))$
 - Quadratic loss: $\mathbb{E}[\ell(Y, f(\underline{X}))] = \mathbb{E}[|Y - f(\underline{X})|^2]$

- **Beware:** As \hat{f} depends on \mathcal{D}_n , $\mathcal{R}(\hat{f})$ is a random variable!

- The best solution f^* (which is independent of \mathcal{D}_n) is

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{R}(f) = \arg \min_{f \in \mathcal{F}} \mathbb{E}[\ell(Y, f(\underline{X}))] = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\underline{X}} \left[\mathbb{E}_{Y|\underline{X}}[\ell(Y, f(\underline{X}))] \right]$$

Bayes Predictor (explicit solution)

- In binary classification with 0 – 1 loss:

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = +1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ & \Leftrightarrow \mathbb{P}(Y = +1|\underline{X}) \geq 1/2 \\ -1 & \text{otherwise} \end{cases}$$

- In regression with the quadratic loss

$$f^*(\underline{X}) = \mathbb{E}[Y|\underline{X}]$$

Issue: Solution requires to **know** $\mathbb{E}[Y|\underline{X}]$ for all values of \underline{X} !

Machine Learning

- Learn a rule to construct a **predictor** $\hat{f} \in \mathcal{F}$ from the training data \mathcal{D}_n s.t. **the risk** $\mathcal{R}(\hat{f})$ is **small on average** or with high probability with respect to \mathcal{D}_n .
- In practice, the rule should be an algorithm!

Canonical example: Empirical Risk Minimizer

- One restricts f to a subset of functions $\mathcal{S} = \{f_\theta, \theta \in \Theta\}$
- One replaces the minimization of the average loss by the minimization of the empirical loss

$$\hat{f} = f_{\hat{\theta}} = \underset{f_\theta, \theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_\theta(\underline{X}_i))$$

- Examples:
 - Linear regression
 - Linear classification with

$$\mathcal{S} = \{\underline{x} \mapsto \operatorname{sign}\{\underline{x}^\top \beta + \beta^{(0)}\} / \beta \in \mathbb{R}^d, \beta^{(0)} \in \mathbb{R}\}$$

1 Statistical Learning: Introduction, Setting and Risk Estimation

- Introduction
- Machine Learning
- Supervised Learning
- Risk Estimation and Model Selection
- Cross Validation and Test
- References

2 ML Methods: Probabilistic Point of View

- Motivation
- Supervised Learning
- A Probabilistic Point of View
- Parametric Conditional Density Modeling
- Non Parametric Conditional Density Modeling
- Generative Modeling
- Model Selection
- Penalization

3 ML Methods: Optimization Point of View

- Supervised Learning
- Optimization Point of View
- SVM
- Penalization
- Cross Validation and Weights

4 Optimization: Gradient Descent Algorithms

- Introduction
- Gradient Descent
- Proximal Descent
- Coordinate Descent
- Gradient Descent Acceleration
- Stochastic Gradient Descent

● Gradient Descent Step

- Non-Convex Setting
- References

5 ML Methods: Neural Networks and Deep Learning

- Introduction
- From Logistic Regression to NN
- NN Optimization
- NN Regularization
- Image and CNN
- Text, Recurrent Neural Networks and Transformers
- NN Architecture
- References

6 ML Methods: Trees and Ensemble Methods

- Trees
- Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
- Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
- Ensemble Methods
- References

7 Unsupervised Learning: Dimension Reduction and Clustering

- Unsupervised Learning?
- A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling

● Dimension Reduction

- Simplification
- Reconstruction Error
- Relationship Preservation
- Comparing Methods?

● Clustering

- Prototype Approaches
- Contiguity Approaches
- Agglomerative Approaches
- Other Approaches
- Scalability

● Generative Modeling

- (Plain) Parametric Density Estimation
- Latent Variables
- Approximate Simulation
- Diffusion Model
- Generative Adversarial Network

● Applications to Text

- References

8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory

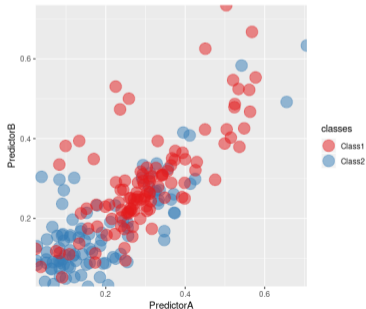
- Supervised Learning
- Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization

- References

9 References

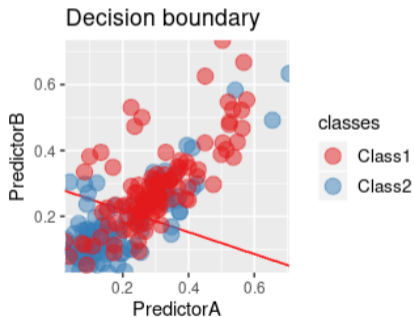
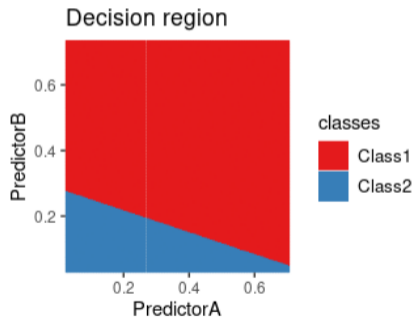
Synthetic Dataset

- Two features/covariates.
- Two classes.
- Dataset from *Applied Predictive Modeling*, M. Kuhn and K. Johnson, Springer
- Numerical experiments with R and the `{caret}` package.



Example: Linear Discrimination

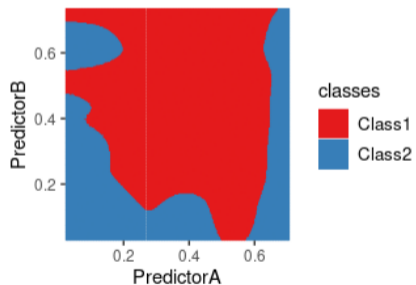
Logistic



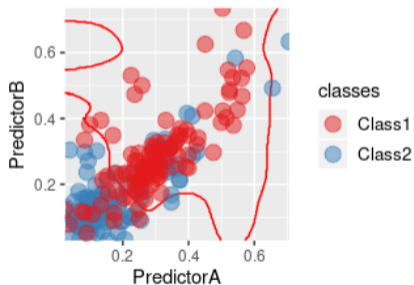
Example: More Complex Model

Naive Bayes with kernel density estimates

Decision region

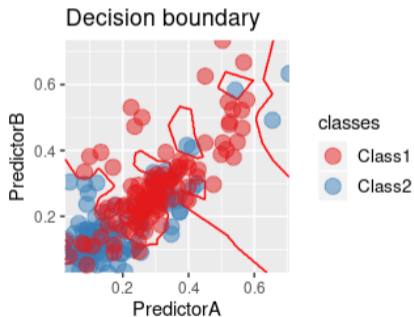
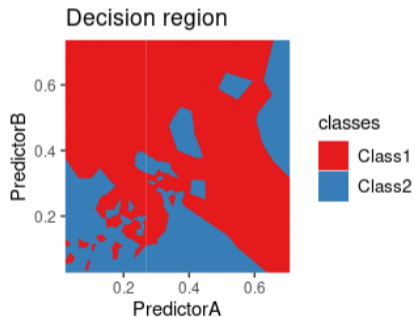


Decision boundary



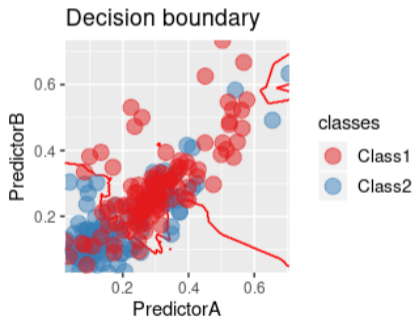
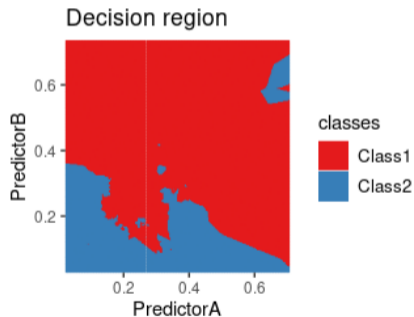
Example: KNN

k-NN with $k=1$



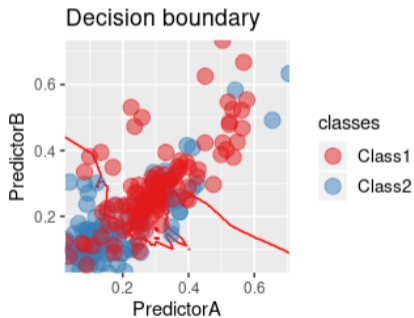
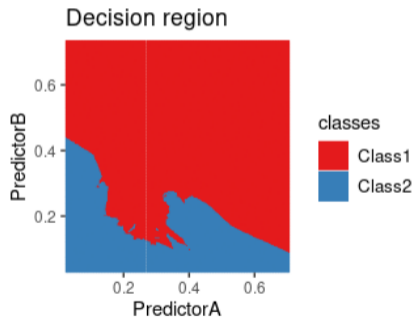
Example: KNN

k-NN with $k=5$



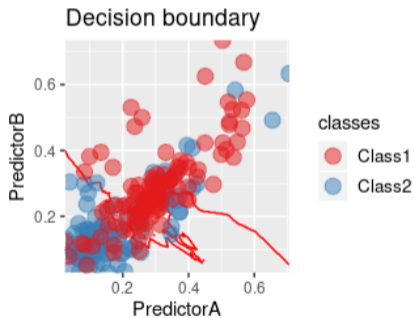
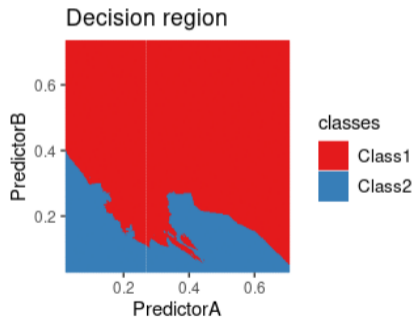
Example: KNN

k-NN with $k=9$



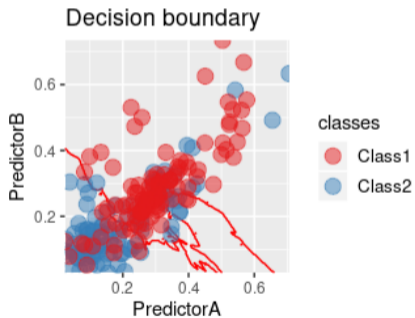
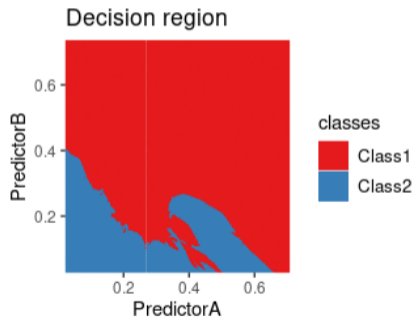
Example: KNN

k-NN with $k=13$



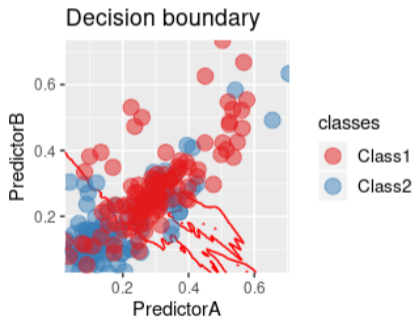
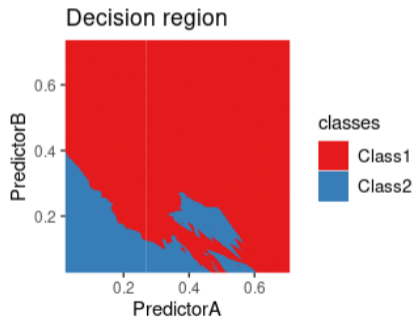
Example: KNN

k-NN with $k=17$



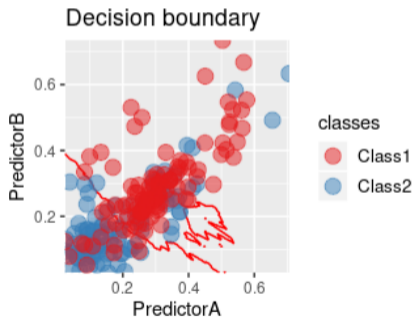
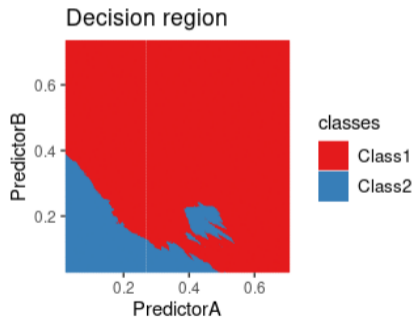
Example: KNN

k-NN with $k=21$



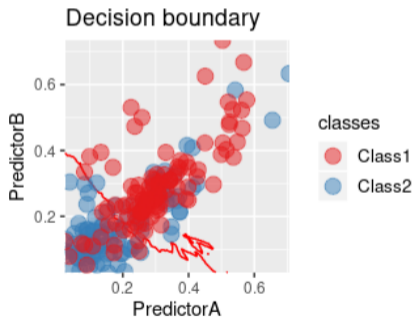
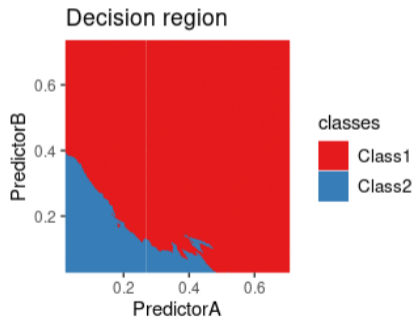
Example: KNN

k-NN with $k=25$



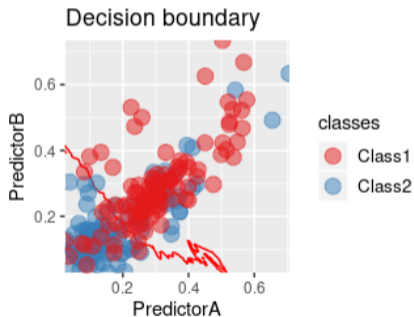
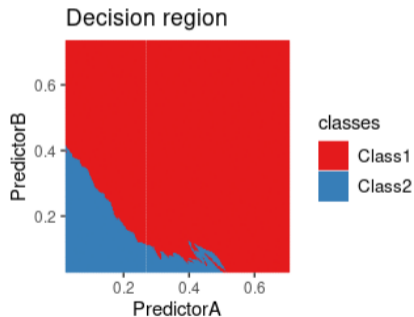
Example: KNN

k-NN with $k=29$



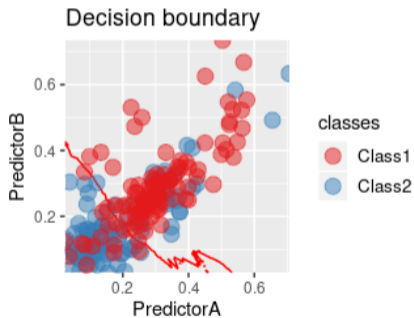
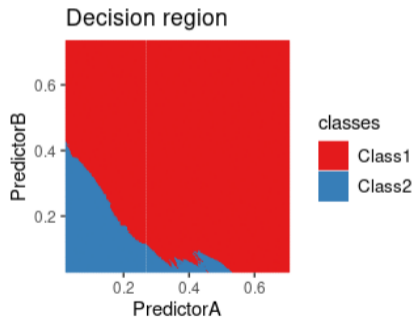
Example: KNN

k-NN with $k=33$



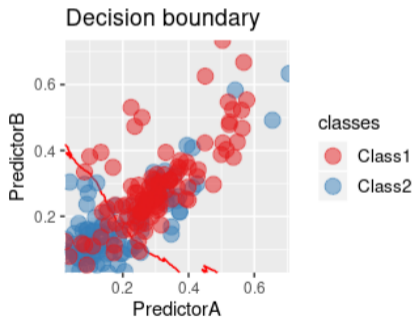
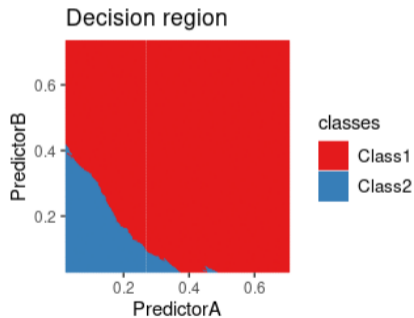
Example: KNN

k-NN with $k=37$



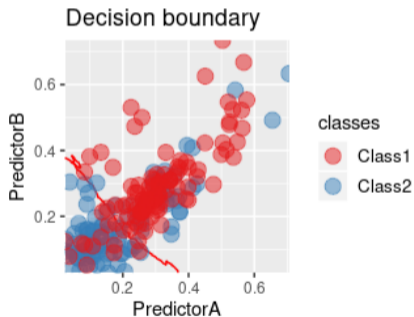
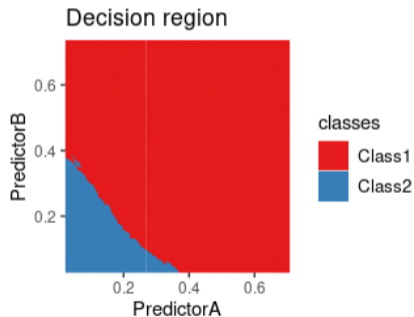
Example: KNN

k-NN with k=45



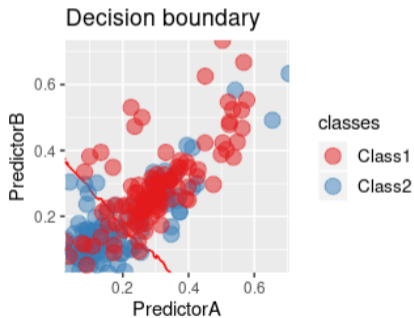
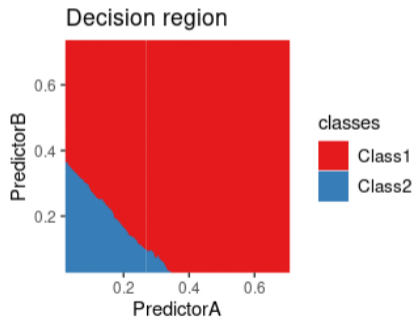
Example: KNN

k-NN with $k=53$



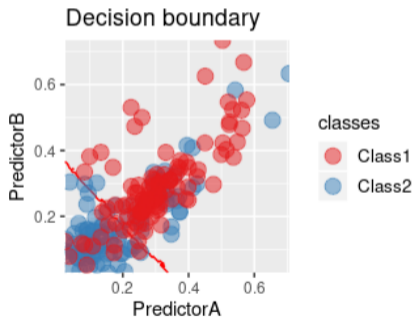
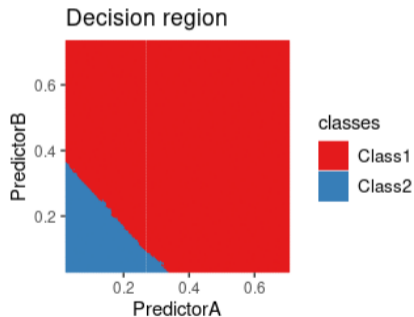
Example: KNN

k-NN with $k=61$



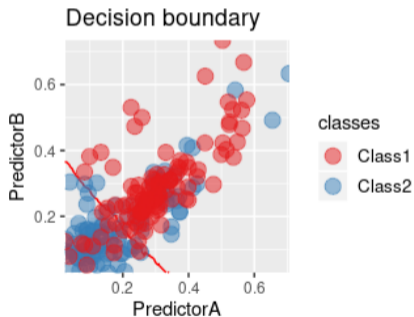
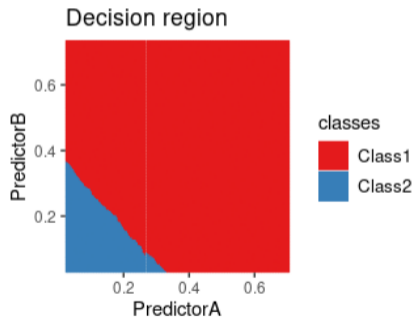
Example: KNN

k-NN with $k=69$



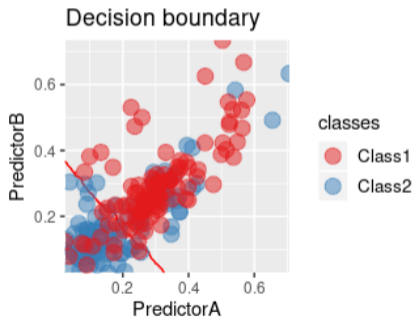
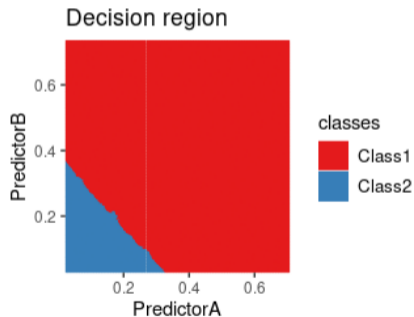
Example: KNN

k-NN with $k=77$



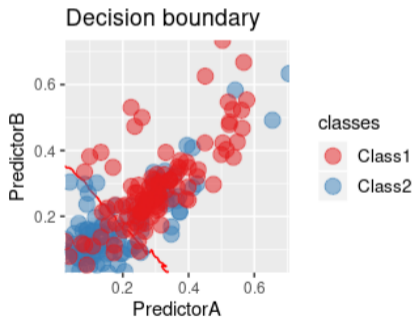
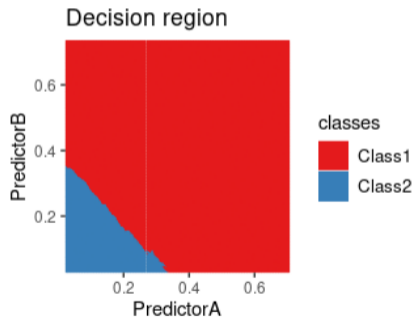
Example: KNN

k-NN with $k=85$



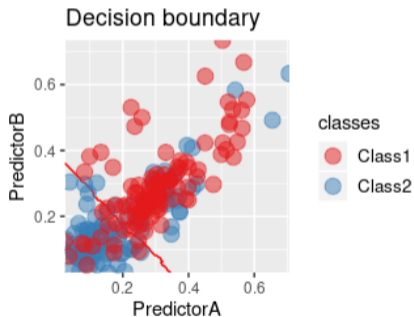
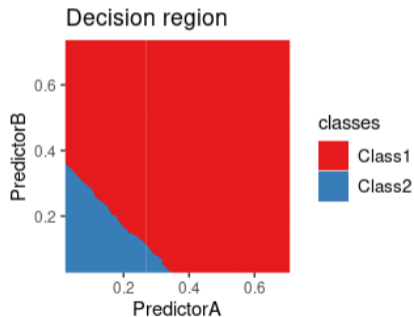
Example: KNN

k-NN with $k=101$



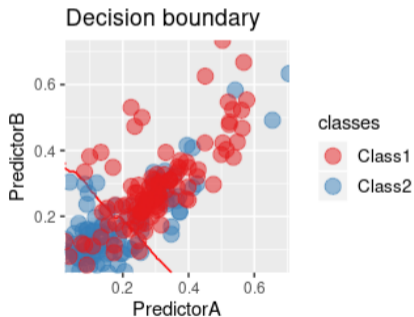
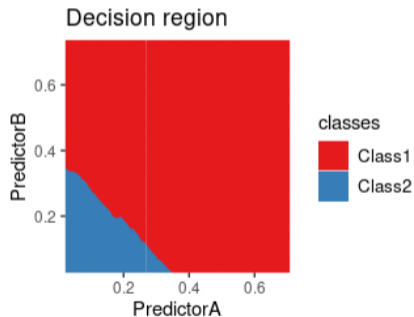
Example: KNN

k-NN with $k=109$



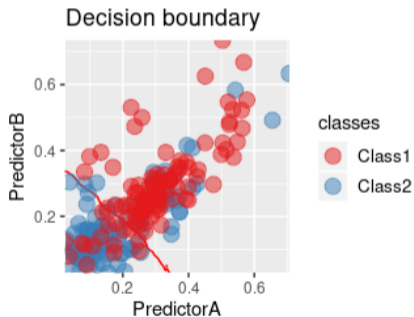
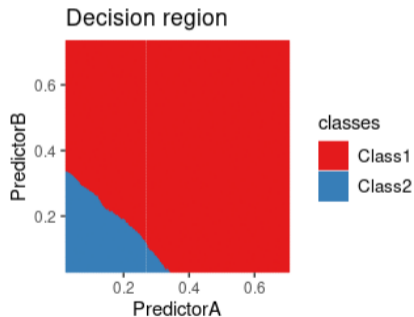
Example: KNN

k-NN with $k=117$



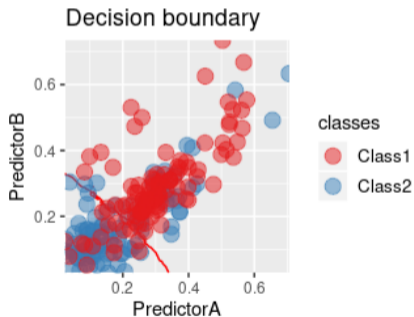
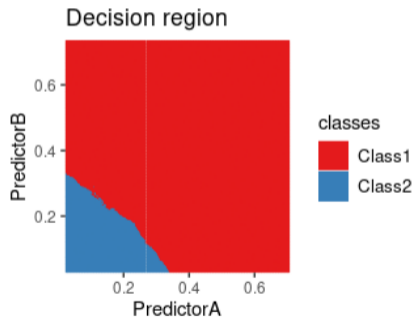
Example: KNN

k-NN with $k=125$



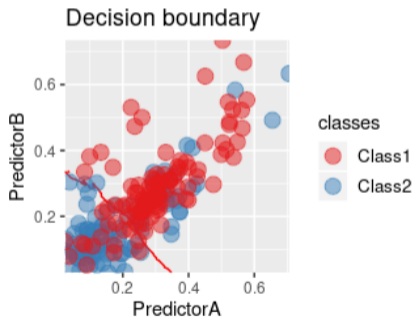
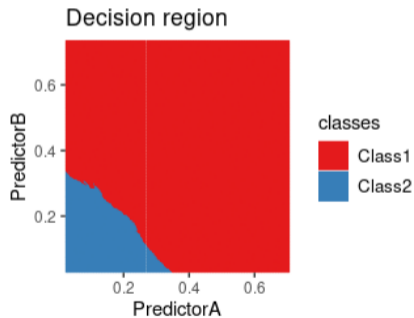
Example: KNN

k-NN with $k=133$



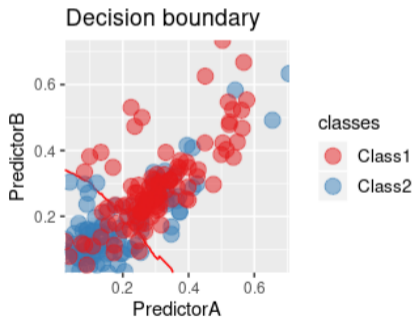
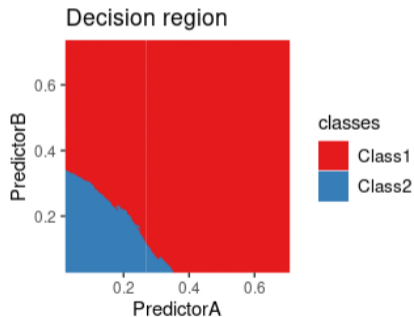
Example: KNN

k-NN with $k=141$



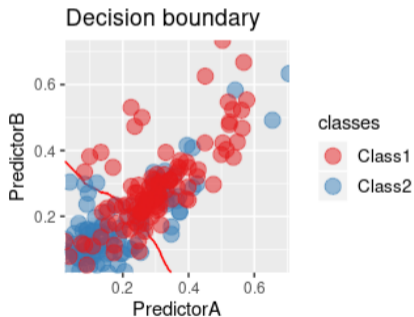
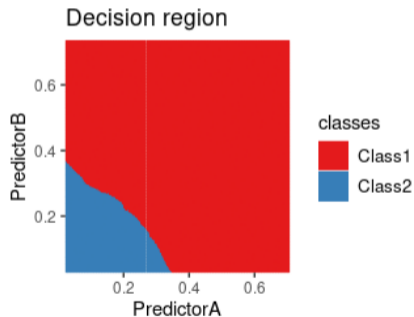
Example: KNN

k-NN with $k=149$



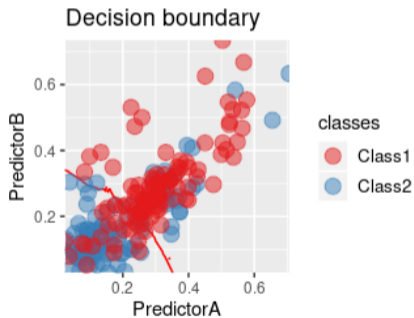
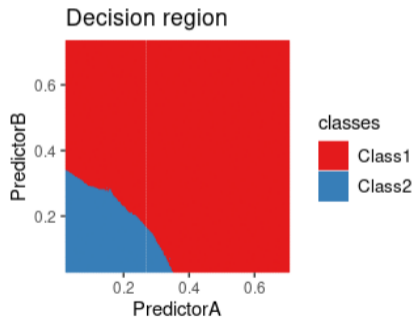
Example: KNN

k-NN with $k=157$



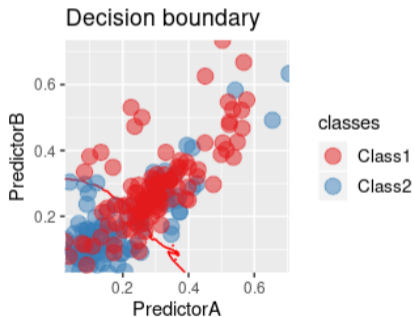
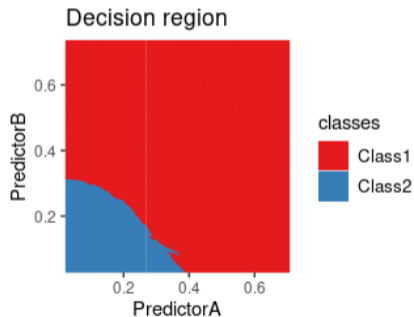
Example: KNN

k-NN with $k=165$



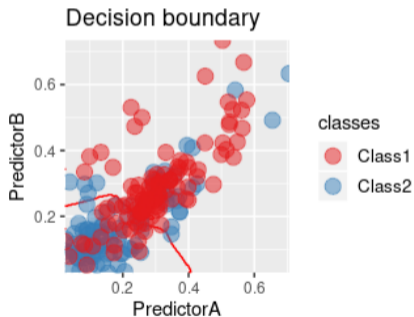
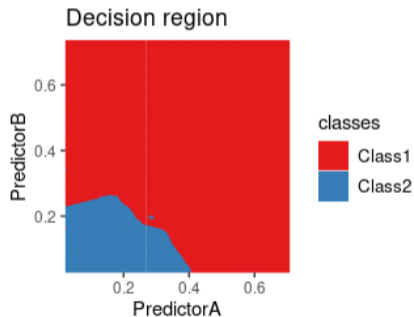
Example: KNN

k-NN with $k=173$



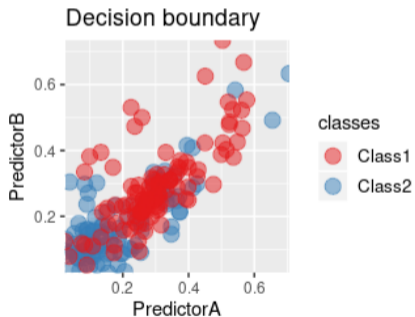
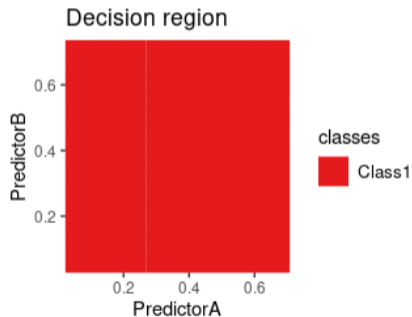
Example: KNN

k-NN with $k=181$



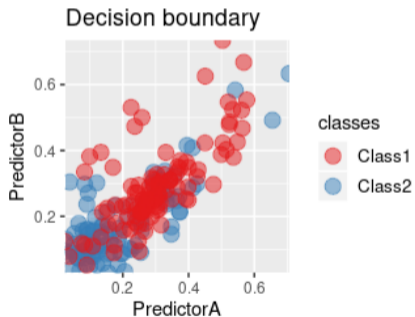
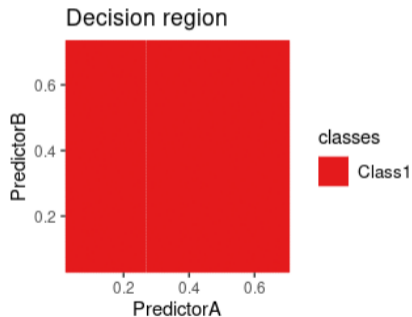
Example: KNN

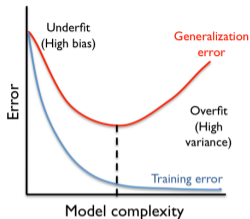
k-NN with $k=189$



Example: KNN

k-NN with $k=197$





Risk behaviour

- Learning/training risk (empirical risk on the learning/training set) decays when the complexity of the **method** increases.
- Quite different behavior when the risk is computed on new observations (generalization risk).
- Overfit for complex methods: parameters learned are too specific to the learning set!
- General situation! (Think of polynomial fit. . .)
- Need to use a different criterion than the training risk!

Predictor Risk Estimation

- **Goal:** Given a predictor f assess its quality.
 - **Method:** Hold-out risk computation (/ Empirical risk correction).
 - **Usage:** Compute an estimate of the risk of a selected f using a **test set** to be used to monitor it in the future.
-
- Basic block very well understood.

Method Selection

- **Goal:** Given a ML method assess its quality.
 - **Method:** Cross Validation (/ Empirical risk correction)
 - **Usage:** Compute risk estimates for several ML methods using **training/validation sets** to choose the most promising one.
-
- Estimates can be pointwise or better intervals.
 - Multiple test issues in method selection.

Two Approaches

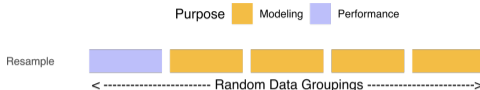
- **Cross validation:** Use empirical risk criterion but on independent data, very efficient (and almost always used in practice!) but slightly biased as its target uses only a fraction of the data.
- **Correction approach:** use empirical risk criterion but *correct* it with a term increasing with the complexity of \mathcal{S}

$$R_n(\hat{f}_S) \rightarrow R_n(\hat{f}_S) + \text{cor}(\mathcal{S})$$

and choose the method with the smallest corrected risk.

Which loss to use?

- The loss used in the risk: most natural!
- The loss used to estimate $\hat{\theta}$: penalized estimation!
- Other performance measure can be used.



- **Very simple idea:** use a second learning/verification set to compute a verification risk.
- Sufficient to remove the dependency issue!
- Implicit random design setting...

Cross Validation

- Use $(1 - \epsilon) \times n$ observations to train and $\epsilon \times n$ to verify!
- Possible issues:
 - Validation for a learning set of size $(1 - \epsilon) \times n$ instead of n ?
 - Unstable risk estimate if ϵn is too small ?
- Most classical variations:
 - Hold Out,
 - Leave One Out,
 - V -fold cross validation.

Principle

- Split the dataset \mathcal{D} in 2 sets $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ of size $n \times (1 - \epsilon)$ and $n \times \epsilon$.
- Learn \hat{f}^{HO} from the subset $\mathcal{D}_{\text{train}}$.
- Compute the empirical risk on the subset $\mathcal{D}_{\text{test}}$:

$$\mathcal{R}_n^{HO}(\hat{f}^{HO}) = \frac{1}{n\epsilon} \sum_{(\underline{X}_i, Y_i) \in \mathcal{D}_{\text{test}}} \ell(Y_i, \hat{f}^{HO}(\underline{X}_i))$$

Predictor Risk Estimation

- Use \hat{f}^{HO} as predictor.
- Use $\mathcal{R}_n^{HO}(\hat{f}^{HO})$ as an estimate of the risk of this estimator.

Method Selection by Cross Validation

- Compute $\mathcal{R}_n^{HO}(\hat{f}_S^{HO})$ for all the considered methods,
- Select the method with the smallest CV risk,
- Reestimate the \hat{f}_S with all the data.

Principle

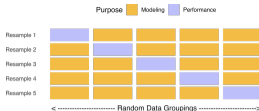
- Split the dataset \mathcal{D} in 2 sets $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ of size $n \times (1 - \epsilon)$ and $n \times \epsilon$.
- Learn \hat{f}^{HO} from the subset $\mathcal{D}_{\text{train}}$.
- Compute the empirical risk on the subset $\mathcal{D}_{\text{test}}$:

$$\mathcal{R}_n^{HO}(\hat{f}^{HO}) = \frac{1}{n\epsilon} \sum_{(\underline{X}_i, Y_i) \in \mathcal{D}_{\text{test}}} \ell(Y_i, \hat{f}^{HO}(\underline{X}_i))$$

- Only possible setting for risk estimation.

Hold Out Limitation for Method Selection

- Biased toward simpler method as the estimation does not use all the data initially.
- Learning variability of $\mathcal{R}_n^{HO}(\hat{f}^{HO})$ not taken into account.



Principle

- Split the dataset \mathcal{D} in V sets \mathcal{D}_v of almost equals size.
- For $v \in \{1, \dots, V\}$:
 - Learn \hat{f}^{-v} from the dataset \mathcal{D} minus the set \mathcal{D}_v .
 - Compute the empirical risk:

$$\mathcal{R}_n^{-v}(\hat{f}^{-v}) = \frac{1}{n_v} \sum_{(\underline{X}_j, Y_j) \in \mathcal{D}_v} \ell(Y_j, \hat{f}^{-v}(\underline{X}_j))$$

- Compute the average empirical risk:

$$\mathcal{R}_n^{CV}(\hat{f}) = \frac{1}{V} \sum_{v=1}^V \mathcal{R}_n^{-v}(\hat{f}^{-v})$$

- Estimation of the quality of a method not of a given predictor.
- Leave One Out : $V = n$.

Analysis (when n is a multiple of V)

- The $\mathcal{R}_n^{-v}(\hat{f}^{-v})$ are identically distributed variable but are not independent!
- Consequence:

$$\begin{aligned}\mathbb{E} \left[\mathcal{R}_n^{CV}(\hat{f}) \right] &= \mathbb{E} \left[\mathcal{R}_n^{-v}(\hat{f}^{-v}) \right] \\ \text{Var} \left[\mathcal{R}_n^{CV}(\hat{f}) \right] &= \frac{1}{V} \text{Var} \left[\mathcal{R}_n^{-v}(\hat{f}^{-v}) \right] \\ &\quad + \left(1 - \frac{1}{V} \right) \text{Cov} \left[\mathcal{R}_n^{-v}(\hat{f}^{-v}), \mathcal{R}_n^{-v'}(\hat{f}^{-v'}) \right]\end{aligned}$$

- Average risk for a sample of size $(1 - \frac{1}{V})n$.
 - Variance term much more complex to analyze!
 - Fine analysis shows that the larger V the better...
-
- Accuracy/Speed tradeoff: $V = 5$ or $V = 10$...

- Leave One Out = V fold for $V = n$: very expensive in general.

A fast LOO formula for the linear regression

- **Prop:** for the least squares linear regression,

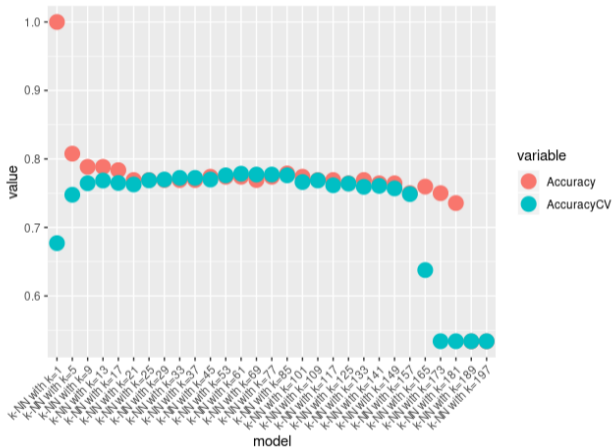
$$\hat{f}^{-i}(\underline{X}_i) = \frac{\hat{f}(\underline{X}_i) - h_{ii} Y_i}{1 - h_{ii}}$$

with h_{ii} the i th diagonal coefficient of the **hat** (projection) matrix.

- Proof based on linear algebra!
- Leads to a fast formula for LOO:

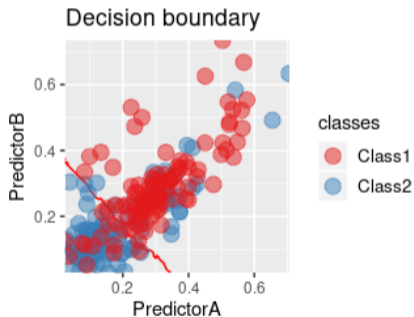
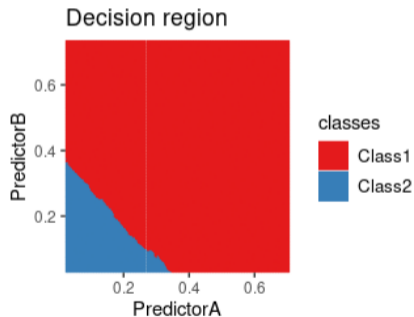
$$\mathcal{R}_n^{LOO}(\hat{f}) = \frac{1}{n} \sum_{i=1}^n \frac{|Y_i - \hat{f}(\underline{X}_i)|^2}{(1 - h_{ii})^2}$$

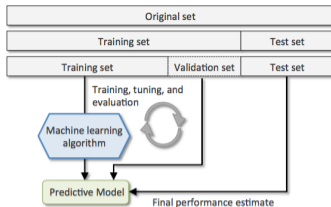
Cross Validation



Example: KNN ($\hat{k} = 61$ using cross-validation)

k-NN with k=61





- **Selection Bias Issue:**
 - After method selection, the cross validation is biased.
 - Furthermore, it qualifies the method and not the final predictor.
- Need to (re)estimate the risk of the final predictor.

(Train/Validation)/Test strategy

- **Split** the dataset in two a (Train/Validation) and Test.
 - Use **CV** with the (Train/Validation) to **select a method**.
 - Train this method on (Train/Validation) to **obtain a single predictor**.
 - Estimate the **performance of this predictor** on Test.
-
- Every choice made from the data is part of the method!

- Empirical loss of an estimator computed on the dataset used to choose it is biased!
- Empirical loss is an optimistic estimate of the true loss.

Risk Correction Heuristic

- Estimate an upper bound of this optimism for a given family.
- Correct the empirical loss by adding this upper bound.
- **Rk:** Finding such an upper bound can be complicated!
- Correction often called a **penalty**.

Penalized Loss

- Minimization of

$$\operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_{\theta}(\underline{X}_i)) + \operatorname{pen}(\theta)$$

where $\operatorname{pen}(\theta)$ is a risk correction (penalty).

Penalties

- Upper bound of the optimism of the empirical loss
- Depends on the loss and the framework!

Instantiation

- Mallows Cp: Least Squares with $\operatorname{pen}(\theta) = 2\frac{d}{n}\sigma^2$.
- AIC Heuristics: Maximum Likelihood with $\operatorname{pen}(\theta) = \frac{d}{n}$.
- BIC Heuristics: Maximum Likelihood with $\operatorname{pen}(\theta) = \log(n)\frac{d}{n}$.
- Structural Risk Minimization: Pred. loss and clever penalty.

1 Statistical Learning: Introduction, Setting and Risk Estimation

- Introduction
- Machine Learning
- Supervised Learning
- Risk Estimation and Model Selection
- **Cross Validation and Test**
- References

2 ML Methods: Probabilistic Point of View

- Motivation
- Supervised Learning
- A Probabilistic Point of View
- Parametric Conditional Density Modeling
- Non Parametric Conditional Density Modeling
- Generative Modeling
- Model Selection
- Penalization

3 ML Methods: Optimization Point of View

- Supervised Learning
- Optimization Point of View
- SVM
- Penalization
- Cross Validation and Weights

4 Optimization: Gradient Descent Algorithms

- Introduction
- Gradient Descent
- Proximal Descent
- Coordinate Descent
- Gradient Descent Acceleration
- Stochastic Gradient Descent

● Gradient Descent Step

- Non-Convex Setting
- References

5 ML Methods: Neural Networks and Deep Learning

- Introduction
- From Logistic Regression to NN
- NN Optimization
- NN Regularization
- Image and CNN
- Text, Recurrent Neural Networks and Transformers
- NN Architecture
- References

6 ML Methods: Trees and Ensemble Methods

- Trees
- Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
- Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
- Ensemble Methods
- References

7 Unsupervised Learning: Dimension Reduction and Clustering

- Unsupervised Learning?
- A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling

● Dimension Reduction

- Simplification
- Reconstruction Error
- Relationship Preservation
- Comparing Methods?

● Clustering

- Prototype Approaches
- Contiguity Approaches
- Agglomerative Approaches
- Other Approaches
- Scalability

● Generative Modeling

- (Plain) Parametric Density Estimation
- Latent Variables
- Approximate Simulation
- Diffusion Model
- Generative Adversarial Network

● Applications to Text

- References

8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory

- Supervised Learning
- Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization

- References

9 References

Means

- **Setting:** r.v. $e_i^{(l)}$ with $1 \leq i \leq n_l$ and $l \in \{1, 2\}$ and their means

$$\overline{e^{(l)}} = \frac{1}{n_l} \sum_{i=1}^{n_l} e_i^{(l)}$$

- **Question:** are the means $\overline{e^{(l)}}$ statistically different?

Classical i.i.d setting

- **Assumption:** $e_i^{(l)}$ are i.i.d. for each l .
- **Test formulation:** Can we reject the null hypothesis that $\mathbb{E}[e^{(1)}] = \mathbb{E}[e^{(2)}]$?
- **Methods:**
 - Gaussian (Student) test using asymptotic normality of a mean.
 - Non-parametric permutation test.

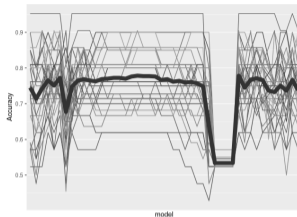
- Gaussian approach is linked to confidence intervals.
- The larger n_l the smaller the confidence intervals.

Non i.i.d. case

- **Assumption:** $e_i^{(l)}$ are i.d. for each l but not necessarily independent.
- **Test formulation:** Can we reject the null hypothesis that $\mathbb{E}[e^{(1)}] = \mathbb{E}[e^{(2)}]$?
- **Methods:**
 - Gaussian (Student) test using asymptotic normality of a mean but variance is hard to estimate.
 - Non-parametric permutation test but no confidence intervals.
- Setting for Cross Validation (other than holdout).
- Much more complicated than the i.i.d. case

Several means

- **Assumption:** $e_i^{(l)}$ are i.i.d. for each l but not necessarily independent.
- **Tests formulation:**
 - Can we reject the null hypothesis that the $\mathbb{E}[e^{(l)}]$ are different?
 - Is the smaller mean statistically smaller than the second one?
- **Methods:**
 - Gaussian (Student) test using asymptotic normality of a mean with multiple tests correction.
 - Non-parametric permutation test but no confidence intervals.
- Setting for Cross Validation (other than holdout).
- The more models one compares:
 - the larger the confidence intervals
 - the most probable the best model is a lucky winner
- Justify the fallback to the simplest model that could be the best one.



CV Risk, Methods and Predictors

- Cross-Validation risk: estimate of the average risk of a ML method.
- No risk bound on the predictor obtained in practice.

Probably-Approximately-Correct (PAC) Approach

- Replace the control on the average risk by a probabilistic bound

$$\mathbb{P}\left(\mathbb{E}\left[\ell(Y, \hat{f}(X))\right] > R\right) \leq \epsilon$$

- Requires estimating quantiles of the risk.

Cross Validation and Confidence Interval

- How to replace pointwise estimation by a confidence interval?
- Can we use the variability of the CV estimates?
- **Negative result:** No unbiased estimate of the variance!

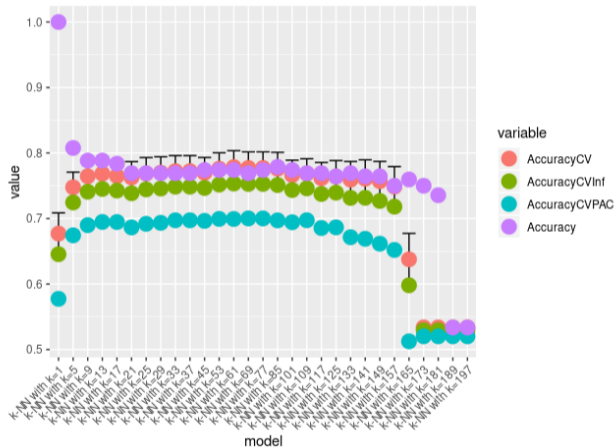
Gaussian Interval (Comparison of the means and \sim indep.)

- Compute the empirical variance and divide it by the number of folds to construct an asymptotic Gaussian confidence interval,
- Select the simplest model whose value falls into the confidence interval of the model having the smallest CV risk.

PAC approach (Quantile, \sim indep. and small risk estim. error)

- Compute the raw medians (or a larger raw quantiles)
- Select the model having the smallest quantiles to ensure a small risk with high probability.
- Always reestimate the chosen model with all the data.
- To obtain an unbiased risk estimate of the final predictor: hold out risk on untouched test data.

Cross Validation



1 Statistical Learning: Introduction, Setting and Risk Estimation

- Introduction
- Machine Learning
- Supervised Learning
- Risk Estimation and Model Selection
- Cross Validation and Test
- References

2 ML Methods: Probabilistic Point of View

- Motivation
- Supervised Learning
- A Probabilistic Point of View
- Parametric Conditional Density Modeling
- Non Parametric Conditional Density Modeling
- Generative Modeling
- Model Selection
- Penalization

3 ML Methods: Optimization Point of View

- Supervised Learning
- Optimization Point of View
- SVM
- Penalization
- Cross Validation and Weights

4 Optimization: Gradient Descent Algorithms

- Introduction
- Gradient Descent
- Proximal Descent
- Coordinate Descent
- Gradient Descent Acceleration
- Stochastic Gradient Descent

● Gradient Descent Step

- Non-Convex Setting
- References

5 ML Methods: Neural Networks and Deep Learning

- Introduction
- From Logistic Regression to NN
- NN Optimization
- NN Regularization
- Image and CNN
- Text, Recurrent Neural Networks and Transformers
- NN Architecture
- References

6 ML Methods: Trees and Ensemble Methods

- Trees
- Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
- Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
- Ensemble Methods
- References

7 Unsupervised Learning: Dimension Reduction and Clustering

- Unsupervised Learning?
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling

● Dimension Reduction

- Simplification
- Reconstruction Error
- Relationship Preservation
- Comparing Methods?

● Clustering

- Prototype Approaches
- Contiguity Approaches
- Agglomerative Approaches
- Other Approaches
- Scalability

● Generative Modeling

- (Plain) Parametric Density Estimation
- Latent Variables
- Approximate Simulation
- Diffusion Model
- Generative Adversarial Network

● Applications to Text

- References

8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory

- Supervised Learning
- Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization

- References

9 References

References



T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning.
Springer Series in Statistics, 2009



M. Mohri, A. Rostamizadeh, and A. Talwalkar.
Foundations of Machine Learning (2nd ed.)
MIT Press, 2018



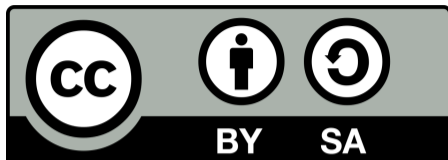
A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



Ch. Giraud.
Introduction to High-Dimensional Statistics.
CRC Press, 2014



S. Bubeck.
Convex Optimization: Algorithms and Complexity.
Now Publisher, 2015



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

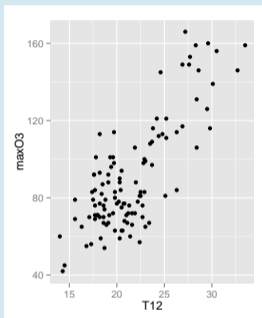
Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 6 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 6 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 6 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 **ML Methods: Probabilistic Point of View**
 - **Motivation**
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- Gradient Descent Step
- Non-Convex Setting
- References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Ozone pollution



- **Data:** Air Breizh, Summer 2001
- **Input:** Temperature at 12h00
- **Output:** max Ozone concentration

Credit Default, Credit Score, Bank Risk, Market Risk Management



- **Data:** Client profile, Client credit history...
- **Input:** Client profile
- **Output:** Credit risk

Marketing: advertisement, recommendation...

More Ideas Based on Your Browsing History

You looked at



[Thriving in the Knowledge Age: New...](#) Paperback by John H. Falk
~~\$29.95~~

You might also consider



[Museum Administration: An Introduction](#) Paperback by Hugh H. Genoways
~~\$34.95~~ **\$28.75**

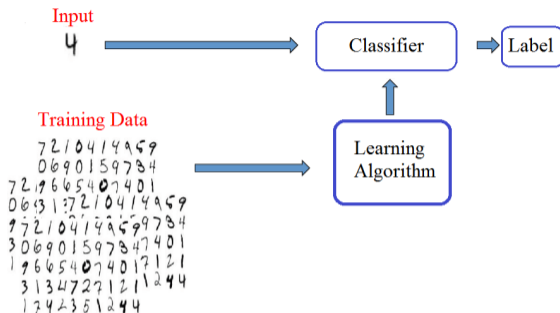


[Exhibit Labels: An Interpretive Approach](#) Paperback by Beverly Serrell
~~\$34.95~~ **\$27.85**

[Find similar items](#)

Recommendations don't have to be about showing you more of the same...

- **Data:** User profile, Web site history...
- **Input:** User profile, Current web page
- **Output:** Advertisement with price, recommendation...



A definition by Tom Mitchell

A computer program is said to learn from **experience E** with respect to some **class of tasks T** and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 **ML Methods: Probabilistic Point of View**
 - Motivation
 - **Supervised Learning**
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- Gradient Descent Step
- Non-Convex Setting
- References
- 5 **ML Methods: Neural Networks and Deep Learning**
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 **ML Methods: Trees and Ensemble Methods**
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 **Unsupervised Learning: Dimension Reduction and Clustering**
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- Applications to Text
- References
- 8 **Statistical Learning: PAC-Bayesian Approach and Complexity Theory**
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 **References**

Supervised Learning Framework

- Input measurement $\underline{X} \in \mathcal{X}$
- Output measurement $Y \in \mathcal{Y}$.
- $(\underline{X}, Y) \sim \mathbb{P}$ with \mathbb{P} unknown.
- **Training data** : $\mathcal{D}_n = \{(\underline{X}_1, Y_1), \dots, (\underline{X}_n, Y_n)\}$ (i.i.d. $\sim \mathbb{P}$)
- Often
 - $\underline{X} \in \mathbb{R}^d$ and $Y \in \{-1, 1\}$ (classification)
 - or $\underline{X} \in \mathbb{R}^d$ and $Y \in \mathbb{R}$ (regression).
- A **predictor** is a function in $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \text{ meas.}\}$

Goal

- Construct a **good** predictor \hat{f} from the training data.
- Need to specify the meaning of good.
- Classification and regression are almost the **same** problem!

Loss function for a generic predictor

- **Loss function:** $\ell(Y, f(\underline{X}))$ measures the goodness of the prediction of Y by $f(\underline{X})$
- Examples:
 - 0/1 loss: $\ell(Y, f(\underline{X})) = \mathbf{1}_{Y \neq f(\underline{X})}$
 - Quadratic loss: $\ell(Y, f(\underline{X})) = |Y - f(\underline{X})|^2$

Risk function

- Risk measured as the average loss for a new couple:

$$\mathcal{R}(f) = \mathbb{E}_{(X, Y) \sim \mathbb{P}}[\ell(Y, f(\underline{X}))]$$

- Examples:
 - 0/1 loss: $\mathbb{E}[\ell(Y, f(\underline{X}))] = \mathbb{P}(Y \neq f(\underline{X}))$
 - Quadratic loss: $\mathbb{E}[\ell(Y, f(\underline{X}))] = \mathbb{E}[|Y - f(\underline{X})|^2]$

- **Beware:** As \hat{f} depends on \mathcal{D}_n , $\mathcal{R}(\hat{f})$ is a random variable!

- The best solution f^* (which is independent of \mathcal{D}_n) is

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{R}(f) = \arg \min_{f \in \mathcal{F}} \mathbb{E}[\ell(Y, f(\underline{X}))] = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\underline{X}} \left[\mathbb{E}_{Y|\underline{X}}[\ell(Y, f(\underline{X}))] \right]$$

Bayes Predictor (explicit solution)

- In binary classification with 0 – 1 loss:

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = +1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ & \Leftrightarrow \mathbb{P}(Y = +1|\underline{X}) \geq 1/2 \\ -1 & \text{otherwise} \end{cases}$$

- In regression with the quadratic loss

$$f^*(\underline{X}) = \mathbb{E}[Y|\underline{X}]$$

Issue: Solution requires to **know** $\mathbb{E}[Y|\underline{X}]$ for all values of \underline{X} !

Machine Learning

- Learn a rule to construct a **predictor** $\hat{f} \in \mathcal{F}$ from the training data \mathcal{D}_n s.t. **the risk** $\mathcal{R}(\hat{f})$ is **small on average** or with high probability with respect to \mathcal{D}_n .
- In practice, the rule should be an algorithm!

Canonical example: Empirical Risk Minimizer

- One restricts f to a subset of functions $\mathcal{S} = \{f_\theta, \theta \in \Theta\}$
- One replaces the minimization of the average loss by the minimization of the empirical loss

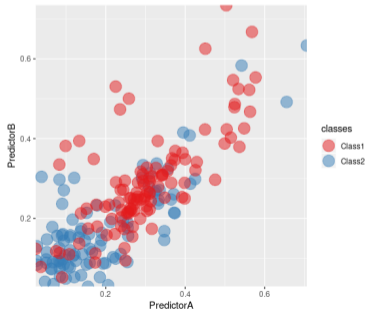
$$\hat{f} = f_{\hat{\theta}} = \operatorname{argmin}_{f_\theta, \theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_\theta(\underline{X}_i))$$

- Examples:
 - Linear regression
 - Linear classification with

$$\mathcal{S} = \{\underline{x} \mapsto \operatorname{sign}\{\underline{x}^\top \beta + \beta^{(0)}\} / \beta \in \mathbb{R}^d, \beta^{(0)} \in \mathbb{R}\}$$

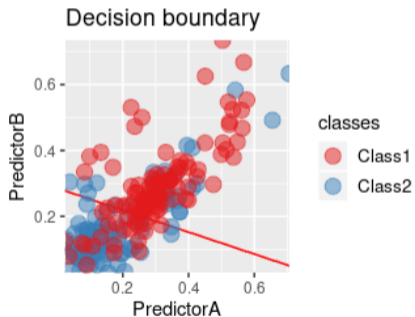
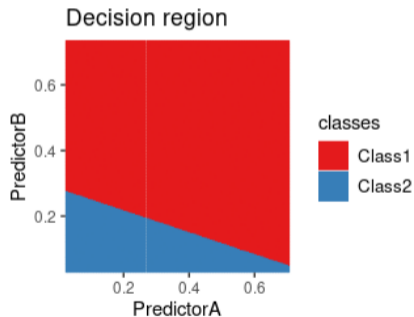
Synthetic Dataset

- Two features/covariates.
- Two classes.
- Dataset from *Applied Predictive Modeling*, M. Kuhn and K. Johnson, Springer
- Numerical experiments with R and the `{caret}` package.



Example: Linear Discrimination

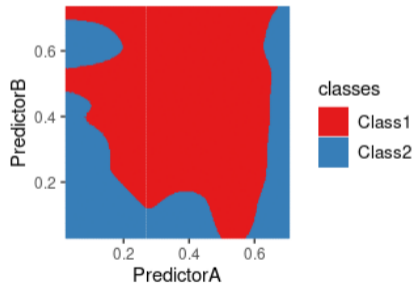
Logistic



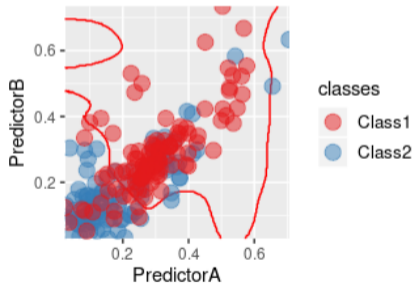
Example: More Complex Model

Naive Bayes with kernel density estimates

Decision region

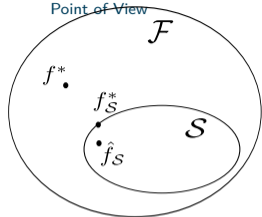


Decision boundary



Bias-Variance Dilemma

- General setting:
 - $\mathcal{F} = \{\text{measurable functions } \mathcal{X} \rightarrow \mathcal{Y}\}$
 - Best solution: $f^* = \operatorname{argmin}_{f \in \mathcal{F}} \mathcal{R}(f)$
 - Class $\mathcal{S} \subset \mathcal{F}$ of functions
 - Ideal target in \mathcal{S} : $f_S^* = \operatorname{argmin}_{f \in \mathcal{S}} \mathcal{R}(f)$
 - Estimate in \mathcal{S} : \hat{f}_S obtained with some procedure

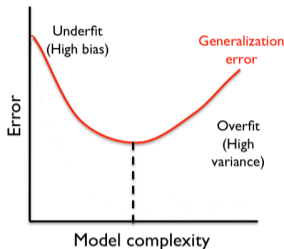


Approximation error and estimation error (Bias/Variance)

$$\mathcal{R}(\hat{f}_S) - \mathcal{R}(f^*) = \underbrace{\mathcal{R}(f_S^*) - \mathcal{R}(f^*)}_{\text{Approximation error}} + \underbrace{\mathcal{R}(\hat{f}_S) - \mathcal{R}(f_S^*)}_{\text{Estimation error}}$$

- Approx. error can be large if the model \mathcal{S} is not suitable.
- Estimation error can be large if the model is complex.

Under-fitting / Over-fitting Issue



- Different behavior for different model complexity
- **Low complexity model** are easily learned but the approximation error (**bias**) may be large (**Under-fit**).
- **High complexity model** may contain a good ideal target but the estimation error (**variance**) can be large (**Over-fit**)

Bias-variance trade-off \iff avoid **overfitting** and **underfitting**

- **Rk:** Better to think in term of method (including feature engineering and specific algorithm) rather than only of model.

Probabilistic and Optimization Framework

How to find a good function f with a *small* risk

$$\mathcal{R}(f) = \mathbb{E}[\ell(Y, f(\underline{X}))] \quad ?$$

Canonical approach: $\hat{f}_S = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(\underline{X}_i))$

Problems

- How to choose \mathcal{S} ?
- How to compute the minimization?

A Probabilistic Point of View

Solution: For \underline{X} , estimate $Y|\underline{X}$ plug this estimate in the Bayes classifier:
(Generalized) Linear Models, Kernel methods, k -nn, Naive Bayes, Tree, Bagging...

An Optimization Point of View

Solution: If necessary replace the loss ℓ by an upper bound $\bar{\ell}$ and minimize the empirical loss: **SVR, SVM, Neural Network, Tree, Boosting...**

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 **ML Methods: Probabilistic Point of View**
 - Motivation
 - Supervised Learning
 - **A Probabilistic Point of View**
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- Gradient Descent Step
- Non-Convex Setting
- References
- 5 **ML Methods: Neural Networks and Deep Learning**
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 **ML Methods: Trees and Ensemble Methods**
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 **Unsupervised Learning: Dimension Reduction and Clustering**
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- Applications to Text
- References
- 8 **Statistical Learning: PAC-Bayesian Approach and Complexity Theory**
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 **References**

Logistic Regression

- Let $f_{\theta}(\underline{X}) = \underline{X}^{\top} \beta + \beta^{(0)}$ with $\theta = (\beta, \beta^{(0)})$.
- Let $\mathbb{P}_{\theta}(Y = 1|\underline{X}) = e^{-f_{\theta}(\underline{X})} / (1 + e^{f_{\theta}(\underline{X})})$
- Estimate θ by $\hat{\theta}$ using a Maximum Likelihood.
- Classify using $\mathbb{P}_{\hat{\theta}}(Y = 1|\underline{X}) > 1/2$

k Nearest Neighbors

- For any \underline{X}' , define $\mathcal{V}_{\underline{X}'}$ as the k closest samples X_i from the dataset.
- Compute a score $g_k = \sum_{X_i \in \mathcal{V}_{\underline{X}'}} \mathbf{1}_{Y_i=k}$
- Classify using $\arg \max g_k$ (majority vote).

Quadratic Discriminant Analysis

- For each class, estimate the mean μ_k and the covariance matrix Σ_k .
- Estimate the proportion $\mathbb{P}(Y = k)$ of each class.
- Compute a score $\ln(\mathbb{P}(\underline{X}|Y = k)) + \ln(\mathbb{P}(Y = k))$

$$g_k(\underline{X}) = -\frac{1}{2}(\underline{X} - \mu_k)^\top \Sigma_k^{-1}(\underline{X} - \mu_k) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma_k|) + \ln(\mathbb{P}(Y = k))$$

- Classify using $\arg \max g_k$
- Those three methods rely on a similar heuristic: the probabilistic point of view!

- The best solution f^* (which is independent of \mathcal{D}_n) is

$$f^* = \arg \min_{f \in \mathcal{F}} R(f) = \arg \min_{f \in \mathcal{F}} \mathbb{E}[\ell(Y, f(\underline{X}))] = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\underline{X}} \left[\mathbb{E}_{Y|\underline{X}}[\ell(Y, f(\underline{X}))] \right]$$

Bayes Predictor (explicit solution)

- In binary classification with 0 – 1 loss:

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = +1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ & \Leftrightarrow \mathbb{P}(Y = +1|\underline{X}) \geq 1/2 \\ -1 & \text{otherwise} \end{cases}$$

- In regression with the quadratic loss

$$f^*(\underline{X}) = \mathbb{E}[Y|\underline{X}]$$

Issue: Explicit solution requires to **know** $Y|\underline{X}$ (or $\mathbb{E}[Y|\underline{X}]$) for all values of \underline{X} !

- **Idea:** Estimate $Y|\underline{X}$ by $\widehat{Y|\underline{X}}$ and plug it the Bayes classifier.

Plugin Bayes Predictor

- In binary classification with 0 – 1 loss:

$$\widehat{f}(\underline{X}) = \begin{cases} +1 & \text{if } \overline{\mathbb{P}(Y = +1|\underline{X})} \geq \overline{\mathbb{P}(Y = -1|\underline{X})} \\ & \Leftrightarrow \overline{\mathbb{P}(Y = +1|\underline{X})} \geq 1/2 \\ -1 & \text{otherwise} \end{cases}$$

- In regression with the quadratic loss

$$\widehat{f}(\underline{X}) = \mathbb{E}[\widehat{Y|\underline{X}}]$$

- **Rk:** Direct estimation of $\mathbb{E}[Y|\underline{X}]$ by $\widehat{\mathbb{E}[Y|\underline{X}]}$ also possible. . .

- How to estimate $Y|\underline{X}$?

Three main heuristics

- **Parametric Conditional modeling:** Estimate the law of $Y|\underline{X}$ by a **parametric** law $\mathcal{L}_\theta(\underline{X})$: *(generalized) linear regression...*
- **Non Parametric Conditional modeling:** Estimate the law of $Y|\underline{X}$ by a **non parametric** estimate: *kernel methods, loess, nearest neighbors...*
- **Fully Generative modeling:** Estimate the law of (\underline{X}, Y) and use the **Bayes formula** to deduce an estimate of $Y|\underline{X}$: *LDA/QDA, Naive Bayes...*
- **Rk:** Direct estimation of $\mathbb{E}[Y|\underline{X}]$ by $\widehat{\mathbb{E}[Y|\underline{X}]}$ also possible...

- **Input:** a data set \mathcal{D}_n
Learn $Y|\underline{X}$ or equivalently $\mathbb{P}(Y = k|\underline{X})$ (using the data set) and plug this estimate in the Bayes classifier
- **Output:** a classifier $\hat{f} : \mathbb{R}^d \rightarrow \{-1, 1\}$

$$\hat{f}(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(\widehat{Y} = 1|\underline{X}) \geq \mathbb{P}(\widehat{Y} = -1|\underline{X}) \\ -1 & \text{otherwise} \end{cases}$$

- Can we guaranty that the classifier is good if $Y|\underline{X}$ is well estimated?

Theorem

- If $\hat{f} = \text{sign}(2\hat{p}_{+1} - 1)$ then

$$\begin{aligned}\mathbb{E} \left[\ell^{0,1}(Y, \hat{f}(\underline{X})) \right] - \mathbb{E} \left[\ell^{0,1}(Y, f^*(\underline{X})) \right] \\ \leq \mathbb{E} \left[\|\widehat{Y|\underline{X}} - Y|\underline{X}\|_1 \right] \\ \leq \left(\mathbb{E} \left[2\text{KL}(Y|\underline{X}, \widehat{Y|\underline{X}}) \right] \right)^{1/2}\end{aligned}$$

- If one estimates $\mathbb{P}(Y = 1|\underline{X})$ well then one estimates f^* well!
- Link between a *conditional density estimation* task and a *classification* one!
- **Rk:** In general, the conditional density estimation task is more complicated as one should be good for all values of $\mathbb{P}(Y = 1|\underline{X})$ while the classification task focus on values around 1/2 for the 0/1 loss!
- In **regression**, (often) direct control of the quadratic loss. . .

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 **ML Methods: Probabilistic Point of View**
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - **Parametric Conditional Density Modeling**
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- **Idea:** Estimate directly $Y|\underline{X}$ by a parametric conditional density $\mathbb{P}_\theta(Y|\underline{X})$.

Maximum Likelihood Approach

- Classical choice for θ :

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} - \sum_{i=1}^n \log \mathbb{P}_\theta(Y_i|\underline{X}_i)$$

- **Goal:** *Minimize* the Kullback-Leibler divergence between the conditional law of $Y|\underline{X}$ and $\mathbb{P}_\theta(Y|\underline{X})$

$$\mathbb{E}[\text{KL}(Y|\underline{X}, \mathbb{P}_\theta(Y|\underline{X}))]$$

- **Rk:** This is often not (exactly) the learning task!
- Large choice for the family $\{\mathbb{P}_\theta(Y|\underline{X})\}$ but depends on \mathcal{Y} (and \mathcal{X}).
- **Regression:** One can also model directly $\mathbb{E}[Y|\underline{X}]$ by $f_\theta(\underline{X})$ and estimate it with a least-squares criterion...

Linear Models

- **Classical choice:** $\theta = (\theta', \varphi)$

$$\mathbb{P}_{\theta}(Y|\underline{X}) = \mathbb{P}_{\underline{X}^{\top}\beta, \varphi}(Y)$$

- **Very strong assumption!**
- Classical examples:
 - Binary variable: logistic, probit...
 - Discrete variable: multinomial logistic regression...
 - Integer variable: Poisson regression...
 - Continuous variable: Gaussian regression...

Plugin Linear Classification

- Model $\mathbb{P}(Y = +1|\underline{X})$ by $h(\underline{X}^\top \beta + \beta^{(0)})$ with h non decreasing.
- $h(\underline{X}^\top \beta + \beta^{(0)}) > 1/2 \Leftrightarrow \underline{X}^\top \beta + \beta^{(0)} - h^{-1}(1/2) > 0$
- Linear Classifier: $\text{sign}(\underline{X}^\top \beta + \beta^{(0)} - h^{-1}(1/2))$

Plugin Linear Classifier Estimation

- Classical choice for h :

$$h(t) = \frac{e^t}{1 + e^t}$$

logit or logistic

$$h(t) = F_{\mathcal{N}}(t)$$

probit

$$h(t) = 1 - e^{-e^t}$$

log-log

- Choice of the *best* β from the data.

Probabilistic Model

- By construction, $Y|\underline{X}$ follows $\mathcal{B}(\mathbb{P}(Y = +1|\underline{X}))$
- Approximation of $Y|\underline{X}$ by $\mathcal{B}(h(\underline{x}^\top \beta + \beta^{(0)}))$
- *Natural* probabilistic choice for β : maximum likelihood estimate.
- *Natural* probabilistic choice for β : β approximately minimizing a distance between $\mathcal{B}(h(\underline{x}^\top \beta))$ and $\mathcal{B}(\mathbb{P}(Y = 1|\underline{X}))$.

Maximum Likelihood Approach

- Minimization of the negative log-likelihood:

$$-\sum_{i=1}^n \log(\mathbb{P}(Y_i|\underline{X}_i)) = -\sum_{i=1}^n \left(\mathbf{1}_{Y_i=1} \log(h(\underline{X}_i^\top \beta)) + \mathbf{1}_{Y_i=-1} \log(1 - h(\underline{X}_i^\top \beta)) \right)$$

- Minimization possible if h is regular...

KL Distance and negative log-likelihood

- *Natural* distance: Kullback-Leibler divergence

$$\begin{aligned} & \text{KL}(\mathcal{B}(\mathbb{P}(Y = 1|\underline{X})), \mathcal{B}(h(\underline{X}^\top \beta))) \\ &= \mathbb{E}_{\underline{X}} \left[\mathbb{P}(Y = 1|\underline{X}) \log \frac{\mathbb{P}(Y = 1|\underline{X})}{h(\underline{X}^\top \beta)} \right. \\ & \quad \left. + \mathbb{P}(Y = -1|\underline{X}) \log \frac{1 - \mathbb{P}(Y = 1|\underline{X})}{1 - h(\underline{X}^\top \beta)} \right] \\ &= \mathbb{E}_{\underline{X}} \left[-\mathbb{P}(Y = 1|\underline{X}) \log(h(\underline{X}^\top \beta)) \right. \\ & \quad \left. - \mathbb{P}(Y = -1|\underline{X}) \log(1 - h(\underline{X}^\top \beta)) \right] + C_{\underline{X}, Y} \end{aligned}$$

- Empirical counterpart = negative log-likelihood (up to $1/n$ factor):

$$-\frac{1}{n} \sum_{i=1}^n \left(\mathbf{1}_{Y_i=1} \log(h(\underline{X}_i^\top \beta)) + \mathbf{1}_{Y_i=-1} \log(1 - h(\underline{X}_i^\top \beta)) \right)$$

Logistic Regression and Odd

- Logistic model: $h(t) = \frac{e^t}{1+e^t}$ (most *natural* choice...)

- The Bernoulli law $\mathcal{B}(h(t))$ satisfies then

$$\frac{\mathbb{P}(Y = 1)}{\mathbb{P}(Y = -1)} = e^t \Leftrightarrow \log \frac{\mathbb{P}(Y = 1)}{\mathbb{P}(Y = -1)} = t$$

- Interpretation in term of odd.
- Logistic model: linear model on the logarithm of the odd

$$\log \frac{\mathbb{P}(Y = 1|\underline{X})}{\mathbb{P}(Y = -1|\underline{X})} = \underline{X}^\top \beta$$

Associated Classifier

- Plugin strategy:

$$f_\beta(\underline{X}) = \begin{cases} 1 & \text{if } \frac{e^{\underline{X}^\top \beta}}{1+e^{\underline{X}^\top \beta}} > 1/2 \Leftrightarrow \underline{X}^\top \beta > 0 \\ -1 & \text{otherwise} \end{cases}$$

Likelihood Rewriting

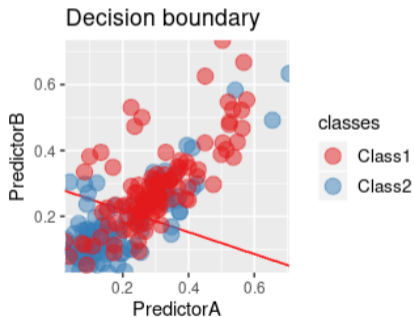
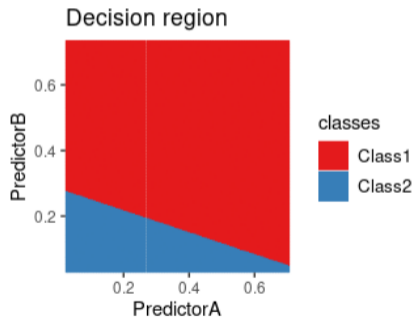
- Negative log-likelihood:

$$\begin{aligned} & -\frac{1}{n} \sum_{i=1}^n \left(\mathbf{1}_{Y_i=1} \log(h(\underline{X}_i^\top \beta)) + \mathbf{1}_{Y_i=-1} \log(1 - h(\underline{X}_i^\top \beta)) \right) \\ & = -\frac{1}{n} \sum_{i=1}^n \left(\mathbf{1}_{Y_i=1} \log \frac{e^{\underline{X}_i^\top \beta}}{1 + e^{\underline{X}_i^\top \beta}} + \mathbf{1}_{Y_i=-1} \log \frac{1}{1 + e^{\underline{X}_i^\top \beta}} \right) \\ & = \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-Y_i(\underline{X}_i^\top \beta)} \right) \end{aligned}$$

- Convex and smooth function of β
- Easy optimization.

Example: Logistic

Logistic



Transformed Representation

- From \underline{X} to $\Phi(\underline{X})$!
- New description of \underline{X} leads to a different **linear** model:

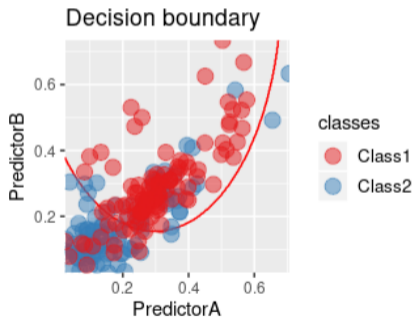
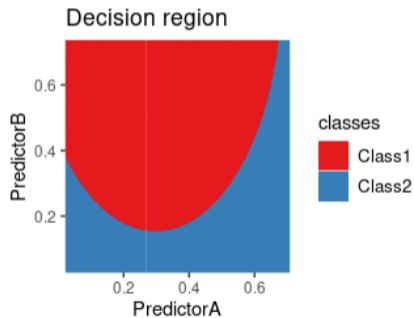
$$f_{\beta}(\underline{X}) = \Phi(\underline{X})^{\top} \beta$$

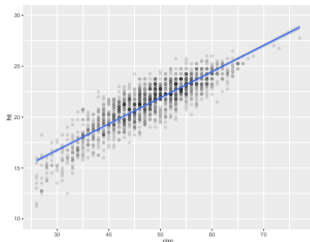
Feature Design

- Art of choosing Φ .
- Examples:
 - Renormalization, (domain specific) transform
 - Basis decomposition
 - Interaction between different variables. . .

Example: Quadratic Logistic

Quadratic Logistic





Gaussian Linear Model

- **Model:** $Y|\underline{X} \sim \mathcal{N}(\underline{X}^\top \beta, \sigma^2)$ plus independence
- Probably the most classical model of all time!
- Maximum Likelihood with explicit formulas for the two parameters.
- In regression, estimation of $\mathbb{E}[Y|\underline{X}]$ is sufficient: other/no model for the noise possible.

Generalized Linear Model

- Model entirely characterized by its mean (up to a scalar nuisance parameter) ($v(\mathbb{E}_\theta[Y]) = \theta$ with v invertible).
- Exponential family: Probability law family P_θ such that the density can be written

$$f(y, \theta, \varphi) = e^{\frac{y\theta - v(\theta)}{\varphi} + w(y, \varphi)}$$

where φ is a nuisance parameter and w a function independent of θ .

- Examples:
 - Gaussian: $f(y, \theta, \varphi) = e^{-\frac{y\theta - \theta^2/2}{\varphi} - \frac{y^2/2}{\varphi}}$
 - Bernoulli: $f(y, \theta) = e^{y\theta - \ln(1+e^\theta)}$ ($\theta = \ln p/(1-p)$)
 - Poisson: $f(y, \theta) = e^{(y\theta - e^\theta) + \ln(y!)}$ ($\theta = \ln \lambda$)
- Linear Conditional model: $Y|\underline{X} \sim P_{\underline{x}^\top \beta} \dots$

- ML fit of the parameters

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 **ML Methods: Probabilistic Point of View**
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - **Non Parametric Conditional Density Modeling**
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- Gradient Descent Step
- Non-Convex Setting
- References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References



- **Idea:** Estimate $Y|\underline{X}$ or $\mathbb{E}[Y|\underline{X}]$ directly without resorting to an explicit parametric model.

Non Parametric Conditional Estimation

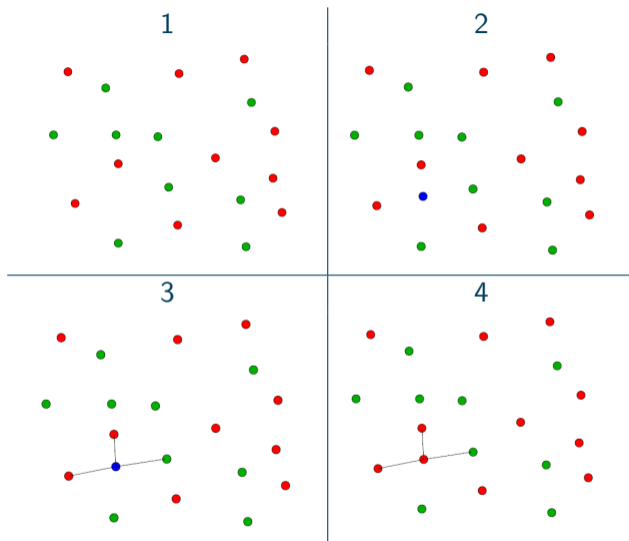
- Two heuristics:
 - $Y|\underline{X}$ (or $\mathbb{E}[Y|\underline{X}]$) is almost constant (or simple) in a neighborhood of \underline{X} . (Kernel methods)
 - $Y|\underline{X}$ (or $\mathbb{E}[Y|\underline{X}]$) can be approximated by a model whose dimension depends on the complexity and the number of observation. (Quite similar to parametric model plus model selection...)
- Focus on **kernel methods!**

- **Idea:** The behavior of $Y|\underline{X}$ is locally *constant* or simple!

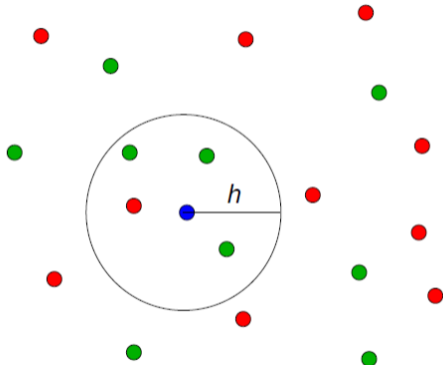
Kernel

- Choose a kernel K (think of a weighted neighborhood).
 - For each \tilde{X} , compute a simple localized estimate of $Y|\underline{X}$
 - Use this local estimate to take the decision
-
- In regression, estimation of $\mathbb{E}[Y|\underline{X}]$ is sufficient.

Example: k Nearest-Neighbors (with $k = 3$)



Example: k Nearest-Neighbors (with $k = 4$)



- Neighborhood $\mathcal{V}_{\underline{x}}$ of \underline{x} : k learning samples closest from \underline{x} .

k -NN as local conditional density estimate

$$\mathbb{P}(\widehat{Y} = 1 | \underline{X}) = \frac{\sum_{\underline{X}_i \in \mathcal{V}_{\underline{X}}} \mathbf{1}_{\{Y_i = +1\}}}{|\mathcal{V}_{\underline{X}}|}$$

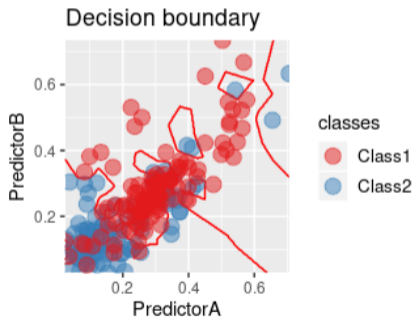
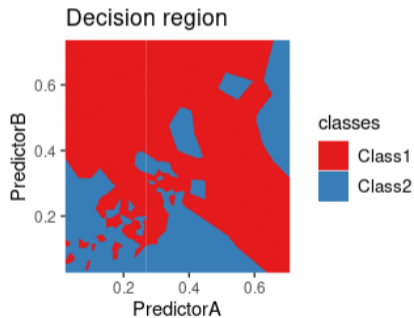
- KNN Classifier:

$$\widehat{f}_{KNN}(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(\widehat{Y} = 1 | \underline{X}) \geq \mathbb{P}(\widehat{Y} = -1 | \underline{X}) \\ -1 & \text{otherwise} \end{cases}$$

- **Lazy learning:** all the computations have to be done at prediction time.
- **Remark:** You can also use your favorite kernel estimator...

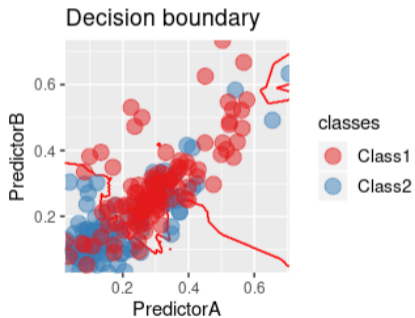
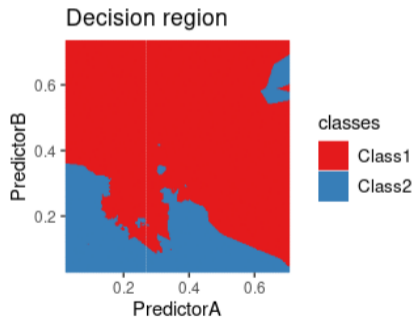
Example: KNN

k-NN with k=1



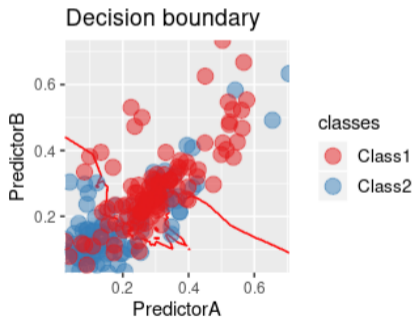
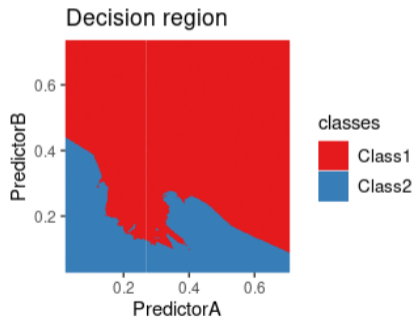
Example: KNN

k-NN with k=5



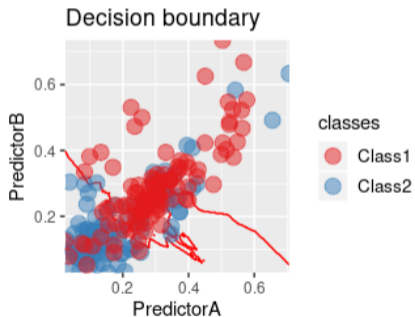
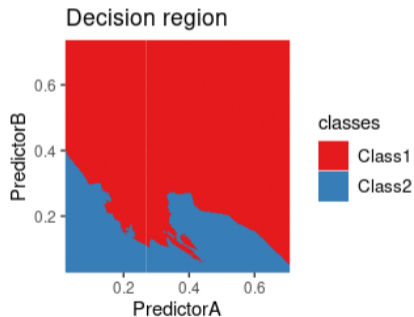
Example: KNN

k-NN with k=9



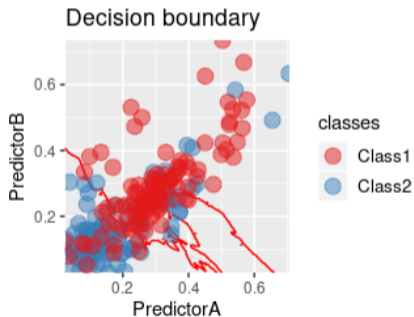
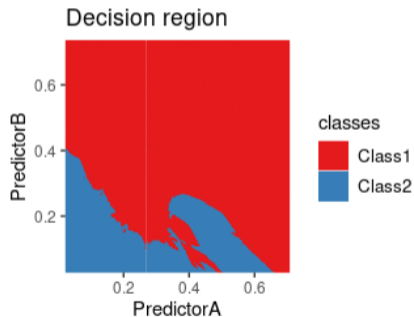
Example: KNN

k-NN with $k=13$



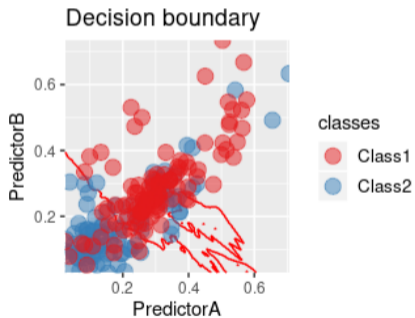
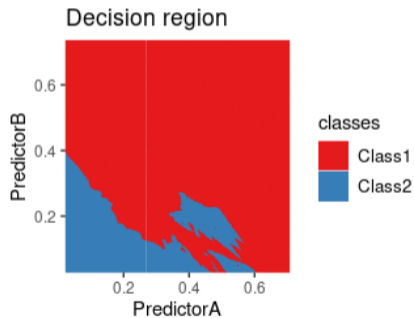
Example: KNN

k-NN with $k=17$



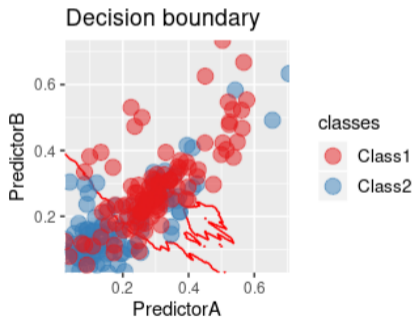
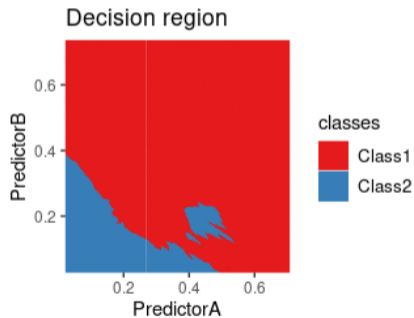
Example: KNN

k-NN with k=21



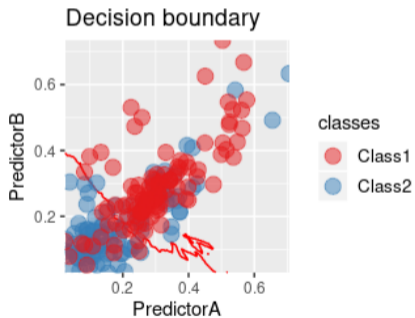
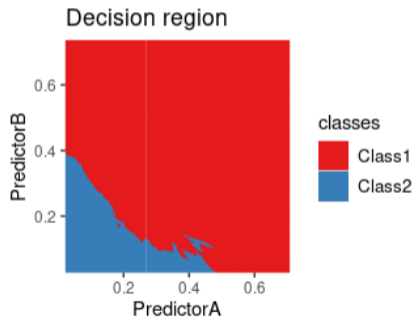
Example: KNN

k-NN with k=25



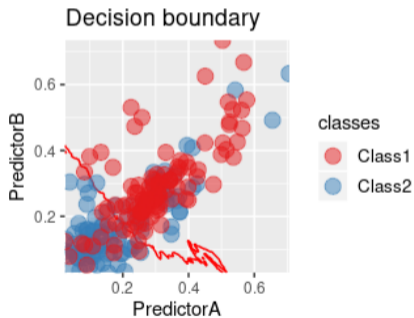
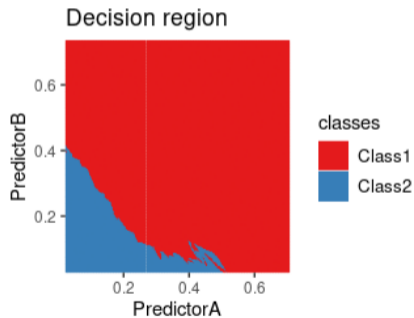
Example: KNN

k-NN with k=29



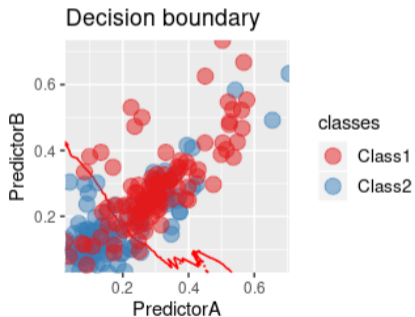
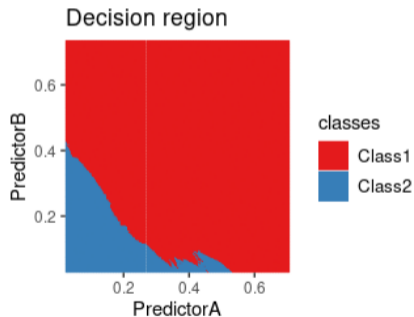
Example: KNN

k-NN with k=33



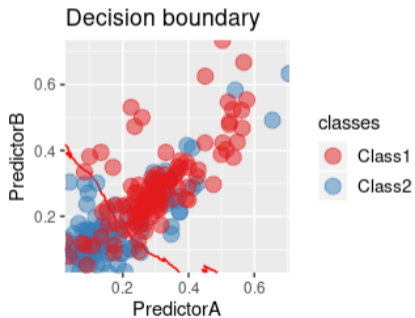
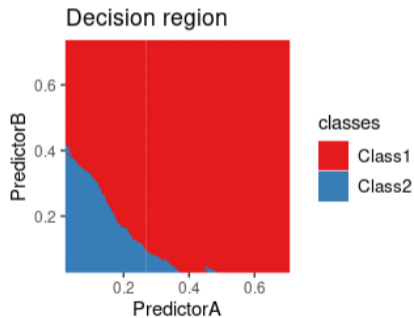
Example: KNN

k-NN with k=37



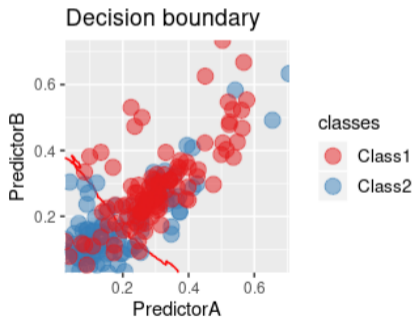
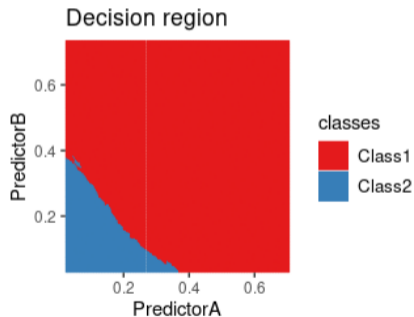
Example: KNN

k-NN with k=45



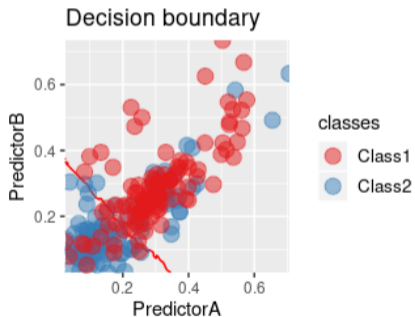
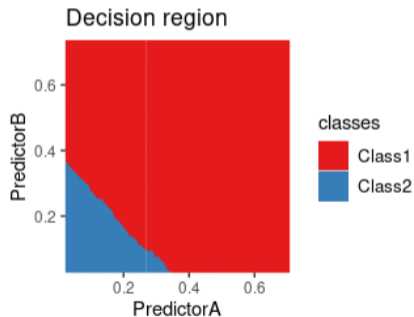
Example: KNN

k-NN with k=53



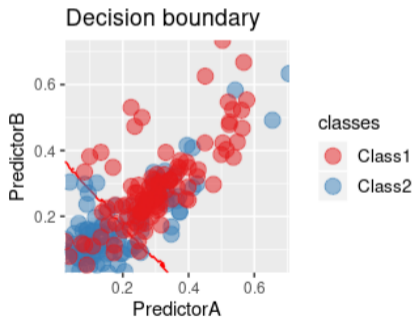
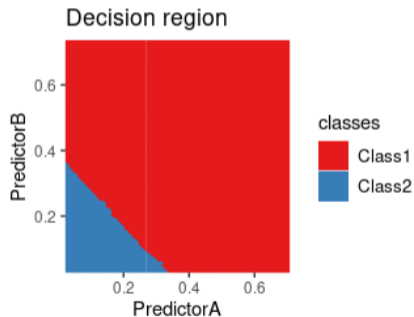
Example: KNN

k-NN with k=61



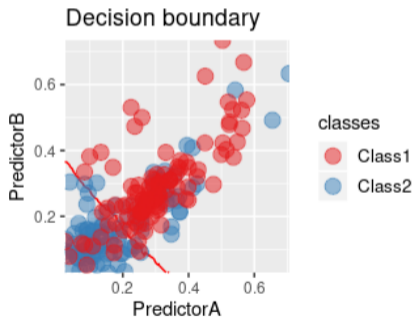
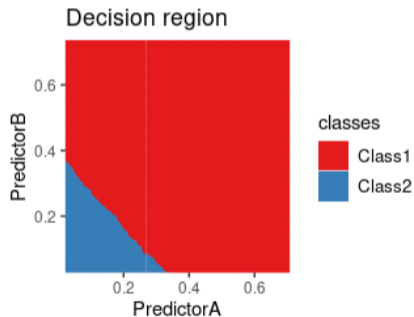
Example: KNN

k-NN with k=69



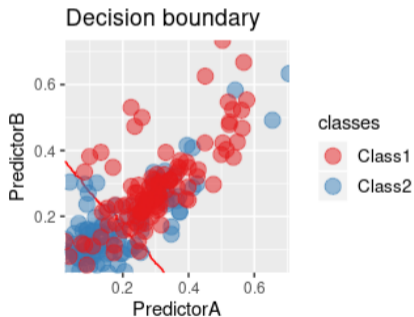
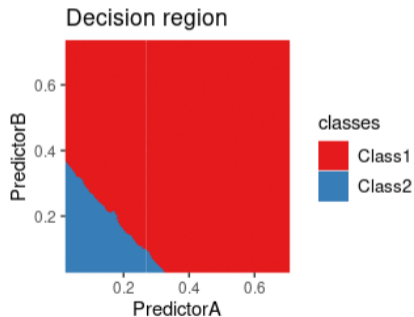
Example: KNN

k-NN with $k=77$



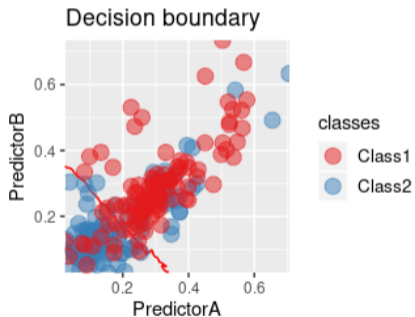
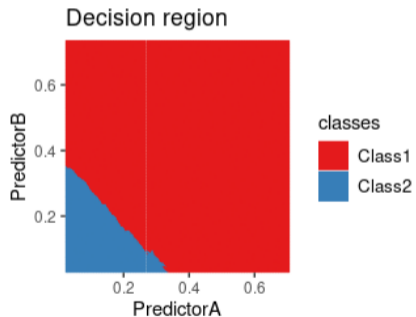
Example: KNN

k-NN with k=85



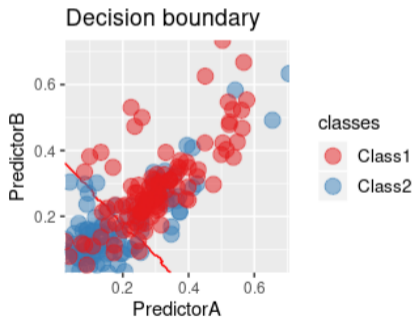
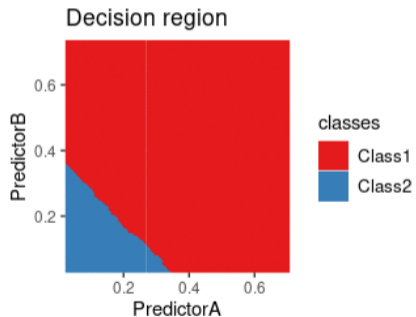
Example: KNN

k-NN with $k=101$



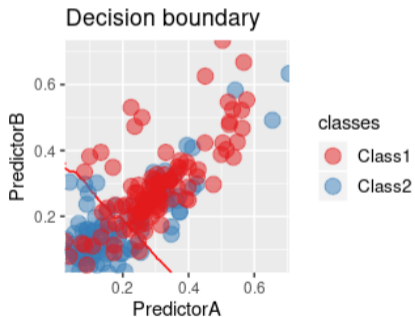
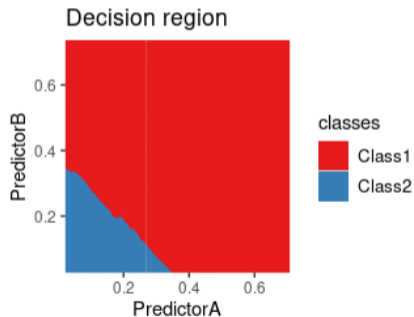
Example: KNN

k-NN with $k=109$



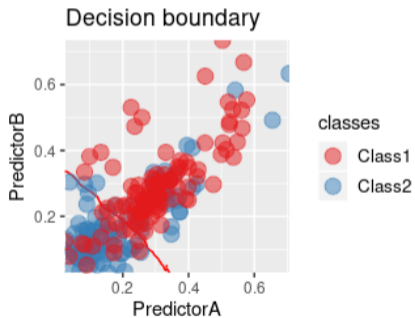
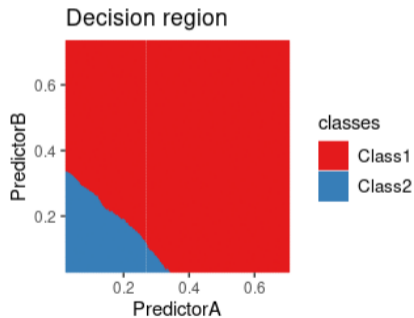
Example: KNN

k-NN with $k=117$



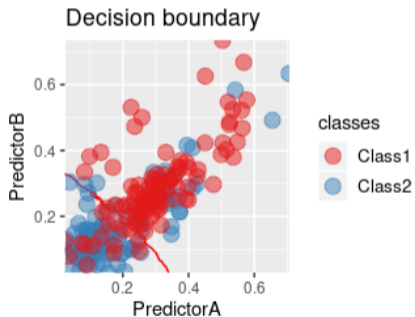
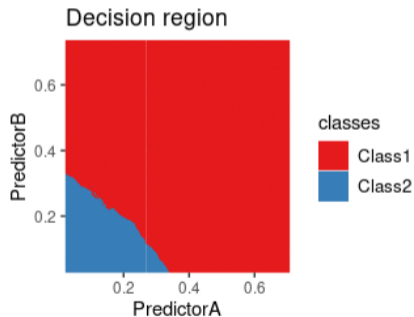
Example: KNN

k-NN with $k=125$



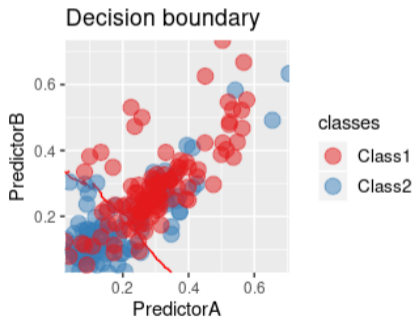
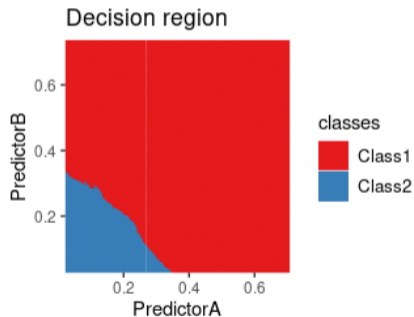
Example: KNN

k-NN with k=133



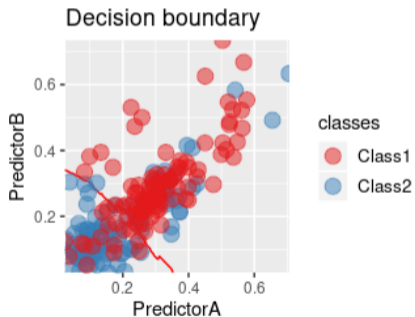
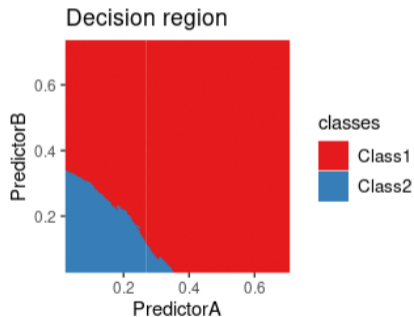
Example: KNN

k-NN with k=141



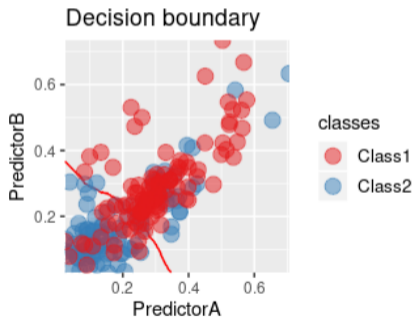
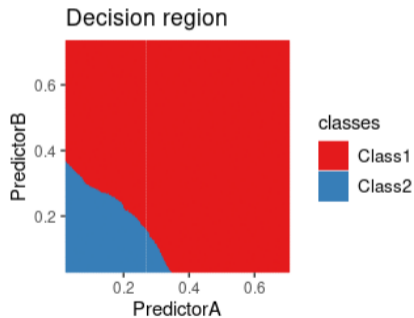
Example: KNN

k-NN with $k=149$



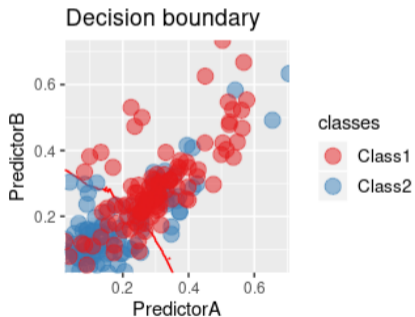
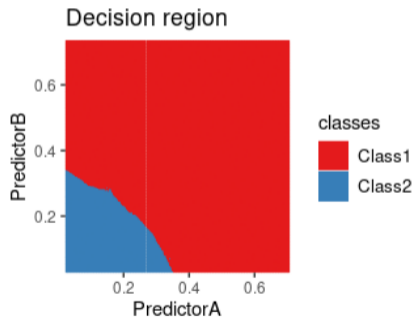
Example: KNN

k-NN with $k=157$



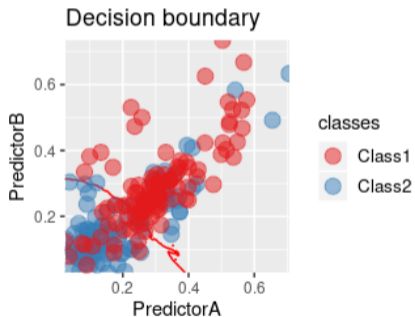
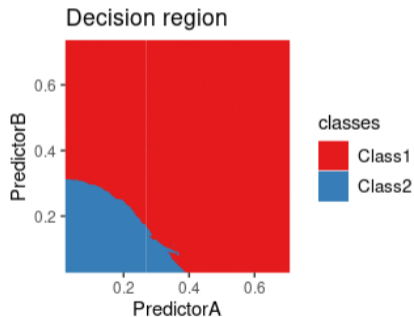
Example: KNN

k-NN with k=165



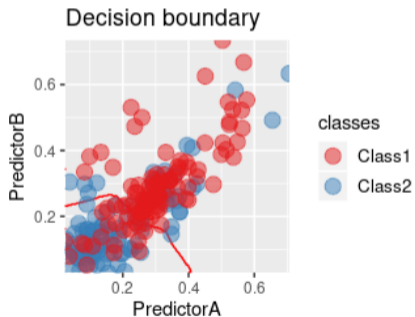
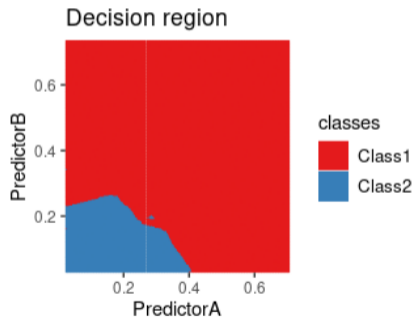
Example: KNN

k-NN with $k=173$



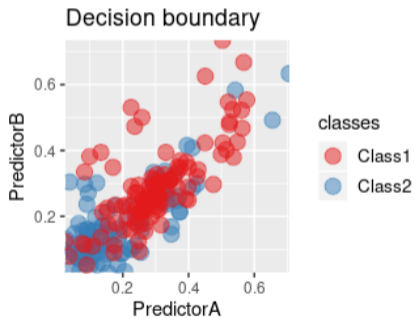
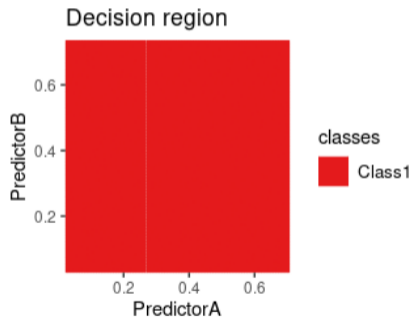
Example: KNN

k-NN with $k=181$



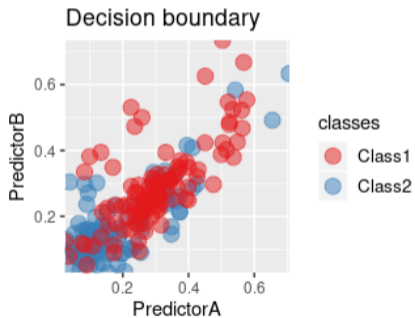
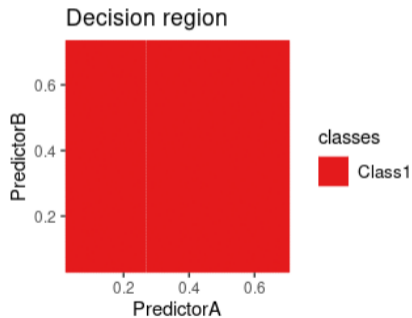
Example: KNN

k-NN with $k=189$



Example: KNN

k-NN with $k=197$



A naive idea

- $\mathbb{E}[Y|\underline{X}]$ can be approximated by a local average:

$$\hat{f}(\underline{X}) = \frac{1}{|\{\underline{X}_i \in \mathcal{N}(\underline{X})\}|} \sum_{\underline{X}_i \in \mathcal{N}(\underline{X})} Y_i$$

where $\mathcal{B}(\underline{X})$ is a neighborhood of \underline{X} .

- **Heuristic:**

- If $\underline{X} \rightarrow \mathbb{E}[Y|\underline{X}]$ is regular then

$$\mathbb{E}[Y|\underline{X}] \simeq \mathbb{E}[\mathbb{E}[Y|\underline{X}'] | \underline{X}' \in \mathcal{N}(\underline{X})] = \mathbb{E}[Y | \underline{X}' \in \mathcal{N}(\underline{X})]$$

- Replace an expectation by an empirical average:

$$\mathbb{E}[Y | \underline{X}' \in \mathcal{N}(\underline{X})] \simeq \frac{1}{|\{\underline{X}_i \in \mathcal{N}(\underline{X})\}|} \sum_{\underline{X}_i \in \mathcal{N}(\underline{X})} Y_i$$

Neighborhood and Size

- Most classical choice: $\mathcal{N}(\underline{X}) = \{ \underline{X}', \|\underline{X} - \underline{X}'\| \leq h \}$ where $\|\cdot\|$ is a (pseudo) norm and h a size (bandwidth) parameter.
- In principle, the norm and h could vary with \underline{X} , and the norm can be replaced by a (pseudo) distance.
- Focus here on a fixed distance with a fixed bandwidth h cased.

Bandwidth Heuristic

- A **large bandwidth** ensures that the average is taken on many samples and thus the **variance is small**...
- A **small bandwidth** is thus that the approximation $\mathbb{E}[Y|\underline{X}] \simeq \mathbb{E}[Y|\underline{X}' \in \mathcal{N}(\underline{X})]$ is more accurate (**small bias**).

Weighted Local Average

- Replace the neighborhood $\mathcal{N}(\underline{X})$ by a decaying **window function** $w(\underline{X}, \underline{X}')$.
- $\mathbb{E}[Y|\underline{X}]$ can be approximated by a **weighted local average**:

$$\hat{f}(\underline{X}) = \frac{\sum_i w(\underline{X}, \underline{X}'_i) Y_i}{\sum_i w(\underline{X}, \underline{X}'_i)}.$$

Kernel

- Most classical choice: $w(\underline{X}, \underline{X}') = K\left(\frac{\underline{X}-\underline{X}'}{h}\right)$ where h the bandwidth is a scale parameter.
- Examples:
 - **Box kernel**: $K(t) = \mathbf{1}_{\|t\| \leq 1}$ (Neighborhood)
 - **Triangular kernel**: $K(t) = \max(1 - \|t\|, 0)$.
 - **Gaussian kernel**: $K(t) = e^{-t^2/2}$
- **Rk**: K and λK yields the same estimate.

Density Estimation

- How to estimate the density p of \underline{X} with respect to the Lebesgue measure from an i.i.d. sample $(\underline{X}_1, \dots, \underline{X}_n)$.
- **Parametric approach:** density has a known parameterized shape and estimate those parameters.
- **Nonparametric approach:** density has a no known parameterized shape and
 - Approximate it by a parametric one, whose parameters can be estimated
 - Estimate directly the density
- Important **nonparametric statistic topic!**
- Used in generative modeling. . .

Kernel Density Estimation (Parzen)

- Choose a positive kernel K such that $\int K(x)dx = 1$
- Use as an estimate

$$\hat{p}(\underline{X}) = \frac{1}{n} \sum_{i=1}^n K(\underline{X} - \underline{X}_i)$$

- If $K = \frac{1}{Z_h} \mathbf{1}_{\|t\| \leq h}$, easy interpretation as a **local empirical density** of samples!
- General K corresponds to a **smoothed version**.
- Often $K_h(t) = \frac{1}{h^d} K(t/h)$ and let

$$\hat{p}_h(\underline{X}) = \frac{1}{n} \sum_{i=1}^n K_h(\underline{X} - \underline{X}_i)$$

Properties

- **Error decomposition:**

$$\mathbb{E} \left[|p(\underline{X}) - \hat{p}_h(\underline{X})|^2 \right] = \mathbb{E} [p(\underline{X}) - \hat{p}_h(\underline{X})]^2 + \text{Var} [p(\underline{X}) - \hat{p}_h(\underline{X})]$$

- **Bias:**

$$\mathbb{E} [p(\underline{X}) - \hat{p}_h(\underline{X})] = p(\underline{X}) - (K_h * p)(\underline{X})$$

- **Variance:** if p is upper bounded by p_{\max} then

$$\text{Var} [p(\underline{X}) - \hat{p}_h(\underline{X})] \leq \frac{p_{\max} \int K_h^2(x) dx}{nh^d}$$



Bandwidth choice

- A small h leads to a small bias but a large variance. . .
- A large h leads to a small variance but a large bias. . .
- Theoretical analysis possible!

Nadaraya-Watson Heuristic

- Provided all the **densities** exist

$$\mathbb{E}[Y|\underline{X}] = \frac{\int Y p(\underline{X}, Y) dY}{\int p(Y, \underline{X}) dY} = \frac{\int Y p(\underline{X}, Y) dY}{p(\underline{X})}$$

- Replace the unknown densities by their **estimates**:

$$\hat{p}(\underline{X}) = \frac{1}{n} \sum_{i=1}^n K(\underline{X} - \underline{X}_i)$$

$$\hat{p}(\underline{X}, Y) = \frac{1}{n} \sum_{i=1}^n K(\underline{X} - \underline{X}_i) K'(Y - Y_i)$$

- Now if K' is a kernel such that $\int Y K'(Y) dY = 0$ then

$$\int Y \hat{p}(\underline{X}, Y) dY = \frac{1}{n} \sum_{i=1}^n K(\underline{X} - \underline{X}_i) Y_i$$

Nadaraya-Watson

- Resulting estimator of $\mathbb{E}[Y|\underline{X}]$

$$\hat{f}(\underline{X}) = \frac{\sum_{i=1}^n Y_i K_h(\underline{X} - \underline{X}_i)}{\sum_{i=1}^n K_h(\underline{X} - \underline{X}_i)}$$

- Same **local weighted average** estimator!

Bandwidth Choice

- Bandwidth h of K allows to **balance between bias and variance**.
- Theoretical analysis of the error is possible.
- The smoother the densities the easier the estimation but the optimal bandwidth depends on the unknown regularity!

Another Point of View on Kernel

- Nadaraya-Watson estimator:

$$\hat{f}(\underline{X}) = \frac{\sum_{i=1}^n Y_i K_h(\underline{X} - \underline{X}_i)}{\sum_{i=1}^n K_h(\underline{X} - \underline{X}_i)}$$

- Can be view as a **minimizer** of

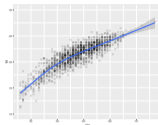
$$\sum_{i=1}^n |Y_i - \beta|^2 K_h(\underline{X} - \underline{X}_i)$$

- **Local regression** of order 0.

Local Linear Model

- Estimate $\mathbb{E}[Y|\underline{X}]$ by $\hat{f}(\underline{X}) = \phi(\underline{X})^\top \hat{\beta}(\underline{X})$ where ϕ is any function of \underline{X} and $\hat{\beta}(\underline{X})$ is the minimizer of

$$\sum_{i=1}^n |Y_i - \phi(\underline{X}_i)^\top \beta|^2 K_h(\underline{X} - \underline{X}_i).$$



1D Nonparametric Regression

- Assume that $\underline{X} \in \mathbb{R}$ and let $\phi(\underline{X}) = (1, \underline{X}, \dots, \underline{X}^d)$.
- **LOESS estimate:** $\hat{f}(\underline{X}) = \sum_{j=0}^d \hat{\beta}(\underline{X}^{(j)}) \underline{X}^j$ with $\hat{\beta}(\underline{X})$ minimizing

$$\sum_{i=1}^n |Y_i - \sum_{j=0}^d \beta^{(j)} \underline{X}_i^j|^2 K_h(\underline{X} - \underline{X}_i).$$

- Most classical kernel used: Tricubic kernel

$$K(t) = \max(1 - |t|^3, 0)^3$$

- Most classical degree: 2...
- Local bandwidth choice such that a proportion of points belongs to the window.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 **ML Methods: Probabilistic Point of View**
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - **Generative Modeling**
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 **ML Methods: Neural Networks and Deep Learning**
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 **ML Methods: Trees and Ensemble Methods**
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 **Unsupervised Learning: Dimension Reduction and Clustering**
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 6 **Dimension Reduction**
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 6 **Clustering**
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 6 **Generative Modeling**
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 **Applications to Text**
 - References
- 6 **Statistical Learning: PAC-Bayesian Approach and Complexity Theory**
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 6 **References**

- **Idea:** If one knows the law of (\underline{X}, Y) everything is easy!

Bayes formula

- With a slight abuse of notation,

$$\begin{aligned}\mathbb{P}(Y|\underline{X}) &= \frac{\mathbb{P}((\underline{X}, Y))}{\mathbb{P}(\underline{X})} \\ &= \frac{\mathbb{P}(\underline{X}|Y)\mathbb{P}(Y)}{\mathbb{P}(\underline{X})}\end{aligned}$$

- **Generative Modeling:**

- Propose a model for (\underline{X}, Y) (or equivalently $\underline{X}|Y$ and Y),
- Estimate it as a density estimation problem,
- Plug the estimate in the Bayes formula
- Plug the conditional estimate in the Bayes *classifier*.
- **Rk:** Require to estimate (\underline{X}, Y) rather than only $Y|\underline{X}$!
- Great flexibility in the model design but may lead to complex computation.

- Simpler setting in classification!

Bayes formula

$$\mathbb{P}(Y = k|\underline{X}) = \frac{\mathbb{P}(\underline{X}|Y = k)\mathbb{P}(Y = k)}{\mathbb{P}(\underline{X})}$$

- Binary Bayes classifier (the best solution)

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = 1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ -1 & \text{otherwise} \end{cases}$$

- **Heuristic:** Estimate those quantities and plug the estimations.
- By using different models/estimators for $\mathbb{P}(\underline{X}|Y)$, we get different classifiers.
- **Rk:** No need to renormalize by $\mathbb{P}(\underline{X})$ to take the decision!

Discriminant Analysis (Gaussian model)

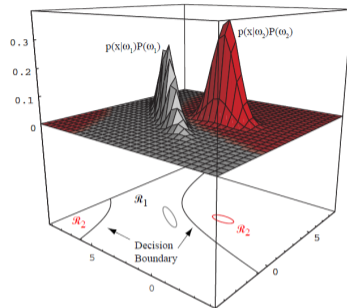
- The densities are modeled as multivariate normal, i.e.,

$$\mathbb{P}(\underline{X}|Y = k) \sim \mathcal{N}_{\mu_k, \Sigma_k}$$

- Discriminant functions: $\mathbf{g}_k(\underline{X}) = \ln(\mathbb{P}(\underline{X}|Y = k)) + \ln(\mathbb{P}(Y = k))$

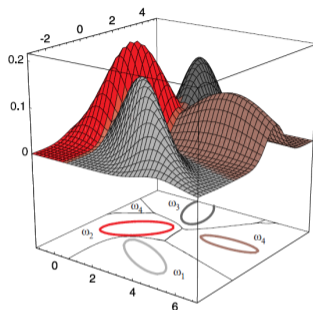
$$\begin{aligned} g_k(\underline{X}) = & -\frac{1}{2}(\underline{X} - \mu_k)^\top \Sigma_k^{-1}(\underline{X} - \mu_k) \\ & - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma_k|) + \ln(\mathbb{P}(Y = k)) \end{aligned}$$

- QDA (different Σ_k in each class) and LDA ($\Sigma_k = \Sigma$ for all k)
- **Beware: this model can be false but the methodology remains valid!**



Quadratic Discriminant Analysis

- The probability densities are Gaussian
- The effect of any decision rule is to divide the feature space into some decision regions $\mathcal{R}_1, \mathcal{R}_2$
- The regions are separated by decision boundaries



Quadratic Discriminant Analysis

- The probability densities are Gaussian
- The effect of any decision rule is to divide the feature space into some decision regions $\mathcal{R}_1, \mathcal{R}_2, \dots, \mathcal{R}_c$
- The regions are separated by decision boundaries

Estimation

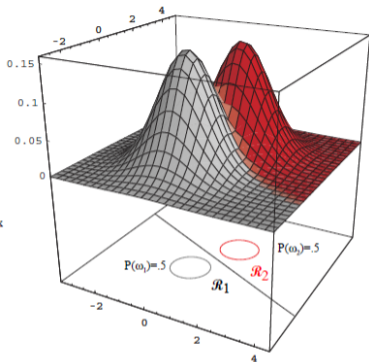
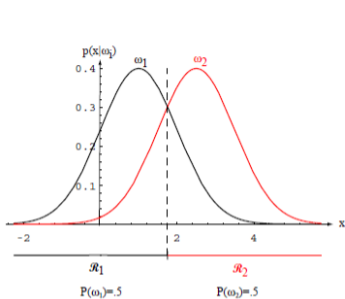
In practice, we will need to estimate μ_k , Σ_k and $\mathbb{P}_k := \mathbb{P}(Y = k)$

- The estimate proportion $\mathbb{P}(\widehat{Y} = k) = \frac{n_k}{n} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{Y_i=k\}}$
- Maximum likelihood estimate of $\widehat{\mu}_k$ and $\widehat{\Sigma}_k$ (explicit formulas)

- DA classifier

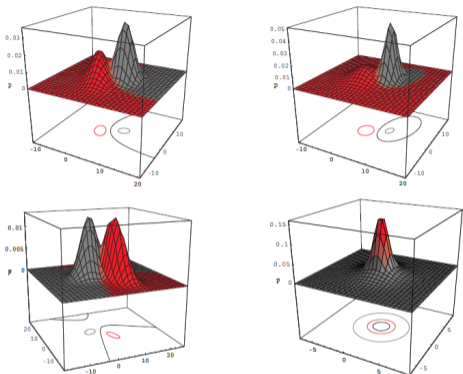
$$\widehat{f}_G(\underline{X}) = \begin{cases} +1 & \text{if } \widehat{g}_{+1}(\underline{X}) \geq \widehat{g}_{-1}(\underline{X}) \\ -1 & \text{otherwise} \end{cases}$$

- Decision boundaries: quadratic = degree 2 polynomials.
- If one imposes $\Sigma_{-1} = \Sigma_1 = \Sigma$ then the decision boundaries is a linear hyperplane.



Linear Discriminant Analysis

- $\Sigma_{\omega_1} = \Sigma_{\omega_2} = \Sigma$
- The decision boundaries are linear hyperplanes

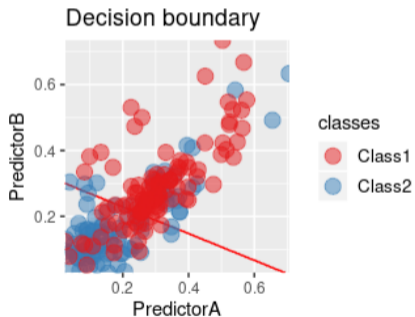
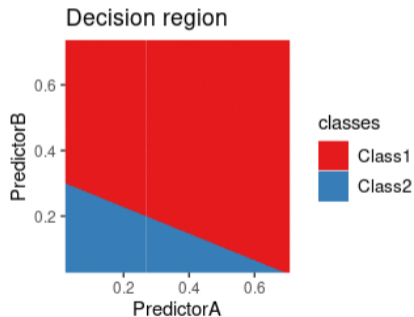


Quadratic Discriminant Analysis

- $\Sigma_{\omega_1} \neq \Sigma_{\omega_2}$
- Arbitrary Gaussian distributions lead to Bayes decision boundaries that are general quadratics.

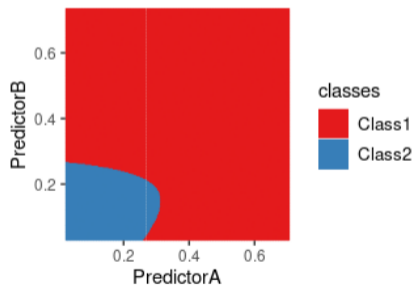
Example: LDA

Linear Discriminant Analysis

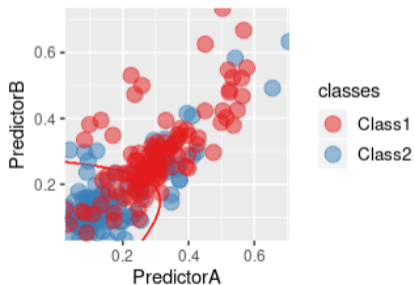


Quadratic Discriminant Analysis

Decision region



Decision boundary



Naive Bayes

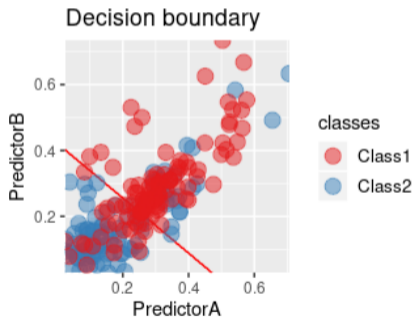
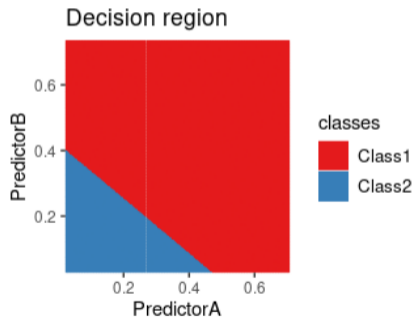
- Classical algorithm using a crude modeling for $\mathbb{P}(\underline{X}|Y)$:
 - Feature **independence** assumption:

$$\mathbb{P}(\underline{X}|Y) = \prod_{l=1}^d \mathbb{P}(X^{(l)}|Y)$$

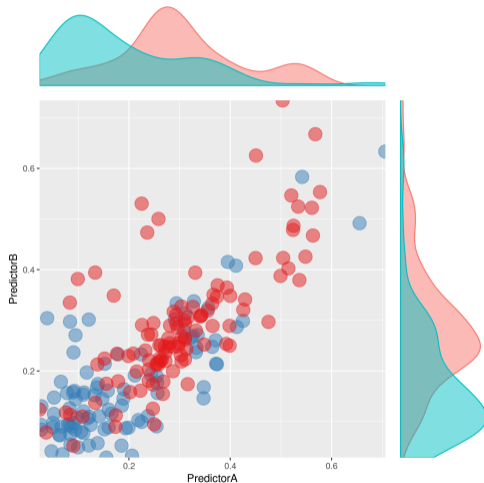
- Simple featurewise model: binomial if binary, multinomial if finite and Gaussian if continuous
- If all features are continuous, similar to the previous Gaussian but with a **diagonal covariance matrix!**
- Very simple learning even in **very high dimension!**

Example: Naive Bayes

Naive Bayes with Gaussian model



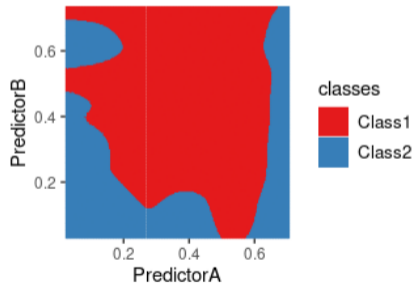
Naive Bayes with Density Estimation



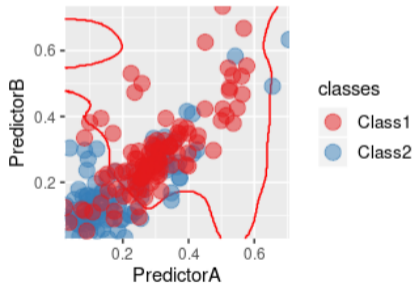
Example: Naive Bayes

Naive Bayes with kernel density estimates

Decision region



Decision boundary



- Other models of the world!

Bayesian Approach

- Generative Model plus prior on the parameters
- Inference thanks to the Bayes formula

Graphical Models

- Markov type models on Graphs

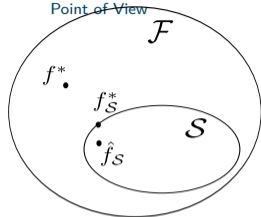
Gaussian Processes

- Multivariate Gaussian models
- ...

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 **ML Methods: Probabilistic Point of View**
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - **Model Selection**
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 **ML Methods: Neural Networks and Deep Learning**
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 **ML Methods: Trees and Ensemble Methods**
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 **Unsupervised Learning: Dimension Reduction and Clustering**
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 6 **Dimension Reduction**
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 6 **Clustering**
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 6 **Generative Modeling**
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 **Applications to Text**
 - References
- 6 **Statistical Learning: PAC-Bayesian Approach and Complexity Theory**
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 6 **References**

Bias-Variance Dilemma

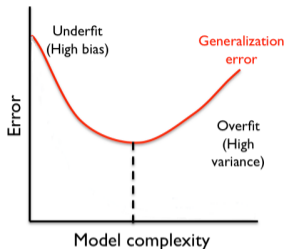
- General setting:
 - $\mathcal{F} = \{\text{measurable functions } \mathcal{X} \rightarrow \mathcal{Y}\}$
 - Best solution: $f^* = \operatorname{argmin}_{f \in \mathcal{F}} \mathcal{R}(f)$
 - Class $\mathcal{S} \subset \mathcal{F}$ of functions
 - Ideal target in \mathcal{S} : $f_S^* = \operatorname{argmin}_{f \in \mathcal{S}} \mathcal{R}(f)$
 - Estimate in \mathcal{S} : \hat{f}_S obtained with some procedure



Approximation error and estimation error (Bias/Variance)

$$\mathcal{R}(\hat{f}_S) - \mathcal{R}(f^*) = \underbrace{\mathcal{R}(f_S^*) - \mathcal{R}(f^*)}_{\text{Approximation error}} + \underbrace{\mathcal{R}(\hat{f}_S) - \mathcal{R}(f_S^*)}_{\text{Estimation error}}$$

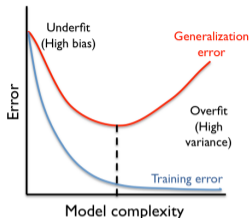
- Approx. error can be large if the model \mathcal{S} is not suitable.
- Estimation error can be large if the model is complex.



- Different behavior for different model complexity
- **Low complexity model** are easily learned but the approximation error (**bias**) may be large (**Under-fit**).
- **High complexity model** may contain a good ideal target but the estimation error (**variance**) can be large (**Over-fit**)

Bias-variance trade-off \iff avoid **overfitting** and **underfitting**

- **Rk:** Better to think in term of method (including feature engineering and specific algorithm) rather than only of model.



Risk behaviour

- Learning/training risk (empirical risk on the learning/training set) decays when the complexity of the **method** increases.
- Quite different behavior when the risk is computed on new observations (generalization risk).
- Overfit for complex methods: parameters learned are too specific to the learning set!
- General situation! (Think of polynomial fit. . .)
- Need to use a different criterion than the training risk!

Predictor Risk Estimation

- **Goal:** Given a predictor f assess its quality.
 - **Method:** Hold-out risk computation (/ Empirical risk correction).
 - **Usage:** Compute an estimate of the risk of a selected f using a **test set** to be used to monitor it in the future.
- Basic block very well understood.

Method Selection

- **Goal:** Given a ML method assess its quality.
 - **Method:** Cross Validation (/ Empirical risk correction)
 - **Usage:** Compute risk estimates for several ML methods using **training/validation sets** to choose the most promising one.
- Estimates can be pointwise or better intervals.
- Multiple test issues in method selection.

Two Approaches

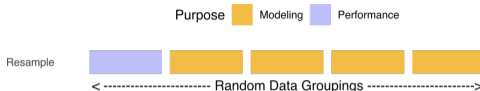
- **Cross validation:** Use empirical risk criterion but on independent data, very efficient (and almost always used in practice!) but slightly biased as its target uses only a fraction of the data.
- **Correction approach:** use empirical risk criterion but *correct* it with a term increasing with the complexity of \mathcal{S}

$$R_n(\hat{f}_S) \rightarrow R_n(\hat{f}_S) + \text{cor}(\mathcal{S})$$

and choose the method with the smallest corrected risk.

Which loss to use?

- The loss used in the risk: most natural!
- The loss used to estimate $\hat{\theta}$: penalized estimation!
- Other performance measure can be used.



- **Very simple idea:** use a second learning/verification set to compute a verification risk.
- Sufficient to remove the dependency issue!
- Implicit random design setting...

Cross Validation

- Use $(1 - \epsilon) \times n$ observations to train and $\epsilon \times n$ to verify!
- Possible issues:
 - Validation for a learning set of size $(1 - \epsilon) \times n$ instead of n ?
 - Unstable risk estimate if ϵn is too small ?
- Most classical variations:
 - Hold Out,
 - Leave One Out,
 - V -fold cross validation.

Principle

- Split the dataset \mathcal{D} in 2 sets $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ of size $n \times (1 - \epsilon)$ and $n \times \epsilon$.
- Learn \hat{f}^{HO} from the subset $\mathcal{D}_{\text{train}}$.
- Compute the empirical risk on the subset $\mathcal{D}_{\text{test}}$:

$$\mathcal{R}_n^{HO}(\hat{f}^{HO}) = \frac{1}{n\epsilon} \sum_{(\underline{X}_i, Y_i) \in \mathcal{D}_{\text{test}}} \ell(Y_i, \hat{f}^{HO}(\underline{X}_i))$$

Predictor Risk Estimation

- Use \hat{f}^{HO} as predictor.
- Use $\mathcal{R}_n^{HO}(\hat{f}^{HO})$ as an estimate of the risk of this estimator.

Method Selection by Cross Validation

- Compute $\mathcal{R}_n^{HO}(\hat{f}_S^{HO})$ for all the considered methods,
- Select the method with the smallest CV risk,
- Reestimate the \hat{f}_S with all the data.

Principle

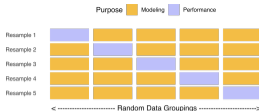
- Split the dataset \mathcal{D} in 2 sets $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ of size $n \times (1 - \epsilon)$ and $n \times \epsilon$.
- Learn \hat{f}^{HO} from the subset $\mathcal{D}_{\text{train}}$.
- Compute the empirical risk on the subset $\mathcal{D}_{\text{test}}$:

$$\mathcal{R}_n^{HO}(\hat{f}^{HO}) = \frac{1}{n\epsilon} \sum_{(\underline{X}_i, Y_i) \in \mathcal{D}_{\text{test}}} \ell(Y_i, \hat{f}^{HO}(\underline{X}_i))$$

- Only possible setting for risk estimation.

Hold Out Limitation for Method Selection

- Biased toward simpler method as the estimation does not use all the data initially.
- Learning variability of $\mathcal{R}_n^{HO}(\hat{f}^{HO})$ not taken into account.



Principle

- Split the dataset \mathcal{D} in V sets \mathcal{D}_v of almost equals size.
- For $v \in \{1, \dots, V\}$:
 - Learn \hat{f}^{-v} from the dataset \mathcal{D} minus the set \mathcal{D}_v .
 - Compute the empirical risk:

$$\mathcal{R}_n^{-v}(\hat{f}^{-v}) = \frac{1}{n_v} \sum_{(\underline{X}_i, Y_i) \in \mathcal{D}_v} \ell(Y_i, \hat{f}^{-v}(\underline{X}_i))$$

- Compute the average empirical risk:

$$\mathcal{R}_n^{CV}(\hat{f}) = \frac{1}{V} \sum_{v=1}^V \mathcal{R}_n^{-v}(\hat{f}^{-v})$$

- Estimation of the quality of a method not of a given predictor.
- Leave One Out : $V = n$.

Analysis (when n is a multiple of V)

- The $\mathcal{R}_n^{-v}(\hat{f}^{-v})$ are identically distributed variable but are not independent!
- Consequence:

$$\begin{aligned}\mathbb{E} \left[\mathcal{R}_n^{CV}(\hat{f}) \right] &= \mathbb{E} \left[\mathcal{R}_n^{-v}(\hat{f}^{-v}) \right] \\ \text{Var} \left[\mathcal{R}_n^{CV}(\hat{f}) \right] &= \frac{1}{V} \text{Var} \left[\mathcal{R}_n^{-v}(\hat{f}^{-v}) \right] \\ &\quad + \left(1 - \frac{1}{V} \right) \text{Cov} \left[\mathcal{R}_n^{-v}(\hat{f}^{-v}), \mathcal{R}_n^{-v'}(\hat{f}^{-v'}) \right]\end{aligned}$$

- Average risk for a sample of size $(1 - \frac{1}{V})n$.
 - Variance term much more complex to analyze!
 - Fine analysis shows that the larger V the better...
-
- Accuracy/Speed tradeoff: $V = 5$ or $V = 10$...

- Leave One Out = V fold for $V = n$: very expensive in general.

A fast LOO formula for the linear regression

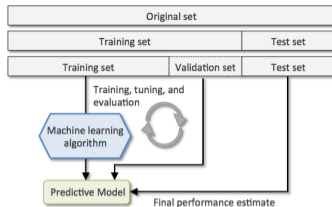
- **Prop:** for the least squares linear regression,

$$\hat{f}^{-i}(\underline{X}_i) = \frac{\hat{f}(\underline{X}_i) - h_{ii} Y_i}{1 - h_{ii}}$$

with h_{ii} the i th diagonal coefficient of the **hat** (projection) matrix.

- Proof based on linear algebra!
- Leads to a fast formula for LOO:

$$\mathcal{R}_n^{LOO}(\hat{f}) = \frac{1}{n} \sum_{i=1}^n \frac{|Y_i - \hat{f}(\underline{X}_i)|^2}{(1 - h_{ii})^2}$$



- **Selection Bias Issue:**
 - After method selection, the cross validation is biased.
 - Furthermore, it qualifies the method and not the final predictor.
- Need to (re)estimate the risk of the final predictor.

(Train/Validation)/Test strategy

- **Split** the dataset in two a (Train/Validation) and Test.
 - Use **CV** with the (Train/Validation) to **select a method**.
 - Train this method on (Train/Validation) to **obtain a single predictor**.
 - Estimate the **performance of this predictor** on Test.
-
- Every choice made from the data is part of the method!

- Empirical loss of an estimator computed on the dataset used to choose it is biased!
- Empirical loss is an optimistic estimate of the true loss.

Risk Correction Heuristic

- Estimate an upper bound of this optimism for a given family.
- Correct the empirical loss by adding this upper bound.
- **Rk:** Finding such an upper bound can be complicated!
- Correction often called a **penalty**.

Penalized Loss

- Minimization of

$$\operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_{\theta}(\underline{X}_i)) + \operatorname{pen}(\theta)$$

where $\operatorname{pen}(\theta)$ is a risk correction (penalty).

Penalties

- Upper bound of the optimism of the empirical loss
- Depends on the loss and the framework!

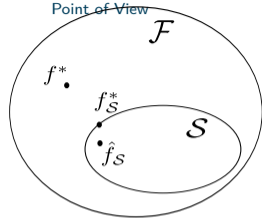
Instantiation

- Mallows Cp: Least Squares with $\operatorname{pen}(\theta) = 2\frac{d}{n}\sigma^2$.
- AIC Heuristics: Maximum Likelihood with $\operatorname{pen}(\theta) = \frac{d}{n}$.
- BIC Heuristics: Maximum Likelihood with $\operatorname{pen}(\theta) = \log(n)\frac{d}{n}$.
- Structural Risk Minimization: Pred. loss and clever penalty.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 **ML Methods: Probabilistic Point of View**
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - **Penalization**
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- Gradient Descent Step
- Non-Convex Setting
- References
- 5 **ML Methods: Neural Networks and Deep Learning**
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 **ML Methods: Trees and Ensemble Methods**
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 **Unsupervised Learning: Dimension Reduction and Clustering**
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- Applications to Text
- References
- 8 **Statistical Learning: PAC-Bayesian Approach and Complexity Theory**
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 **References**

Bias-Variance Dilemma

- General setting:
 - $\mathcal{F} = \{\text{measurable functions } \mathcal{X} \rightarrow \mathcal{Y}\}$
 - Best solution: $f^* = \operatorname{argmin}_{f \in \mathcal{F}} \mathcal{R}(f)$
 - Class $\mathcal{S} \subset \mathcal{F}$ of functions
 - Ideal target in \mathcal{S} : $f_S^* = \operatorname{argmin}_{f \in \mathcal{S}} \mathcal{R}(f)$
 - Estimate in \mathcal{S} : \hat{f}_S obtained with some procedure

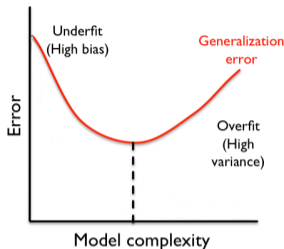


Approximation error and estimation error (Bias/Variance)

$$\mathcal{R}(\hat{f}_S) - \mathcal{R}(f^*) = \underbrace{\mathcal{R}(f_S^*) - \mathcal{R}(f^*)}_{\text{Approximation error}} + \underbrace{\mathcal{R}(\hat{f}_S) - \mathcal{R}(f_S^*)}_{\text{Estimation error}}$$

- Approx. error can be large if the model \mathcal{S} is not suitable.
- Estimation error can be large if the model is complex.

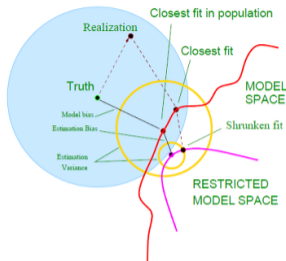
Under-fitting / Over-fitting Issue



- Different behavior for different model complexity
- **Low complexity model** are easily learned but the approximation error (**bias**) may be large (**Under-fit**).
- **High complexity model** may contain a good ideal target but the estimation error (**variance**) can be large (**Over-fit**)

Bias-variance trade-off \iff avoid **overfitting** and **underfitting**

- **Rk:** Better to think in term of method (including feature engineering and specific algorithm) rather than only of model.



Bias-Variance Issue

- Most complex models may not be the best ones due to the variability of the estimate.
- Naive idea: can we *simplify* our model without losing too much?
 - by using only a subset of the variables?
 - by forcing the coefficients to be small?
- Can we do better than exploring all possibilities?

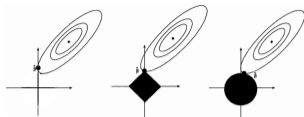
- **Setting:** Gen. linear model = prediction of Y by $h(\underline{x}^\top \beta)$.

Model coefficients

- Model entirely specified by β .
- Coefficientwise:
 - $\beta^{(i)} = 0$ means that the i th covariate is not used.
 - $\beta^{(i)} \sim 0$ means that the i th covariate as a *low* influence. . .
- If some covariates are useless, better use a simpler model. . .

Submodels

- *Simplify* the model through a constraint on β !
- Examples:
 - Support: Impose that $\beta^{(i)} = 0$ for $i \notin I$.
 - Support size: Impose that $\|\beta\|_0 = \sum_{i=1}^d \mathbf{1}_{\beta^{(i)} \neq 0} < C$
 - Norm: Impose that $\|\beta\|_p < C$ with $1 \leq p$ (Often $p = 2$ or $p = 1$)



Sparsity

- β is sparse if its number of non-zero coefficients (l_0) is small. . .
- Easy interpretation in terms of dimension/complexity.

Norm Constraint and Sparsity

- Sparsest solution obtained by definition with the l_0 norm.
- No induced sparsity with the l_2 norm. . .
- Sparsity with the l_1 norm (can even be proved to be the same as with the l_0 norm under some assumptions).
- Geometric explanation.

Constrained Optimization

- Choose a constant C .
- Compute β as

$$\operatorname{argmin}_{\beta \in \mathbb{R}^d, \|\beta\|_p \leq C} \frac{1}{n} \sum_{i=1}^n \bar{\ell}(Y_i, h(\underline{x}_i^\top \beta))$$

Lagrangian Reformulation

- Choose λ and compute β as

$$\operatorname{argmin}_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \bar{\ell}(Y_i, h(\underline{x}_i^\top \beta)) + \lambda \|\beta\|_{p'}$$

with $p' = p$ except if $p = 0$ where $p' = 1$.

- Easier calibration... but no explicit model \mathcal{S} .
- **Rk:** $\|\beta\|_p$ is not scaling invariant if $p \neq 0$...
- Initial rescaling issue.

Penalized Linear Model

- Minimization of

$$\operatorname{argmin}_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \bar{\ell}(Y_i, h(\underline{x}_i^\top \beta)) + \operatorname{pen}(\beta)$$

where $\operatorname{pen}(\beta)$ is a (sparsity promoting) penalty

- Variable selection if β is sparse.

Classical Penalties

- AIC: $\operatorname{pen}(\beta) = \lambda \|\beta\|_0$ (non-convex / sparsity)
- Ridge: $\operatorname{pen}(\beta) = \lambda \|\beta\|_2^2$ (convex / no sparsity)
- Lasso: $\operatorname{pen}(\beta) = \lambda \|\beta\|_1$ (convex / sparsity)
- Elastic net: $\operatorname{pen}(\beta) = \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2$ (convex / sparsity)

- Easy optimization if pen (and the loss) is convex...
- **Need to specify λ to define a ML method!**

Practical Selection Methodology

- Choose a penalty family pen_λ .
 - Compute a CV risk for the penalty pen_λ for all $\lambda \in \Lambda$.
 - Determine $\hat{\lambda}$ the λ minimizing the CV risk.
 - Compute the final model with the penalty $\text{pen}_{\hat{\lambda}}$.
- CV allows to select a ML method, penalized estimation with a penalty $\text{pen}_{\hat{\lambda}}$, not a single predictor hence the need of a final reestimation.

Why not using CV on a grid?

- Grid size scales exponentially with the dimension!
- **If the penalized minimization is easy**, much cheaper to compute the CV risk for all $\lambda \in \Lambda$...
- CV performs best when the set of candidates is not too big (or is structured...)



T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning.
Springer Series in Statistics, 2009



M. Mohri, A. Rostamizadeh, and A. Talwalkar.
Foundations of Machine Learning.
MIT Press, 2012



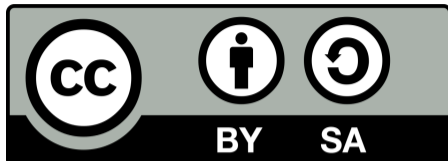
A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (2nd ed.)
O'Reilly, 2019



Ch. Giraud.
Introduction to High-Dimensional Statistics.
CRC Press, 2014



S. Bubeck.
Convex Optimization: Algorithms and Complexity.
Now Publisher, 2015



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 6 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 6 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 6 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - **Supervised Learning**
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 6 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 6 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 6 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Supervised Learning Framework

- Input measurement $\underline{X} \in \mathcal{X}$
- Output measurement $Y \in \mathcal{Y}$.
- $(\underline{X}, Y) \sim \mathbb{P}$ with \mathbb{P} unknown.
- **Training data** : $\mathcal{D}_n = \{(\underline{X}_1, Y_1), \dots, (\underline{X}_n, Y_n)\}$ (i.i.d. $\sim \mathbb{P}$)
- Often
 - $\underline{X} \in \mathbb{R}^d$ and $Y \in \{-1, 1\}$ (classification)
 - or $\underline{X} \in \mathbb{R}^d$ and $Y \in \mathbb{R}$ (regression).
- A **predictor** is a function in $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \text{ meas.}\}$

Goal

- Construct a **good** predictor \hat{f} from the training data.
- Need to specify the meaning of good.
- Classification and regression are almost the **same** problem!

Loss function for a generic predictor

- **Loss function:** $\ell(Y, f(\underline{X}))$ measures the goodness of the prediction of Y by $f(\underline{X})$
- Examples:
 - 0/1 loss: $\ell(Y, f(\underline{X})) = \mathbf{1}_{Y \neq f(\underline{X})}$
 - Quadratic loss: $\ell(Y, f(\underline{X})) = |Y - f(\underline{X})|^2$

Risk function

- Risk measured as the average loss for a new couple:

$$\mathcal{R}(f) = \mathbb{E}_{(X, Y) \sim \mathbb{P}}[\ell(Y, f(\underline{X}))]$$

- Examples:
 - 0/1 loss: $\mathbb{E}[\ell(Y, f(\underline{X}))] = \mathbb{P}(Y \neq f(\underline{X}))$
 - Quadratic loss: $\mathbb{E}[\ell(Y, f(\underline{X}))] = \mathbb{E}[|Y - f(\underline{X})|^2]$

- **Beware:** As \hat{f} depends on \mathcal{D}_n , $\mathcal{R}(\hat{f})$ is a random variable!

- The best solution f^* (which is independent of \mathcal{D}_n) is

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{R}(f) = \arg \min_{f \in \mathcal{F}} \mathbb{E}[\ell(Y, f(\underline{X}))] = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\underline{X}} \left[\mathbb{E}_{Y|\underline{X}}[\ell(Y, f(\underline{X}))] \right]$$

Bayes Predictor (explicit solution)

- In binary classification with 0 – 1 loss:

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = +1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ & \Leftrightarrow \mathbb{P}(Y = +1|\underline{X}) \geq 1/2 \\ -1 & \text{otherwise} \end{cases}$$

- In regression with the quadratic loss

$$f^*(\underline{X}) = \mathbb{E}[Y|\underline{X}]$$

Issue: Solution requires to **know** $\mathbb{E}[Y|\underline{X}]$ for all values of \underline{X} !

Machine Learning

- Learn a rule to construct a **predictor** $\hat{f} \in \mathcal{F}$ from the training data \mathcal{D}_n s.t. **the risk** $\mathcal{R}(\hat{f})$ is **small on average** or with high probability with respect to \mathcal{D}_n .
- In practice, the rule should be an algorithm!

Canonical example: Empirical Risk Minimizer

- One restricts f to a subset of functions $\mathcal{S} = \{f_\theta, \theta \in \Theta\}$
- One replaces the minimization of the average loss by the minimization of the empirical loss

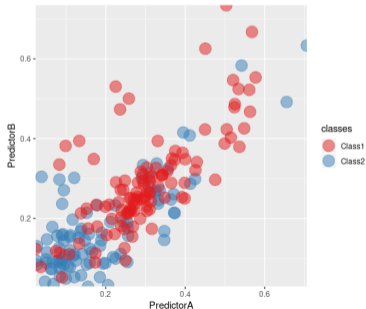
$$\hat{f} = f_{\hat{\theta}} = \operatorname{argmin}_{f_\theta, \theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_\theta(\underline{X}_i))$$

- Examples:
 - Linear regression
 - Linear classification with

$$\mathcal{S} = \{\underline{x} \mapsto \operatorname{sign}\{\underline{x}^\top \beta + \beta^{(0)}\} / \beta \in \mathbb{R}^d, \beta^{(0)} \in \mathbb{R}\}$$

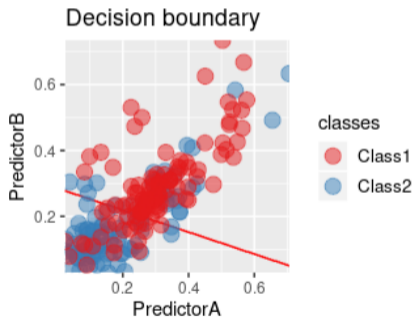
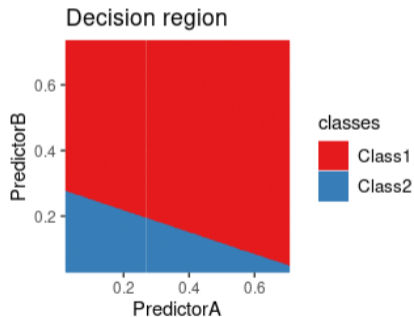
Synthetic Dataset

- Two features/covariates.
- Two classes.
- Dataset from *Applied Predictive Modeling*, M. Kuhn and K. Johnson, Springer
- Numerical experiments with R and the `{caret}` package.



Example: Linear Discrimination

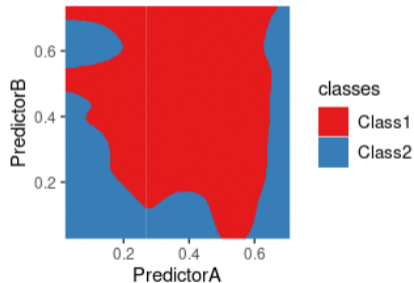
Logistic



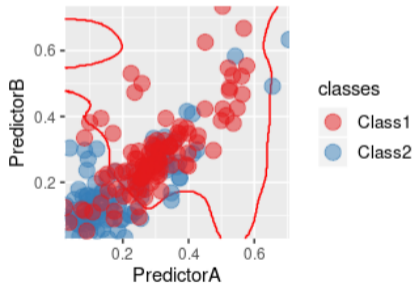
Example: More Complex Model

Naive Bayes with kernel density estimates

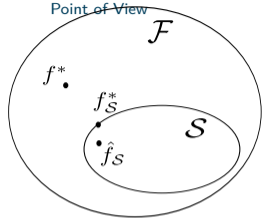
Decision region



Decision boundary



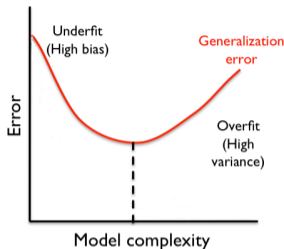
- General setting:
 - $\mathcal{F} = \{\text{measurable functions } \mathcal{X} \rightarrow \mathcal{Y}\}$
 - Best solution: $f^* = \operatorname{argmin}_{f \in \mathcal{F}} \mathcal{R}(f)$
 - Class $\mathcal{S} \subset \mathcal{F}$ of functions
 - Ideal target in \mathcal{S} : $f_S^* = \operatorname{argmin}_{f \in \mathcal{S}} \mathcal{R}(f)$
 - Estimate in \mathcal{S} : \hat{f}_S obtained with some procedure



Approximation error and estimation error (Bias/Variance)

$$\mathcal{R}(\hat{f}_S) - \mathcal{R}(f^*) = \underbrace{\mathcal{R}(f_S^*) - \mathcal{R}(f^*)}_{\text{Approximation error}} + \underbrace{\mathcal{R}(\hat{f}_S) - \mathcal{R}(f_S^*)}_{\text{Estimation error}}$$

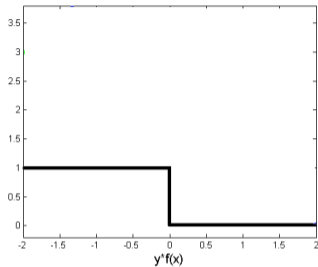
- Approx. error can be large if the model \mathcal{S} is not suitable.
- Estimation error can be large if the model is complex.



- Different behavior for different model complexity
- **Low complexity model** are easily learned but the approximation error (**bias**) may be large (**Under-fit**).
- **High complexity model** may contain a good ideal target but the estimation error (**variance**) can be large (**Over-fit**)

Bias-variance trade-off \iff avoid **overfitting** and **underfitting**

- **Rk:** Better to think in term of method (including feature engineering and specific algorithm) rather than only of model.



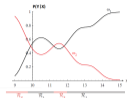
Empirical Risk Minimizer

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f(\underline{X}_i))$$

- Classification loss: $\ell^{0/1}(y, f(\underline{x})) = \mathbf{1}_{y \neq f(\underline{x})}$
- Not convex and not smooth!

Probabilistic Point of View

Ideal Solution and Estimation



- The best solution f^* (which is independent of \mathcal{D}_n) is

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{R}(f) = \arg \min_{f \in \mathcal{F}} \mathbb{E}[\ell(Y, f(\underline{X}))] = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\underline{X}} \left[\mathbb{E}_{Y|\underline{X}}[\ell(Y, f(\underline{x}))] \right]$$

Bayes Predictor (explicit solution)

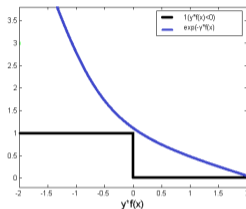
In binary classification with 0 – 1 loss:

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = +1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ -1 & \text{otherwise} \end{cases}$$

- **Issue:** Solution requires to **know** $\mathbb{E}[Y|\underline{X}]$ for all values of \underline{X} !
- **Solution:** Replace it by an estimate.

Optimization Point of View

Loss Convexification



Minimizer of the risk

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f(\underline{X}_i))$$

- **Issue:** Classification loss is not convex or smooth.
- **Solution:** Replace it by a convex majorant.

Probabilistic and Optimization Framework

How to find a good function f with a *small* risk

$$\mathcal{R}(f) = \mathbb{E}[\ell(Y, f(\underline{X}))] \quad ?$$

Canonical approach: $\hat{f}_S = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(\underline{X}_i))$

Problems

- How to choose \mathcal{S} ?
- How to compute the minimization?

A Probabilistic Point of View

Solution: For \underline{X} , estimate $Y|\underline{X}$ plug this estimate in the Bayes classifier:
(Generalized) Linear Models, Kernel methods, k -nn, Naive Bayes, Tree, Bagging...

An Optimization Point of View

Solution: If necessary replace the loss ℓ by an upper bound $\bar{\ell}$ and minimize the empirical loss: **SVR, SVM, Neural Network, Tree, Boosting...**

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - **Optimization Point of View**
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- Gradient Descent Step
- Non-Convex Setting
- References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- Applications to Text
- References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Probabilistic and Optimization Framework

How to find a good function f with a *small* risk

$$\mathcal{R}(f) = \mathbb{E}[\ell(Y, f(\underline{X}))] \quad ?$$

Canonical approach: $\hat{f}_S = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(\underline{X}_i))$

Problems

- How to choose \mathcal{S} ?
- How to compute the minimization?

A Probabilistic Point of View

Solution: For \underline{X} , estimate $Y|\underline{X}$ plug this estimate in the Bayes classifier:
(Generalized) Linear Models, Kernel methods, k -nn, Naive Bayes, Tree, Bagging...

An Optimization Point of View

Solution: If necessary replace the loss ℓ by an upper bound $\bar{\ell}$ and minimize the empirical loss: **SVR, SVM, Neural Network, Tree, Boosting...**

Penalized Logistic Regression

- Let $f_{\theta}(\underline{X}) = \underline{X}^{\top} \beta + \beta^{(0)}$ with $\theta = (\beta, \beta^{(0)})$.
- Find $\hat{\theta} = \arg \min \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-Y_i f_{\theta}(\underline{X}_i)} \right) + \lambda \|\beta\|_1$
- Classify using $\text{sign}(f_{\hat{\theta}})$

Deep Learning

- Let $f_{\theta}(\underline{X})$ with f a feed forward neural network outputting two values with a softmax layer as a last layer.
- Optimize by gradient descent the cross-entropy $-\frac{1}{n} \sum_{i=1}^n \log \left(f_{\theta}(\underline{X}_i)^{(Y_i)} \right)$
- Classify using $\text{sign}(f_{\hat{\theta}})$

Support Vector Machine

- Let $f_{\theta}(\underline{X}) = \underline{X}^{\top} \beta + \beta^{(0)}$ with $\theta = (\beta, \beta^{(0)})$.
 - Find $\hat{\theta} = \arg \min \frac{1}{n} \sum_{i=1}^n \max(1 - Y_i f_{\theta}(\underline{X}_i), 0) + \lambda \|\beta\|_2^2$
 - Classify using $\text{sign}(f_{\hat{\theta}})$
- Those three methods rely on a similar heuristic: the optimization point of view!

- The best solution f^* is the one minimizing

$$f^* = \arg \min R(f) = \arg \min \mathbb{E}[\ell(Y, f(\underline{X}))]$$

Empirical Risk Minimization

- One restricts f to a subset of functions $\mathcal{S} = \{f_\theta, \theta \in \Theta\}$
- One replaces the minimization of the average loss by the minimization of the average empirical loss

$$\hat{f} = f_{\hat{\theta}} = \operatorname{argmin}_{f_\theta, \theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_\theta(\underline{X}_i))$$

- Intractable for the $\ell^{0/1}$ loss!

Risk Convexification

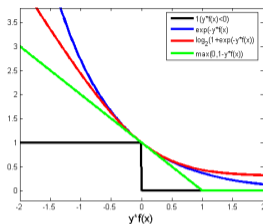
- Replace the loss $\ell(Y, f_\theta(\underline{X}))$ by a convex upperbound $\bar{\ell}(Y, f_\theta(\underline{X}))$ (surrogate loss).
- Minimize the average of the surrogate empirical loss

$$\tilde{f} = f_{\hat{\theta}} = \operatorname{argmin}_{f_\theta, \theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \bar{\ell}(Y_i, f_\theta(\underline{X}_i))$$

- Use $\hat{f} = \operatorname{sign}(\tilde{f})$
- Much easier optimization.

Instantiation

- Logistic (Revisited)
- Support Vector Machine
- (Deep) Neural Network
- Boosting



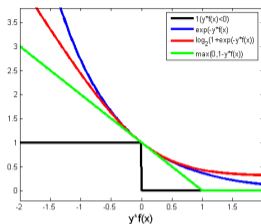
Convexification

- Replace the loss $\ell^{0/1}(Y, f(\underline{X}))$ by

$$\bar{\ell}(Y, f(\underline{X})) = l(Yf(\underline{X}))$$

with l a convex function.

- **Further mild assumption:** l is decreasing, differentiable at 0 and $l'(0) < 0$.



Classical convexification

- Logistic loss: $\bar{\ell}(Y, f(\underline{X})) = \log_2(1 + e^{-Yf(\underline{X})})$ (Logistic / NN)
- Hinge loss: $\bar{\ell}(Y, f(\underline{X})) = (1 - Yf(\underline{X}))_+$ (SVM)
- Exponential loss: $\bar{\ell}(Y, f(\underline{X})) = e^{-Yf(\underline{X})}$ (Boosting...)

The Target is the Bayes Classifier

- The minimizer of

$$\mathbb{E}[\bar{\ell}(Y, f(\underline{X}))] = \mathbb{E}[I(Yf(\underline{X}))]$$

is the Bayes classifier $f^* = \text{sign}(2\eta(\underline{X}) - 1)$

Control of the Excess Risk

- It exists a convex function Ψ such that

$$\begin{aligned} \Psi \left(\mathbb{E}[\ell^{0/1}(Y, \text{sign}(f(\underline{X})))] - \mathbb{E}[\ell^{0/1}(Y, f^*(\underline{X}))] \right) \\ \leq \mathbb{E}[\bar{\ell}(Y, f(\underline{X}))] - \mathbb{E}[\bar{\ell}(Y, f^*(\underline{X}))] \end{aligned}$$

- Theoretical guarantee!

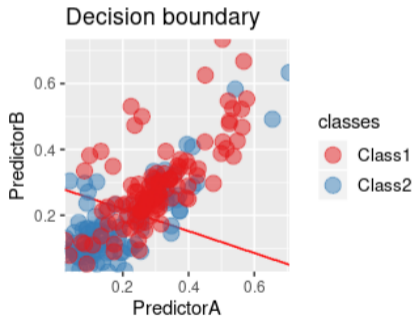
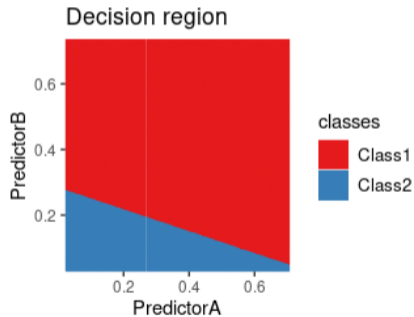
- Ideal solution:

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f(\underline{X}_i))$$

Logistic regression

- Use $f(\underline{X}) = \underline{X}^\top \beta + \beta^{(0)}$.
- Use the logistic loss $\bar{\ell}(y, f) = \log_2(1 + e^{-yf})$, i.e. the negative log-likelihood.
- Different vision than the statistician but same algorithm!

Logistic



- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 **ML Methods: Optimization Point of View**
 - Supervised Learning
 - Optimization Point of View
 - **SVM**
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- Gradient Descent Step
- Non-Convex Setting
- References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- Applications to Text
- References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

$$f_{\theta}(\underline{X}) = \underline{X}^{\top} \beta + \beta^{(0)} \quad \text{with} \quad \theta = (\beta, \beta^{(0)})$$

$$\hat{\theta} = \arg \min \frac{1}{n} \sum_{i=1}^n \max(1 - Y_i f_{\theta}(\underline{X}_i), 0) + \lambda \|\beta\|_2^2$$

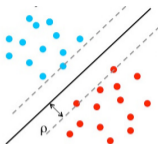
Support Vector Machine

- Convexification of the 0/1-loss with the hinge loss:

$$\mathbf{1}_{Y_i f_{\theta}(\underline{X}_i) < 0} \leq \max(1 - Y_i f_{\theta}(\underline{X}_i), 0)$$

- Penalization by the quadratic norm (Ridge/Tikhonov).
- Solution can be approximated by gradient descent algorithms.

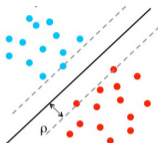
- **Revisit** of the original point of view.
- Original point of view leads to a different optimization algorithm and to some extensions.



- Linear classifier: $\text{sign}(\underline{X}^\top \beta + \beta^{(0)})$
- Separable case: $\exists(\beta, \beta^{(0)}), \forall i, Y_i(\underline{X}_i^\top \beta + \beta^{(0)}) > 0$

How to choose $(\beta, \beta^{(0)})$ so that the separation is maximal?

- Strict separation: $\exists(\beta, \beta^{(0)}), \forall i, Y_i(\underline{X}_i^\top \beta + \beta^{(0)}) \geq 1$
- Distance between $\underline{X}^\top \beta + \beta^{(0)} = 1$ and $\underline{X}^\top \beta + \beta^{(0)} = -1$:
$$\frac{2}{\|\beta\|}$$
- Maximizing this distance is equivalent to minimizing $\frac{1}{2}\|\beta\|^2$.

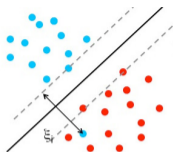


Separable SVM

- Constrained optimization formulation:

$$\min \frac{1}{2} \|\beta\|^2 \quad \text{with} \quad \forall i, Y_i(\underline{X}_i^\top \beta + \beta^{(0)}) \geq 1$$

- Quadratic Programming setting.
- Efficient solver available. . .



- What about the non separable case?

SVM relaxation

- Relax the assumptions

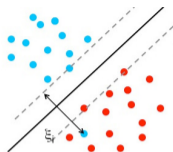
$$\forall i, Y_i(\underline{X}_i^T \beta + \beta^{(0)}) \geq 1 \quad \text{to} \quad \forall i, Y_i(\underline{X}_i^T \beta + \beta^{(0)}) \geq 1 - s_i$$

with the **slack variables** $s_i \geq 0$

- Keep those slack variables as small as possible by minimizing

$$\frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n s_i$$

where $C > 0$ is the **goodness-of-fit strength**



SVM

- Constrained optimization formulation:

$$\min \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n s_i \quad \text{with} \quad \begin{cases} \forall i, Y_i(\underline{X}_i^\top \beta + \beta^{(0)}) \geq 1 - s_i \\ \forall i, s_i \geq 0 \end{cases}$$

- **Hinge Loss** reformulation:

$$\min \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n \underbrace{\max(0, 1 - Y_i(\underline{X}_i^\top \beta + \beta^{(0)}))}_{\text{Hinge Loss}}$$

- Constrained convex optimization algorithms vs gradient descent algorithms.

- Convex relaxation:

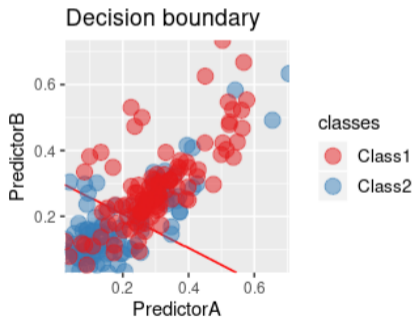
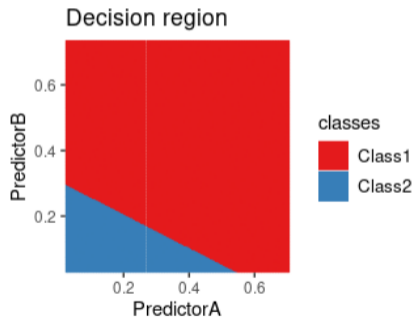
$$\begin{aligned} & \operatorname{argmin} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n \max(1 - Y_i(\underline{X}_i^\top \beta + \beta^{(0)}), 0) \\ & = \operatorname{argmin} \frac{1}{n} \sum_{i=1}^n \max(1 - Y_i(\underline{X}_i^\top \beta + \beta^{(0)}), 0) + \frac{1}{Cn} \frac{1}{2} \|\beta\|^2 \end{aligned}$$

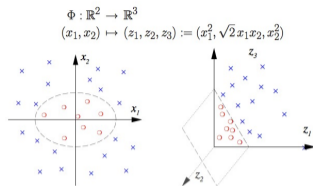
- **Prop:** $\ell^{0/1}(Y_i, \operatorname{sign}(\underline{X}_i^\top \beta + \beta^{(0)})) \leq \max(1 - Y_i(\underline{X}_i^\top \beta + \beta^{(0)}), 0)$

Penalized convex relaxation (Tikhonov!)

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, \operatorname{sign}(\underline{X}_i^\top \beta + \beta^{(0)})) + \frac{1}{Cn} \frac{1}{2} \|\beta\|^2 \\ & \leq \frac{1}{n} \sum_{i=1}^n \max(1 - Y_i(\underline{X}_i^\top \beta + \beta^{(0)}), 0) + \frac{1}{Cn} \frac{1}{2} \|\beta\|^2 \end{aligned}$$

Support Vector Machine





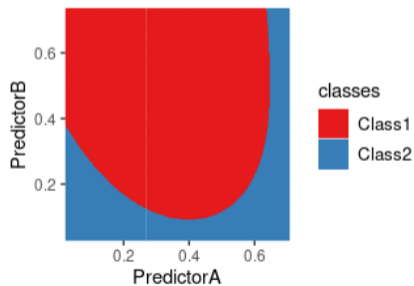
- Non linear separation: just replace \underline{X} by a non linear $\Phi(\underline{X})$...
- Knowing $\phi(\underline{X}_i)^\top \phi(\underline{X}_j)$ is sufficient to compute the SVM solution.

Kernel trick

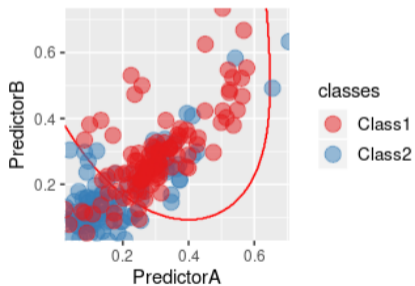
- **Computing $k(\underline{X}, \underline{X}') = \phi(\underline{X})^\top \phi(\underline{X}')$ may be easier than computing $\phi(\underline{X})$, $\phi(\underline{X}')$ and then the scalar product!**
- ϕ can be specified through its definite positive kernel k .
- Examples: Polynomial kernel $k(\underline{X}, \underline{X}') = (1 + \underline{X}^\top \underline{X}')^d$, Gaussian kernel $k(\underline{X}, \underline{X}') = e^{-\|\underline{X} - \underline{X}'\|^2/2}, \dots$
- RKHS setting!
- Can be used in (logistic) regression and more...

Support Vector Machine with polynomial kernel

Decision region

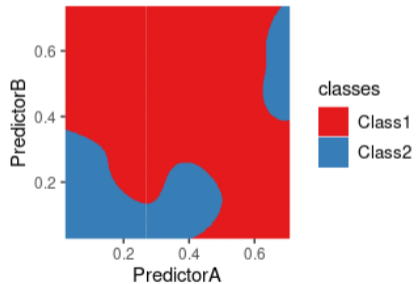


Decision boundary

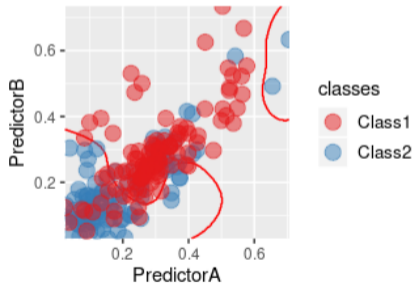


Support Vector Machine with Gaussian kernel

Decision region



Decision boundary



SVM

- Constrained optimization formulation:

$$\min \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n s_i \quad \text{with} \quad \begin{cases} \forall i, Y_i(\underline{X}_i^\top \beta + \beta^{(0)}) \geq 1 - s_i \\ \forall i, s_i \geq 0 \end{cases}$$

SVM Lagrangian

- Lagrangian:

$$\begin{aligned} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu) = & \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n s_i \\ & + \sum_i \alpha_i (1 - s_i - Y_i(\underline{X}_i^\top \beta + \beta^{(0)})) - \sum_i \mu_i s_i \end{aligned}$$

KKT Optimality Conditions

- Stationarity:

$$\nabla_{\beta} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu) = \beta - \sum_i \alpha_i Y_i \underline{X}_i = 0$$

$$\nabla_{\beta^{(0)}} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu) = - \sum_i \alpha_i = 0$$

$$\nabla_{s_i} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu) = C - \alpha_i - \mu_i = 0$$

- Primal and dual admissibility:

$$(1 - s_i - Y_i(\underline{X}_i^{\top} \beta + \beta^{(0)})) \leq 0, \quad s_i \geq 0, \quad \alpha_i \geq 0, \quad \text{and} \quad \mu_i \geq 0$$

- Complementary slackness:

$$\alpha_i(1 - s_i - Y_i(\underline{X}_i^{\top} \beta + \beta^{(0)})) = 0 \quad \text{and} \quad \mu_i s_i = 0$$

Consequence

- $\beta^* = \sum_i \alpha_i Y_i \underline{X}_i$ and $0 \leq \alpha_i \leq C$.
- If $\alpha_i \neq 0$, \underline{X}_i is called a **support vector** and either
 - $s_i = 0$ and $Y_i(\underline{X}_i^{\top} \beta^* + \beta^{(0)*}) = 1$ (margin hyperplane),
 - or $\alpha_i = C$ (outliers).
- $\beta^{(0)*} = Y_i - \underline{X}_i^{\top} \beta^*$ for any support vector with $0 < \alpha_i < C$.

SVM Lagrangian Dual

- Lagrangian Dual:

$$Q(\alpha, \mu) = \min_{\beta, \beta^{(0)}, s} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu)$$

- Prop:

- if $\sum_i \alpha_i Y_i \neq 0$ or $\exists i, \alpha_i + \mu_i \neq C$,

$$Q(\alpha, \mu) = -\infty$$

- if $\sum_i \alpha_i Y_i = 0$ and $\forall i, \alpha_i + \mu_i = C$,

$$Q(\alpha, \mu) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j Y_i Y_j \underline{X}_i^\top \underline{X}_j$$

SVM Dual problem

- Dual problem is a Quadratic Programming problem:

$$\max_{\alpha \geq 0, \mu \geq 0} Q(\alpha, \mu) \Leftrightarrow \max_{0 \leq \alpha \leq C} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j Y_i Y_j \underline{X}_i^\top \underline{X}_j$$

- Involves the \underline{X}_i only through their scalar products.

Mercer Representation Theorem

- For any loss $\bar{\ell}$ and any increasing function Φ , the minimizer in β of

$$\sum_{i=1}^n \bar{\ell}(Y_i, \underline{X}_i^\top \beta + \beta^{(0)}) + \Phi(\|\beta\|_2)$$

is a linear combination of the input points $\beta^* = \sum_{i=1}^n \alpha'_i \underline{X}_i$.

- Minimization problem in α' :

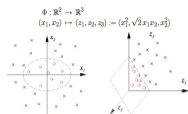
$$\sum_{i=1}^n \bar{\ell}(Y_i, \sum_j \alpha'_j \underline{X}_i^\top \underline{X}_j + \beta^{(0)}) + \Phi(\|\beta\|_2)$$

involving only the scalar product of the data.

- Optimal predictor requires only to compute scalar products.

$$\hat{f}^*(\underline{X}) = \underline{X}^\top \beta^* + \beta^{(0),*} = \sum_i \alpha'_i \underline{X}_i^\top \underline{X}$$

- Transform a problem in dimension $\dim(\mathcal{X})$ in a problem in dimension n .
- Direct minimization in β can be more efficient...

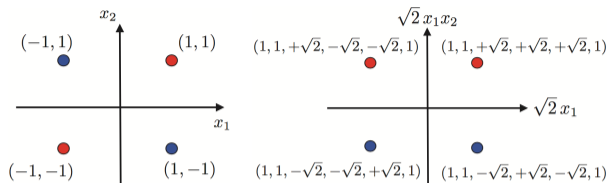


Feature Engineering

- Art of creating **new features** from the existing one \underline{X} .
- Example: add monomials $(\underline{X}^{(j)})^2, \underline{X}^{(j)}\underline{X}^{(j')}\dots$
- Adding feature increases the dimension.

Feature Map

- Application $\phi: \mathcal{X} \rightarrow \mathbb{H}$ with \mathbb{H} an Hilbert space.
- Linear decision boundary in \mathbb{H} : $\phi(\underline{X})^\top \beta + \beta^{(0)} = 0$ is **not an hyperplane anymore** in \mathcal{X} .
- **Heuristic:** Increasing dimension allows to make data almost linearly separable.



Polynomial Mapping of order 2

- $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^6$

$$\phi(\underline{X}) = \left((\underline{X}^{(1)})^2, (\underline{X}^{(2)})^2, \sqrt{2}\underline{X}^{(1)}\underline{X}^{(2)}, \sqrt{2}\underline{X}^{(1)}, \sqrt{2}\underline{X}^{(2)}, 1 \right)$$

- Allow to solve the XOR classification problem with the *hyperplane* $\underline{X}^{(1)}\underline{X}^{(2)} = 0$.

Polynomial Mapping and Scalar Product

- **Prop:**

$$\phi(\underline{X})^\top \phi(\underline{X}') = (1 + \underline{X}^\top \underline{X}')^2$$

Primal, Lagrangian and Dual

- Primal:

$$\min \|\beta\|^2 + C \sum_{i=1}^n s_i \quad \text{with} \quad \begin{cases} \forall i, Y_i(\phi(\underline{X}_i)^\top \beta + \beta^{(0)}) \geq 1 - s_i \\ \forall i, s_i \geq 0 \end{cases}$$

- Lagrangian:

$$\begin{aligned} \mathcal{L}(\beta, \beta^{(0)}, s, \alpha, \mu) &= \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n s_i \\ &\quad + \sum_i \alpha_i (1 - s_i - Y_i(\phi(\underline{X}_i)^\top \beta + \beta^{(0)})) - \sum_i \mu_i s_i \end{aligned}$$

- Dual:

$$\max_{\alpha \geq 0, \mu \geq 0} Q(\alpha, \mu) \Leftrightarrow \max_{0 \leq \alpha \leq C} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j Y_i Y_j \phi(\underline{X}_i)^\top \phi(\underline{X}_j)$$

- Optimal $\phi(\underline{X})^\top \beta^* + \beta^{(0),*} = \sum_i \alpha_i Y_i \phi(\underline{X})^\top \phi(\underline{X}_i)$

- Only need to know to compute $\phi(\underline{X})^\top \phi(\underline{X}')$ to obtain the solution.

- Many algorithms (e.g. SVM) require only to be able to compute the scalar product $\phi(\underline{X})^\top \phi(\underline{X}')$.

Kernel

- Any application

$$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

is called a **kernel** over \mathcal{X} .

Kernel Trick

- Computing directly the **kernel** $k(\underline{X}, \underline{X}') = \phi(\underline{X})^\top \phi(\underline{X}')$ may be easier than computing $\phi(\underline{X})$, $\phi(\underline{X}')$ and then the scalar product.
- Here k is defined from ϕ .
- Under some assumption on k , ϕ can be implicitly *defined* from k !

Positive Definite Symmetric Kernels

- A kernel k is PDS if and only if
 - k is symmetric, i.e.

$$k(\underline{X}, \underline{X}') = k(\underline{X}', \underline{X})$$

- for any $N \in \mathbb{N}$ and any $(\underline{X}_1, \dots, \underline{X}_N) \in \mathcal{X}^N$,

$$\mathbf{K} = [k(\underline{X}_i, \underline{X}_j)]_{1 \leq i, j \leq N}$$

is positive semi-definite, i.e. $\forall u \in \mathbb{R}^N$

$$u^T \mathbf{K} u = \sum_{1 \leq i, j \leq N} u^{(i)} u^{(j)} k(\underline{X}_i, \underline{X}_j) \geq 0$$

or equivalently all the eigenvalues of \mathbf{K} are non-negative.

- The matrix \mathbf{K} is called the **Gram matrix** associated to $(\underline{X}_1, \dots, \underline{X}_N)$.

Moore-Aronsajn Theorem

- For any PDS kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, it exists a Hilbert space $\mathbb{H} \subset \mathbb{R}^{\mathcal{X}}$ with a scalar product $\langle \cdot, \cdot \rangle_{\mathbb{H}}$ such that
 - it exists a mapping $\phi : \mathcal{X} \rightarrow \mathbb{H}$ satisfying
$$k(\underline{X}, \underline{X}') = \langle \phi(\underline{X}), \phi(\underline{X}') \rangle_{\mathbb{H}}$$
 - the **reproducing property** holds, i.e. for any $h \in \mathbb{H}$ and any $\underline{X} \in \mathcal{X}$
$$h(\underline{X}) = \langle h, k(\underline{X}, \cdot) \rangle_{\mathbb{H}}.$$
- By def., \mathbb{H} is a **reproducing kernel Hilbert space** (RKHS).
- \mathbb{H} is called the **feature space** associated to k and ϕ the **feature mapping**.
- No unicity in general.
- **Rk:** if $k(\underline{X}, \underline{X}') = \phi'(\underline{X})^{\top} \phi'(\underline{X}')$ with $\phi' : \mathcal{X} \rightarrow \mathbb{R}^p$ then
 - \mathbb{H} can be chosen as $\{\underline{X} \mapsto \phi'(\underline{X})^{\top} \beta, \beta \in \mathbb{R}^p\}$ and $\|\underline{X} \mapsto \phi'(\underline{X})^{\top} \beta\|_{\mathbb{H}}^2 = \|\beta\|_2^2$.
 - $\phi(\underline{X}') : \underline{X} \mapsto \phi'(\underline{X})^{\top} \phi'(\underline{X}')$.

Separable Kernel

- For any function $\Psi : \mathcal{X} \rightarrow \mathbb{R}$, $k(\underline{X}, \underline{X}') = \Psi(\underline{X})\Psi(\underline{X}')$ is PDS.

Kernel Stability

- For any PDS kernels k_1 and k_2 , $k_1 + k_2$ and $k_1 k_2$ are PDS kernels.
- For any sequence of PDS kernels k_n converging pointwise to a kernel k , k is a PDS kernel.
- For any PDS kernel k such that $|k| \leq r$ and any power series $\sum_n a_n z^n$ with $a_n \geq 0$ and a convergence radius larger than r , $\sum_n a_n k^n$ is a PDS kernel.
- For any PDS kernel k , the renormalized kernel $k'(\underline{X}, \underline{X}') = \frac{k(\underline{X}, \underline{X}')}{\sqrt{k(\underline{X}, \underline{X})k(\underline{X}', \underline{X}')}} is a PDS kernel.$
- Cauchy-Schwartz for k PDS: $k(\underline{X}, \underline{X}')^2 \leq k(\underline{X}, \underline{X})k(\underline{X}', \underline{X}')$

PDS Kernels

- Vanilla kernel:

$$k(\underline{X}, \underline{X}') = \underline{X}^\top \underline{X}'$$

- Polynomial kernel:

$$k(\underline{X}, \underline{X}') = (1 + \underline{X}^\top \underline{X}')^k$$

- Gaussian RBF kernel:

$$k(\underline{X}, \underline{X}') = \exp(-\gamma \|\underline{X} - \underline{X}'\|^2)$$

- Tanh kernel:

$$k(\underline{X}, \underline{X}') = \tanh(a \underline{X}^\top \underline{X}' + b)$$

- Most classical is the Gaussian RBF kernel...
- Lots of freedom to construct kernel for non classical data.



Representer Theorem

- Let k be a PDS kernel and \mathbb{H} its corresponding RKHS, for any increasing function Φ and any function $L : \mathbb{R}^n \rightarrow \mathbb{R}$, the optimization problem

$$\operatorname{argmin}_{h \in \mathbb{H}} L(h(\underline{X}_1), \dots, h(\underline{X}_n)) + \Phi(\|h\|)$$

admits only solutions of the form

$$\sum_{i=1}^n \alpha'_i k(\underline{X}_i, \cdot).$$

- Examples:
 - (kernelized) SVM
 - (kernelized) Penalized Logistic Regression (Ridge)
 - (kernelized) Penalized Regression (Ridge)

Primal

- Constrained Optimization:

$$\min_{f \in \mathbb{H}, \beta^{(0)}, s} \|f\|_{\mathbb{H}}^2 + C \sum_{i=1}^n s_i \quad \text{with} \quad \begin{cases} \forall i, Y_i(f(\underline{X}_i) + \beta^{(0)}) \geq 1 - s_i \\ \forall i, s_i \geq 0 \end{cases}$$

- Hinge loss:

$$\min_{f \in \mathbb{H}, \beta^{(0)}} \|f\|_{\mathbb{H}}^2 + C \sum_{i=1}^n \max(0, 1 - Y_i(f(\underline{X}_i) + \beta^{(0)}))$$

- Representer:

$$\min_{\alpha', \beta^{(0)}} \sum_{i,j} \alpha'_i \alpha'_j k(\underline{X}_i, \underline{X}_j) + C \sum_{i=1}^n \max(0, 1 - Y_i(\sum_j \alpha'_j k(\underline{X}_j, \underline{X}_i) + \beta^{(0)}))$$

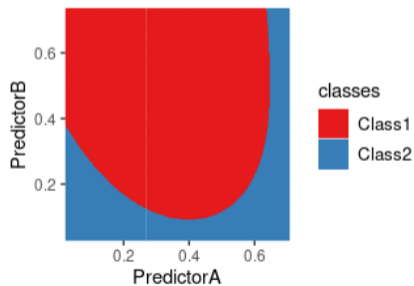
Dual

- Dual:

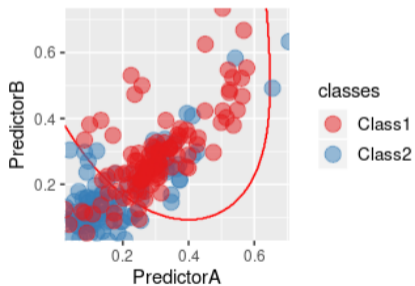
$$\max_{\alpha \geq 0, \mu \geq 0} Q(\alpha, \mu) \Leftrightarrow \max_{0 \leq \alpha \leq C} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j Y_i Y_j k(\underline{X}_i, \underline{X}_j)$$

Support Vector Machine with polynomial kernel

Decision region

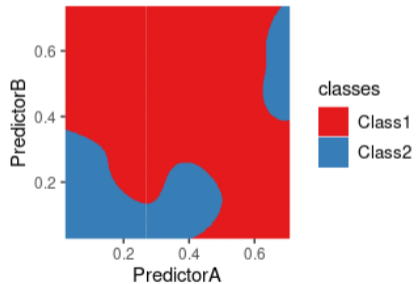


Decision boundary

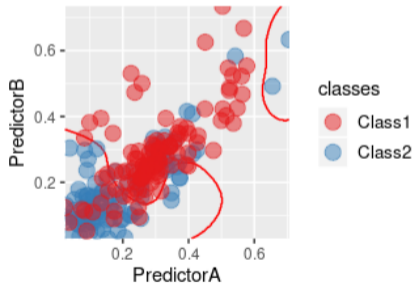


Support Vector Machine with Gaussian kernel

Decision region

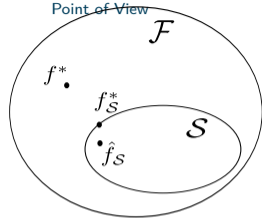


Decision boundary



- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 8 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 9 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 10 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 11 Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- General setting:
 - $\mathcal{F} = \{\text{measurable functions } \mathcal{X} \rightarrow \mathcal{Y}\}$
 - Best solution: $f^* = \operatorname{argmin}_{f \in \mathcal{F}} \mathcal{R}(f)$
 - Class $\mathcal{S} \subset \mathcal{F}$ of functions
 - Ideal target in \mathcal{S} : $f_S^* = \operatorname{argmin}_{f \in \mathcal{S}} \mathcal{R}(f)$
 - Estimate in \mathcal{S} : \hat{f}_S obtained with some procedure

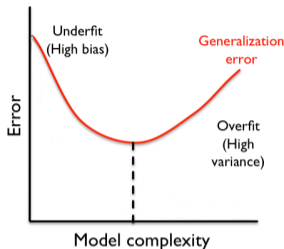


Approximation error and estimation error (Bias/Variance)

$$\mathcal{R}(\hat{f}_S) - \mathcal{R}(f^*) = \underbrace{\mathcal{R}(f_S^*) - \mathcal{R}(f^*)}_{\text{Approximation error}} + \underbrace{\mathcal{R}(\hat{f}_S) - \mathcal{R}(f_S^*)}_{\text{Estimation error}}$$

- Approx. error can be large if the model \mathcal{S} is not suitable.
- Estimation error can be large if the model is complex.

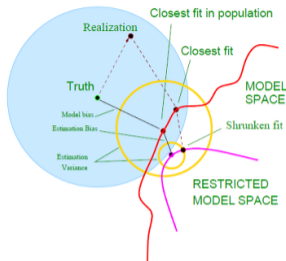
Under-fitting / Over-fitting Issue



- Different behavior for different model complexity
- **Low complexity model** are easily learned but the approximation error (**bias**) may be large (**Under-fit**).
- **High complexity model** may contain a good ideal target but the estimation error (**variance**) can be large (**Over-fit**)

Bias-variance trade-off \iff avoid **overfitting** and **underfitting**

- **Rk:** Better to think in term of method (including feature engineering and specific algorithm) rather than only of model.



Bias-Variance Issue

- Most complex models may not be the best ones due to the variability of the estimate.
- Naive idea: can we *simplify* our model without losing too much?
 - by using only a subset of the variables?
 - by forcing the coefficients to be small?
- Can we do better than exploring all possibilities?

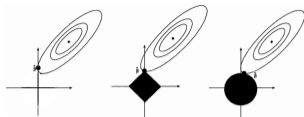
- **Setting:** Gen. linear model = prediction of Y by $h(\underline{x}^T \beta)$.

Model coefficients

- Model entirely specified by β .
- Coefficientwise:
 - $\beta^{(i)} = 0$ means that the i th covariate is not used.
 - $\beta^{(i)} \sim 0$ means that the i th covariate as a *low* influence. . .
- If some covariates are useless, better use a simpler model. . .

Submodels

- *Simplify* the model through a constraint on β !
- Examples:
 - Support: Impose that $\beta^{(i)} = 0$ for $i \notin I$.
 - Support size: Impose that $\|\beta\|_0 = \sum_{i=1}^d \mathbf{1}_{\beta^{(i)} \neq 0} < C$
 - Norm: Impose that $\|\beta\|_p < C$ with $1 \leq p$ (Often $p = 2$ or $p = 1$)



Sparsity

- β is sparse if its number of non-zero coefficients (l_0) is small. . .
- Easy interpretation in terms of dimension/complexity.

Norm Constraint and Sparsity

- Sparsest solution obtained by definition with the l_0 norm.
- No induced sparsity with the l_2 norm. . .
- Sparsity with the l_1 norm (can even be proved to be the same as with the l_0 norm under some assumptions).
- Geometric explanation.

Constrained Optimization

- Choose a constant C .
- Compute β as

$$\operatorname{argmin}_{\beta \in \mathbb{R}^d, \|\beta\|_p \leq C} \frac{1}{n} \sum_{i=1}^n \bar{\ell}(Y_i, h(\underline{x}_i^\top \beta))$$

Lagrangian Reformulation

- Choose λ and compute β as

$$\operatorname{argmin}_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \bar{\ell}(Y_i, h(\underline{x}_i^\top \beta)) + \lambda \|\beta\|_{p'}$$

with $p' = p$ except if $p = 0$ where $p' = 1$.

- Easier calibration... but no explicit model \mathcal{S} .
- **Rk:** $\|\beta\|_p$ is not scaling invariant if $p \neq 0$...
- Initial rescaling issue.

Penalized Linear Model

- Minimization of

$$\operatorname{argmin}_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \bar{\ell}(Y_i, h(\underline{x}_i^\top \beta)) + \operatorname{pen}(\beta)$$

where $\operatorname{pen}(\beta)$ is a (sparsity promoting) penalty

- Variable selection if β is sparse.

Classical Penalties

- AIC: $\operatorname{pen}(\beta) = \lambda \|\beta\|_0$ (non-convex / sparsity)
- Ridge: $\operatorname{pen}(\beta) = \lambda \|\beta\|_2^2$ (convex / no sparsity)
- Lasso: $\operatorname{pen}(\beta) = \lambda \|\beta\|_1$ (convex / sparsity)
- Elastic net: $\operatorname{pen}(\beta) = \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2$ (convex / sparsity)

- Easy optimization if pen (and the loss) is convex...
- **Need to specify λ to define a ML method!**

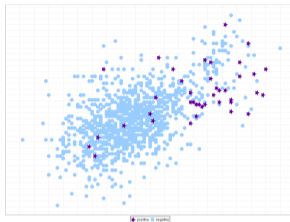
Practical Selection Methodology

- Choose a penalty family pen_λ .
 - Compute a CV risk for the penalty pen_λ for all $\lambda \in \Lambda$.
 - Determine $\hat{\lambda}$ the λ minimizing the CV risk.
 - Compute the final model with the penalty $\text{pen}_{\hat{\lambda}}$.
- CV allows to select a ML method, penalized estimation with a penalty $\text{pen}_{\hat{\lambda}}$, not a single predictor hence the need of a final reestimation.

Why not using CV on a grid?

- Grid size scales exponentially with the dimension!
- **If the penalized minimization is easy**, much cheaper to compute the CV risk for all $\lambda \in \Lambda$...
- CV performs best when the set of candidates is not too big (or is structured...)

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - **Cross Validation and Weights**
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 6 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 6 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 6 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 Applications to Text
 - References
- 6 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 6 References

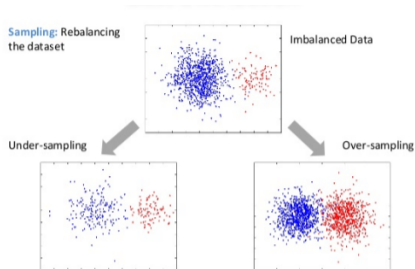


Unbalanced Class

- **Setting:** One of the class is much more present than the other.
- **Issue:** Classifier *too attracted* by the majority class!

Rebalanced Dataset

- **Setting:** Class proportions are different in the training and testing set (stratified sampling)
- **Issue:** Training risks are not estimate of testing risks.



Resampling

- Modify the training dataset so that the classes are more balanced.
- Two flavors:
 - Sub-sampling which spoils data,
 - Over-sampling which needs to create *new* examples.
- **Issues:** Training data is not anymore representative of testing data
- **Hard to do it right!**

Testing

- Testing class prob.: $\pi_t(k)$
- Testing risk target:

$$\begin{aligned}\mathbb{E}_{\pi_t}[\ell(Y, f(\underline{X}))] &= \\ &\sum_k \pi_t(k) \mathbb{E}[\ell(Y, f(\underline{X})) | Y = k]\end{aligned}$$

Training

- Training class prob.: $\pi_{tr}(k)$
- Training risk target:

$$\begin{aligned}\mathbb{E}_{\pi_{tr}}[\ell(Y, f(\underline{X}))] &= \\ &\sum_k \pi_{tr}(k) \mathbb{E}[\ell(Y, f(\underline{X})) | Y = k]\end{aligned}$$

Implicit Testing Risk Using the Training One

- Amounts to use a weighted loss:

$$\begin{aligned}\mathbb{E}_{\pi_{tr}}[\ell(Y, f(\underline{X}))] &= \sum_k \pi_{tr}(k) \mathbb{E}[\ell(Y, f(\underline{X})) | Y = k] \\ &= \sum_k \pi_t(k) \mathbb{E}\left[\frac{\pi_{tr}(k)}{\pi_t(k)} \ell(Y, f(\underline{X})) \mid Y = k\right] \\ &= \mathbb{E}_{\pi_t}\left[\frac{\pi_{tr}(Y)}{\pi_t(Y)} \ell(Y, f(\underline{X}))\right]\end{aligned}$$

- Put more weight on less probable classes!



- In unbalanced situation, often the **cost** of misprediction is not the same for all classes (e.g. medical diagnosis, credit lending...)
- Much better to use this explicitly than to do blind resampling!

Weighted Loss

- **Weighted loss:**

$$\ell(Y, f(\underline{X})) \rightarrow C(Y)\ell(Y, f(\underline{X}))$$

- Weighted risk target:

$$\mathbb{E}[C(Y)\ell(Y, f(\underline{X}))]$$

- **Rk:** Strong link with ℓ as C is independent of f .
- Often allow reusing algorithm constructed for ℓ .
- C may also depend on \underline{X} ...

- The Bayes classifier is now:

$$f^* = \operatorname{argmin} \mathbb{E}[C(Y)\ell(Y, f(\underline{X}))] = \operatorname{argmin} \mathbb{E}_{\underline{X}}[\mathbb{E}_{Y|\underline{X}}[C(Y)\ell(Y, f(\underline{X}))]]$$

Bayes Predictor

- For $\ell^{0/1}$ loss,

$$f^*(\underline{X}) = \operatorname{argmax}_k C(k)\mathbb{P}(Y = k|\underline{X})$$

- Same effect than a threshold modification for the binary setting!
- Allow putting more emphasis on some classes than others.

Cost and Proportions

- Testing risk target:

$$\mathbb{E}_{\pi_t}[C_t(Y)\ell(Y, f(\underline{X}))] = \sum_k \pi_t(k)C_t(k)\mathbb{E}[\ell(Y, f(\underline{X}))|Y = k]$$

- Training risk target

$$\mathbb{E}_{\pi_{tr}}[C_{tr}(Y)\ell(Y, f(\underline{X}))] = \sum_k \pi_{tr}(k)C_{tr}(k)\mathbb{E}[\ell(Y, f(\underline{X}))|Y = k]$$

- **Coincide if**

$$\pi_t(k)C_t(k) = \pi_{tr}(k)C_{tr}(k)$$

- Lots of flexibility in the choice of C_t , C_{tr} or π_{tr} .

Weighted Loss and Resampling

- **Weighted loss:** choice of a weight $C_t \neq 1$.
- **Resampling:** use a $\pi_{tr} \neq \pi_t$.
- Stratified sampling may be used to reduce the size of a dataset without losing a low probability class!

Combining Weights and Resampling

- **Weighted loss:** use $C_{tr} = C_t$ as $\pi_{tr} = \pi_t$.
- **Resampling:** use an implicit $C_t(k) = \pi_{tr}(k)/\pi_t(k)$.
- **Combined:** use $C_{tr}(k) = C_t(k)\pi_t(k)/\pi_{tr}(k)$
- Most ML methods allow such weights!



T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning.
Springer Series in Statistics, 2009



M. Mohri, A. Rostamizadeh, and A. Talwalkar.
Foundations of Machine Learning (2nd ed.)
MIT Press, 2018



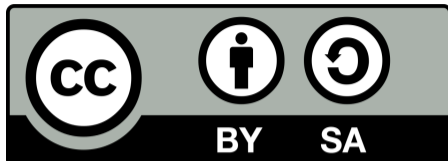
A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



Ch. Giraud.
Introduction to High-Dimensional Statistics.
CRC Press, 2014



S. Bubeck.
Convex Optimization: Algorithms and Complexity.
Now Publisher, 2015



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 **Optimization: Gradient Descent Algorithms**
 - **Introduction**
 - **Gradient Descent**
 - **Proximal Descent**
 - **Coordinate Descent**
 - **Gradient Descent Acceleration**
 - **Stochastic Gradient Descent**
 - **Gradient Descent Step**
 - **Non-Convex Setting**
 - **References**
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

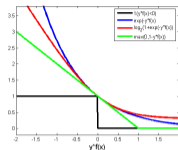
- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 **Optimization: Gradient Descent Algorithms**
 - **Introduction**
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Typical Optimization Problem in ML

- ML Opt.: $\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} G(\mathbf{w})$ with $G(\mathbf{w}) = F(\mathbf{w}) + R(\mathbf{w})$ where
 - F a goodness-of-fit function: $F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_{\mathbf{w}}(\underline{X}_i))$
where ℓ is some loss and $f_{\mathbf{w}}$ a parametric predictor,
 - R a regularizer: $R(\mathbf{w}) = \lambda \operatorname{pen}(\mathbf{w})$
where $\operatorname{pen}(\cdot)$ is some penalization function.

Examples

- Predictor:
 - Linear $f_{\mathbf{w}}(\underline{X}_i) = \underline{X}_i^\top \mathbf{w}$
 - Neural Nets...
- Regularizer:
 - $\operatorname{pen}(\mathbf{w}) = \|\mathbf{w}\|_2^2$ (ridge)
 - $\operatorname{pen}(\mathbf{w}) = \|\mathbf{w}\|_1$ (Lasso)



Classification

- Logistic loss, $\ell(y, f) = \log(1 + e^{-yf})$
- Hinge loss, $\ell(y, f) = (1 - yf)_+$
- Quadratic hinge loss, $\ell(y, f) = \frac{1}{2}(1 - yf)_+^2$
- Huber loss $\ell(y, f) = -4yf\mathbf{1}_{yf < -1} + (1 - yf)_+^2\mathbf{1}_{yf \geq -1}$
- These losses are convex upper bound of the 0/1 loss $\ell(y, y') = \mathbf{1}_{yy' \leq 0}$.

Regression

- Quadratic loss, $\ell(y, f) = \frac{1}{2}(y - f)^2$
- Absolute loss, $\ell(y, f) = \frac{1}{2}|y - f|$

ML Problem

- Minimization of

$$G(\mathbf{w}) = F(\mathbf{w}) + R(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \langle \underline{X}_i, \mathbf{w} \rangle) + \lambda \text{pen}(\mathbf{w})$$

Gradient and Hessian of F

- Gradient:

$$\nabla F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \bar{\ell}(Y_i, \langle \underline{X}_i, \mathbf{w} \rangle) \underline{X}_i$$

with $\bar{\ell}(y, f) = \frac{\partial \ell(y, f)}{\partial f}$

- Hessian matrix:

$$\nabla^2 F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \bar{\ell}'(Y_i, \langle \underline{X}_i, \mathbf{w} \rangle) \underline{X}_i \underline{X}_i^\top$$

with $\bar{\ell}'(y, f) = \frac{\partial^2 \ell(y, f)}{\partial f^2}$

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 **Optimization: Gradient Descent Algorithms**
 - Introduction
 - **Gradient Descent**
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Zero Order Heuristic

- Zero order approximation:

$$G(\mathbf{w}) = G(\mathbf{w}') + O(\|\mathbf{w} - \mathbf{w}'\|)$$

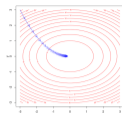
Optimization in a Compact Set

- Optimization problem:

$$\mathbf{w}^* \in \underset{\mathbf{w} \in [0,1]^d}{\operatorname{argmin}} G(\mathbf{w}).$$

Exhaustive Search

- **Algorithm:**
 - Evaluate G on a grid,
 - Approximate the minimizer by the minimizer in the grid.
- If G is C -Lipschitz, evaluating G on a grid of precision $\epsilon/(\sqrt{d}C)$ is sufficient to find a ϵ -minimizer of G .
- Required number of evaluation: $N_\epsilon = O\left((C\sqrt{d}/\epsilon)^d\right)$



First Order Heuristic

- First order approximation:

$$G(\mathbf{w}) = G(\mathbf{w}') + \nabla G(\mathbf{w}')^\top (\mathbf{w} - \mathbf{w}') + o(\|\mathbf{w} - \mathbf{w}'\|)$$

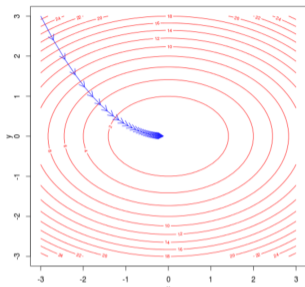
- Best descent direction: $-\nabla G(\mathbf{w}')$

Gradient Descent Algorithm

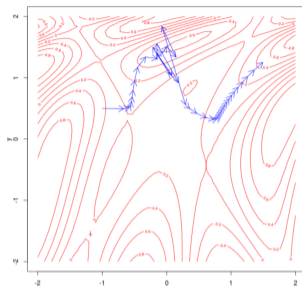
- Start from a point $\mathbf{w}^{[0]}$ and let $k = 0$.
- Repeat until convergence:
 - $\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \alpha^{[k]} \nabla G(\mathbf{w}^{[k]})$
 - $k \rightarrow k + 1$

with $\alpha^{[k]}$ a sequence of (small) steps.

Do The Gradient Descent Converge?



VS



- **Not always!**

Assumptions

- Convexity assumption on G .
- Regularity assumption on G
- Size of the steps $\alpha^{[k]}$.

Convex function

- A function $G : \mathbb{R}^d \rightarrow \mathbb{R}$ is **convex** on \mathbb{R}^d if, for all $x, y \in \mathbb{R}^d$, for all $\lambda \in [0, 1]$,
$$G(\lambda x + (1 - \lambda)y) \leq \lambda G(x) + (1 - \lambda)G(y).$$
- If $f_{\mathbf{w}}$ is linear, $F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_{\mathbf{w}}(X_i))$ is **convex** iff
$$f \mapsto \ell(y_i, f)$$
is convex for any $i = 1, \dots, n$.

Regularity

- G is **L -smooth** if it is continuously differentiable and if
$$\|\nabla G(\mathbf{w}) - \nabla G(\mathbf{w}')\|_2 \leq L \|\mathbf{w} - \mathbf{w}'\|_2 \quad \text{for any } \mathbf{w}, \mathbf{w}' \in \mathbb{R}^d$$
- If G is twice differentiable, this is equivalent to assuming
$$\lambda_{\max}(\nabla^2 G(\mathbf{w})) \leq L \quad \text{for any } \mathbf{w} \in \mathbb{R}^d$$
(largest eigenvalue of the Hessian matrix of G smaller than L)

Formulation

- Minimization of

$$G(\mathbf{w}) = F(\mathbf{w}) + R(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (Y_i - \underline{X}_i^\top \mathbf{w})^2 + \lambda \text{pen}(\mathbf{w})$$

Gradient and Hessian of F

- Gradient:

$$\nabla F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (\underline{X}_i^\top \mathbf{w} - Y_i) \underline{X}_i$$

- Hessian:

$$\nabla^2 F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \underline{X}_i \underline{X}_i^\top$$

so that

$$L = \lambda_{\max} \left(\frac{1}{n} \sum_{i=1}^n \underline{X}_i \underline{X}_i^\top \right)$$

Formulation

- Minimization of

$$G(\mathbf{w}) = F(\mathbf{w}) + R(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{(-Y_i(\underline{X}_i^\top \mathbf{w}))} \right) + \lambda \text{pen}(\mathbf{w})$$

Gradient and Hessian of F

- Gradient:

$$\nabla F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n Y_i (\sigma(Y_i \underline{X}_i^\top \mathbf{w}) - 1) \underline{X}_i$$

- Hessian:

$$\nabla^2 F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \sigma(Y_i \underline{X}_i^\top \mathbf{w}) (1 - \sigma(Y_i \underline{X}_i^\top \mathbf{w})) \underline{X}_i \underline{X}_i^\top$$

so that F is L -smooth with

$$L \leq \frac{1}{4} \lambda_{\max} \left(\frac{1}{n} \sum_{i=1}^n \underline{X}_i \underline{X}_i^\top \right)$$

- $\sigma(t) = e^{-t} / (1 + e^{-t})$.

Theorem

- Let $G : \mathbb{R}^d \rightarrow \mathbb{R}$ be a L -smooth convex function. Let \mathbf{w}^* be the minimum of f on \mathbb{R}^d . Then, Gradient Descent with step size $\alpha \leq 1/L$ satisfies

$$G(\mathbf{w}^{[k]}) - G(\mathbf{w}^*) \leq \frac{\|\mathbf{w}^{[0]} - \mathbf{w}^*\|_2^2}{2\alpha k}.$$

- In particular, for $\alpha = 1/L$,

$$N_\epsilon = O(L\|\mathbf{w}^{[0]} - \mathbf{w}^*\|_2^2 / (2\epsilon))$$

iterations are sufficient to get an ϵ -approximation of the minimal value of G .

- Bound is independent of the dimension d .
- Weak dependency** hidden in the constant L ...

A **key** Property: the Descent Lemma.

- If G is convex and L -smooth, then for any $\mathbf{w}, \mathbf{w}' \in \mathbb{R}^d$

$$G(\mathbf{w}) \leq G(\mathbf{w}') + \nabla G(\mathbf{w}')^\top (\mathbf{w} - \mathbf{w}') + \frac{L}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2.$$

Link with the Gradient Descent Algorithm

- At step k , for any \mathbf{w} ,

$$\begin{aligned} G(\mathbf{w}) &\leq G(\mathbf{w}^{[k]}) + \nabla G(\mathbf{w}^{[k]})^\top (\mathbf{w} - \mathbf{w}^{[k]}) + \frac{L}{2} \|\mathbf{w} - \mathbf{w}^{[k]}\|_2^2 \\ &\leq G(\mathbf{w}^{[k]}) + \frac{L}{2} \left\| \mathbf{w} - \left(\mathbf{w}^{[k]} - \frac{1}{L} \nabla G(\mathbf{w}^{[k]}) \right) \right\|_2^2 - \frac{1}{2L} \left\| \nabla G(\mathbf{w}^{[k]}) \right\|_2^2 \end{aligned}$$

- Optimizing the upper bound in \mathbf{w} leads to

$$\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \frac{1}{L} \nabla G(\mathbf{w}^{[k]})$$

- **Rk:** Newton iteration comes from the approximation

$$G(\mathbf{w}) = G(\mathbf{w}') + \nabla G(\mathbf{w}')^\top (\mathbf{w} - \mathbf{w}') + \frac{1}{2} (\mathbf{w} - \mathbf{w}')^\top \nabla^2 G(\mathbf{w}') (\mathbf{w} - \mathbf{w}') + o(\|\mathbf{w} - \mathbf{w}'\|^2)$$

Strong convexity

- A function $G : \mathbb{R}^d \rightarrow \mathbb{R}$ is **μ -strongly convex** if $x \mapsto G(x) - \frac{\mu}{2} \|x\|_2^2$ is convex.
- If G is differentiable, this is equivalent to, for all $x, y \in \mathbb{R}^d$,
$$G(y) \geq G(x) + \nabla G(x)^\top (y - x) + \frac{\mu}{2} \|y - x\|_2^2.$$
- If G is twice differentiable, this is equivalent to
$$\lambda_{\min}(\nabla^2 G(x)) \geq \mu.$$

Theorem

- Let $G : \mathbb{R}^d \rightarrow \mathbb{R}$ be a L -smooth, μ strongly convex function. Let \mathbf{w}^* be the minimum of G on \mathbb{R}^d . Then, Gradient Descent with step size $\alpha \leq 1/L$ satisfies
$$G(\mathbf{w}^{[k]}) - G(\mathbf{w}^*) \leq \frac{1}{2\alpha} (1 - \alpha\mu)^k \|G(\mathbf{w}^{[0]}) - G(\mathbf{w}^*)\|_2^2.$$
- **Rk:** $N_\epsilon = O(-\log \epsilon / (\alpha\mu))$ iterations are sufficient to obtain an ϵ -minimizer.

Convexity

- If G is a convex function then for any $\mathbf{w} \in \mathbb{R}^d$, there exists a (not necessarily unique) subgradient $\delta G(\mathbf{w})$ such that

$$G(\mathbf{w}') \geq G(\mathbf{w}) + \delta G(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) \quad \text{for any } \mathbf{w}' \in \mathbb{R}^d$$

- If G is differentiable then it is unique and equal to the gradient.
- **Subgradient Descent Algorithm:** Gradient Descent with subgradient instead of gradient.

Regularity

- G is **C-Lipschitz** if

$$|G(\mathbf{w}) - G(\mathbf{w}')| \leq C \|\mathbf{w} - \mathbf{w}'\|_2 \quad \text{for any } \mathbf{w}, \mathbf{w}' \in \mathbb{R}^d$$

- If G is differentiable, this is equivalent to assuming

$$\|\nabla G(\mathbf{w})\| \leq C \quad \text{for any } \mathbf{w} \in \mathbb{R}^d$$

Theorem

- Let $G : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function, C -Lipschitz in $B(\mathbf{w}^*, R)$ where \mathbf{w}^* be the minimizer of f on \mathbb{R}^d . Assume that

$$\alpha^{[k]} > 0, \quad \alpha^{[k]} \rightarrow 0, \quad \sum_k \alpha^{[k]} = +\infty$$

and $\|\mathbf{w}^{[0]} - \mathbf{w}^*\| \leq R$ Then, Subgradient Descent with step size $\alpha^{[k]}$ satisfies

$$\min_{k' \leq k} G(\mathbf{w}^{[k']}) - G(\mathbf{w}^*) \leq C \frac{R^2 + \sum_{k'=0}^k (\alpha^{[k']})^2}{2 \sum_{k'=0}^k \alpha^{[k]}}$$

- In particular, for $\alpha^{[k]} = \frac{r}{\sqrt{k+1}}$

$$\min_{k' \leq k} G(\mathbf{w}^{[k']}) - G(\mathbf{w}^*) \leq C \frac{R^2 + r^2 \log(k+1)}{4r\sqrt{k+1}}$$

and

$$N_\epsilon = O\left(\left(C(-\log \epsilon)/\epsilon\right)^2\right)$$

In practice, how to choose α ?

- **Theoretical setting of $\alpha = 1/L$ is a worst case scenario.**

Exact line search

- At each step, choose the best α by optimizing

$$\alpha^{[k]} = \underset{\alpha > 0}{\operatorname{argmin}} G(\mathbf{w} - \alpha \nabla G(\mathbf{w})).$$

- **Too costly!**

Backtracking line search

- Fix a parameter $0 < \beta < 1$ and set $\alpha = \alpha_{init}$
- At each iteration,
 - while $G(\mathbf{w} - \alpha \nabla G(\mathbf{w})) > G(\mathbf{w}) - \frac{\alpha}{2} \|\nabla G(\mathbf{w})\|^2$,
modify $\alpha \leftarrow \beta \alpha$.
 - Use the final α as the stepsize for this iteration.
- **Simple and work pretty well in practice.**
- Theoretical guarantees available.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 **Optimization: Gradient Descent Algorithms**
 - Introduction
 - Gradient Descent
 - **Proximal Descent**
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Function Decomposition: $G = F + R$

- F convex and L -**smooth**
- R convex and **simple**

Another Descent Algorithm

- At step k , for any \mathbf{w} ,
 $G(\mathbf{w}) \leq F(\mathbf{w}) + R(\mathbf{w})$

$$\begin{aligned} &\leq F(\mathbf{w}^{[k]}) + \nabla F(\mathbf{w}^{[k]})^\top (\mathbf{w} - \mathbf{w}^{[k]}) + \frac{L}{2} \|\mathbf{w} - \mathbf{w}^{[k]}\|_2^2 + R(\mathbf{w}) \\ &\leq F(\mathbf{w}^{[k]}) + \frac{L}{2} \left\| \mathbf{w} - \left(\mathbf{w}^{[k]} - \frac{1}{L} \nabla F(\mathbf{w}^{[k]}) \right) \right\|^2 - \frac{1}{2L} \|\nabla F(\mathbf{w}^{[k]})\|^2 + R(\mathbf{w}) \end{aligned}$$

- Optimizing the upper bound in \mathbf{w} leads to

$$\mathbf{w}^{[k+1]} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} \left\| \mathbf{w} - \left(\mathbf{w}^{[k]} - \frac{1}{L} \nabla F(\mathbf{w}^{[k]}) \right) \right\|^2 + \frac{1}{L} R(\mathbf{w})$$

- R **simple** means that this minimization is easy.

Proximal Operator

- For any convex function R :

$$\text{prox}_\gamma R(\mathbf{w}') = \underset{\mathbf{w}}{\text{argmin}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}'\|^2 + \gamma R(\mathbf{w})$$

- Generalization of the projection operator:

- $R(\mathbf{w}) = \mathbf{1}_\Omega$: $\text{prox}_\gamma R(\mathbf{w}') = P_\Omega \mathbf{w}'$

- $R(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2$: $\text{prox}_\gamma R(\mathbf{w}') = \frac{1}{1+\gamma} \mathbf{w}'$.

- $R(\mathbf{w}) = \|\mathbf{w}\|_1$: $\text{prox}_\gamma R(\mathbf{w}') = T_\gamma(\mathbf{w}')$ with $T_\gamma(\mathbf{w})_i = \text{sign}(\mathbf{w}_i) \max(0, |\mathbf{w}_i| - \gamma)$
(soft thresholding).

Proximal Gradient Descent Algorithm for simple R

- Start from a point $\mathbf{w}^{[0]}$ and let $k = 0$.
- Repeat until convergence:
 - $\mathbf{w}^{[k+1]} = \text{prox}_{1/L} R(\mathbf{w}^{[k]} - \alpha^{[k]} \nabla F(\mathbf{w}^{[k]}))$
 - $k \rightarrow k + 1$

with $\alpha^{[k]}$ a sequence of (small) steps.

- **Same as smooth case!**

Rates

- F L -smooth and R simple:

$$G(\mathbf{w}^{[k]}) - G(\mathbf{w}^*) \leq \frac{\|\mathbf{w}^{[0]} - \mathbf{w}^*\|_2^2}{2\alpha k}.$$

and $N_\epsilon = O(L\|\mathbf{w}^{[0]} - \mathbf{w}^*\|_2^2/2\epsilon)$.

- F L -smooth and μ -convex and R simple:

$$G(\mathbf{w}^{[k]}) - G(\mathbf{w}^*) \leq \frac{1}{2\alpha} (1 - \alpha\mu)^k \|G(\mathbf{w}^{[0]}) - G(\mathbf{w}^*)\|_2^2.$$

and $N_\epsilon = O(-\log \epsilon / (\alpha\mu))$.

- F C -Lipschitz and R is the characteristic function of a convex set:

$$G(\mathbf{w}^{[k]}) - G(\mathbf{w}^*) \leq C \frac{R^2 + r^2 \log(k+1)}{4r\sqrt{k+1}}$$

and $N_\epsilon = O((C(-\log \epsilon)/\epsilon)^2)$.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 **Optimization: Gradient Descent Algorithms**
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - **Coordinate Descent**
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Heuristic

- Is \mathbf{w} a minimizer of a G when

$$G(\mathbf{w} + t\mathbf{e}_i) \geq G(\mathbf{w}) \quad \forall t \in \mathbb{R} \text{ and } 1 \leq i \leq d,$$

with G convex?

- **YES** if G is smooth (\mathcal{C}^1), G is separable or G is a sum of a smooth (\mathcal{C}^1) and a separable function ($G(\mathbf{w}) = \sum_{i=1}^d g_i(\mathbf{w}^{(i)})$).
- **NO** otherwise.

Coordinate Descent Algorithm

- Start from an arbitrary $\mathbf{w}^{[0]}$
- Repeat until convergence:
 - Pick a coordinate i
 - Set $\mathbf{w}^{[k+1],(j)} = \mathbf{w}^{[k],(j)}$ for $j \neq i$.
 - **Optimize** only the i th coordinate to obtain $\mathbf{w}^{[k+1],(i)}$
 - $k \rightarrow k + 1$
- State-of-the art is several ML problems!

Exact Coordinate Descent (CD)

- Start from an arbitrary $\mathbf{w}^{[0]}$
- Repeat until convergence:
 - Pick a coordinate i
 - Set $\mathbf{w}^{[k+1],(j)} = \mathbf{w}^{[k],(j)}$ for $j \neq i$.
 - Compute

$$\mathbf{w}^{[k+1],(i)} = \underset{z \in \mathbb{R}}{\operatorname{argmin}} G(\mathbf{w}^{[k],(1)}, \dots, \mathbf{w}^{[k],(i-1)}, z, \mathbf{w}^{[k],(i+1)}, \dots, \mathbf{w}^{[k],(d)})$$

- $k \rightarrow k + 1$
- Several way to choose the coordinate i : uniform sampling, deterministic cycle...
- Only 1D optimization problems to solve... but probably a lot of them.

Theorem - Warga (1963)

- If G is continuously differentiable and strictly convex, then exact coordinate descent converges to a minimum.

Coordinate Gradient Descent (CGD)

- Start from an arbitrary $\mathbf{w}^{[0]}$
- Repeat until convergence:
 - Pick a coordinate i
 - Set $\mathbf{w}^{[k+1],(j)} = \mathbf{w}^{[k],(j)}$ for $j \neq i$.
 - Compute

$$\mathbf{w}^{[k+1],(i)} = \mathbf{w}^{[k],(i)} - \alpha^{[k]} \nabla^{(i)} G(\mathbf{w}^{[k]})$$

- $k \rightarrow k + 1$

- For smooth function, step-size $\alpha^{[k]}$ can be taken as $\alpha^{[k]} = 1/L_i$ where i is the coordinate chosen and L_i the Lipschitz constant of

$$G^i(z) = G(\mathbf{w} + z\mathbf{e}_i) = G(\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(i-1)}, \mathbf{w}^{(i)} + z, \mathbf{w}^{(i+1)}, \dots, \mathbf{w}^{(d)})$$

- If G is L -smooth, $L_i \leq L$.

Theorem - Nesterov (2012)

- Assume that G is convex and smooth and that each G^i is L_i -smooth. Consider a sequence $\{\mathbf{w}^{[k]}\}$ given by CGD with $\alpha^{[k]} = 1/L_{i_k}$ and coordinates i_1, i_2, \dots chosen at random: i.i.d and uniform distribution in $\{1, \dots, d\}$. Then

$$\mathbb{E} \left[G(\mathbf{w}^{[k+1]}) - G(\mathbf{w}^*) \right] \\ \leq \frac{d}{d+k} \left(\left(1 - \frac{1}{d}\right) (G(\mathbf{w}^{[0]}) - G(\mathbf{w}^*)) + \frac{1}{2} \left\| \mathbf{w}^{[0]} - \mathbf{w}^* \right\|_L^2 \right),$$

with $\left\| \mathbf{w} \right\|_L^2 = \sum_{j=1}^d L_j \mathbf{w}_j^2$.

- Bound in expectation, since coordinates are taken at random.
- For cycling coordinates $i_k = (k \bmod d) + 1$ the bound is much worse.
- Similar result when $G = F + R$ with F smooth and R simple and separable.

Comparison of Gradient Descent and Coordinate Gradient Descent

Gradient Descent

- Cost of an iteration $O(d)$.
- Number of iteration to obtain an ϵ -minimizer:

$$N_\epsilon = \frac{L \left\| \mathbf{w}^{[0]} - \mathbf{w}^* \right\|_2^2}{2\epsilon}$$

Coordinate Gradient Descent

- Cost of an iteration $O(1)$.
- Number of iteration to obtain an ϵ -minimizer:

$$N_\epsilon = \frac{d}{\epsilon} \left(\left(1 - \frac{1}{d}\right) (G(\mathbf{w}^{[0]}) - G(\mathbf{w}^*)) + \frac{1}{2} \left\| \mathbf{w}^{[0]} - \mathbf{w}^* \right\|_L^2 \right)$$

- Same order of complexity but smaller constant for CGD as soon as $(G(\mathbf{w}^{[0]}) - G(\mathbf{w}^*)) \ll L \left\| \mathbf{w}^{[0]} - \mathbf{w}^* \right\|_2^2$.
- In practice, often much faster (especially when dealing with anisotropic regularity).

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 **Optimization: Gradient Descent Algorithms**
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - **Gradient Descent Acceleration**
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

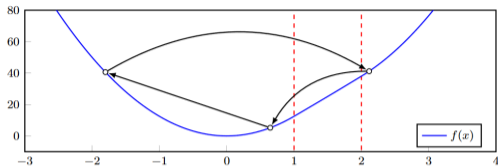
- Can we improve the number of iterations $O(L/\epsilon)$ (L -smooth) and $O(\frac{L}{\mu} \log(1/\epsilon))$ (L -smooth and μ strongly-convex) ?
- **Heuristic:** add some momentum to propagate gradients.

Polyak's momentum algorithm - Heavy ball method

- **Input:** starting point $\mathbf{w}^{[0]}$, step size $\alpha^{[k]} > 0$, momentum $\beta \in [0, 1]$
- While *not converge* do
 - $\mathbf{w}^{[k+1/2]} = \mathbf{w}^{[k]} - \alpha^{[k]} \nabla G(\mathbf{w}^{[k]})$
 - $\mathbf{w}^{[k+1]} = \mathbf{w}^{[k+1/2]} + \beta(\mathbf{w}^{[k]} - \mathbf{w}^{[k-1]})$
 - $k \leftarrow k + 1$
- **Return** last $\mathbf{w}^{[k+1]}$.

- Most classical value for β : 0.9
- When $\alpha^{[k]} = \alpha$, the update equation can be written

$$\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \alpha \sum_{t=1}^k \beta^{k-t} \nabla G(\mathbf{w}^{(t)}).$$



Counter Example

- Polyak's momentum algorithm fails to converge in some specific cases, for instance:

$$\nabla G(x) = \begin{cases} 25x & \text{if } x < 1 \\ x + 24 & \text{if } 1 \leq x < 2 \\ 25x - 24 & \text{if } x \geq 2 \end{cases}$$

- In that case, G is μ strongly convex and L -smooth with $(\mu, L) = (1, 25)$. However, iterations given by Polyak's algorithm cycles.

Nesterov Method

- **Input:** starting point $\mathbf{w}^{[0]} (= \mathbf{w}^{[-1/2]})$, step size $\alpha^{[k]} > 0$, momentum $\beta^{[k]} \in [0, 1]$
- While *not converge* do
 - $\mathbf{w}^{[k+1/2]} = \mathbf{w}^{[k]} - \alpha^{[k]} \nabla G(\mathbf{w}^{[k]})$
 - $\mathbf{w}^{[k+1]} = \mathbf{w}^{[k+1/2]} + \beta^{[k]} (\mathbf{w}^{[k+1/2]} - \mathbf{w}^{[k-1/2]})$
 - $k \leftarrow k + 1$
- **Return** last $\mathbf{w}^{[k+1]}$.

- Subtle difference with momentum method.
- Not much freedom for $\beta^{[k]}$:
 - Nesterov original choice:

$$\beta^{[k]} = \frac{t^{[k]} - 1}{t^{[k+1]}} \text{ with } \begin{cases} t^{[1]} = 1 \\ t^{[k+1]} = \frac{1 + \sqrt{1 + 4(t^{[k]})^2}}{2} \end{cases}$$

- Simpler equivalent choice $\beta^{[k]} = \frac{k}{k+3}$
- Other choice possible for μ convex function.
- Proximal version possible.

Theorem for L -smooth function

- Assume that G is a L -smooth, convex function whose minimum is reached at \mathbf{w}^* . Then, if $\beta^{[k]} = k/(k+3)$,

$$G(\mathbf{w}^{[k]}) - G(\mathbf{w}^*) \leq \frac{2\|\mathbf{w}^{[0]} - \mathbf{w}^*\|_2^2}{\alpha(k+1)^2}.$$

Theorem for μ strongly convex function

- Assume that G is a L -smooth, μ strongly convex function whose minimum is reached at \mathbf{w}^* . Then, if $\beta^{[k]} = \frac{1-\sqrt{\mu/L}}{1+\sqrt{\mu/L}}$,

$$G(\mathbf{w}^{[k]}) - G(\mathbf{w}^*) \leq \frac{\|\mathbf{w}^{[0]} - \mathbf{w}^*\|_2^2}{\alpha} \left(1 - \sqrt{\frac{\mu}{L}}\right)^k.$$

- Faster rates than respectively $O(1/k)$ and $O((1 - \mu/L)^k)$.
- Theorems holds for the proximal algorithm with $G = F + R$ with F smooth and R simple.

General First Order Method

- Any iterative method that generates a sequence $\{\mathbf{w}^{[k]}\}$ s.t.
$$\mathbf{w}^{[k]} \in \mathbf{w}^{[0]} + \text{Span}(\nabla f(\mathbf{w}^{[0]}), \dots, \nabla f(\mathbf{w}^{[k-1]})).$$

Lower Bounds

- For any $\mathbf{w}^{[0]} \in \mathbb{R}^d$ and any k satisfying $1 \leq k \leq (d-1)/2$, there exists a L -smooth convex function f such that for any general first order method

$$G(\mathbf{w}^{[k]}) - G(\mathbf{w}^*) \geq \frac{3L \|\mathbf{w}^{[0]} - \mathbf{w}^*\|_2^2}{32(k+1)^2}.$$

- For any $\mathbf{w}^{[0]} \in \mathbb{R}^d$ and any $k \leq (d-1)/2$, there exists a L -smooth, μ strongly convex function f such that for any general first order method

$$G(\mathbf{w}^{[k]}) - G(\mathbf{w}^*) \geq \frac{\mu}{2} \left(\frac{1 - \sqrt{\mu/L}}{1 + \sqrt{\mu/L}} \right)^{2k} \|\mathbf{w}^{[0]} - \mathbf{w}^*\|_2^2.$$

- Accelerated rates: best possible without further assumptions.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 **Optimization: Gradient Descent Algorithms**
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - **Stochastic Gradient Descent**
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Goodness-of-Fit Optimization

- Minimizer of $G = F + R$ with
 - F smooth and R simple (as before).
 - Empirical average structure for F :

$$F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n F_i(\mathbf{w}) = \mathbb{E}_n[F_i(\mathbf{w})]$$

where each F_i is smooth.

- Cost of evaluating ∇F is proportional to the dataset cardinality n (and the dimension d).
- If the dataset is large, this can be an issue.
- **Stochastic Gradient Heuristic**: replace $\nabla F(\mathbf{w}) = \mathbb{E}_n[\nabla F_i(\mathbf{w})]$ by a Monte Carlo approximation.

Stochastic Gradient

- Empirical average structure

$$\nabla F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \nabla F_i(\mathbf{w}) = \mathbb{E}_n[\nabla F_i(\mathbf{w})]$$

- Monte Carlo approximation of an average:

- draw uniformly $m \ll n$ indices i_1, \dots, i_m
- Replace $\nabla F(\mathbf{w})$ by

$$\hat{\nabla} F(\mathbf{w}) = \frac{1}{m} \sum_{j=1}^m \nabla F_{i_j}(\mathbf{w})$$

- **Prop:** $\mathbb{E}[\hat{\nabla} F(\mathbf{w})] = \nabla F(\mathbf{w})$

- Extreme case $m = 1$: unbiased but quite noisy estimate of the true gradient.
- Evaluation cost of $\hat{\nabla} F(\mathbf{w})$ independent of n .
- **Question:** Can we use this approximation of the true gradient in a descent algorithm?

Stochastic Gradient Descent Algorithm

- Start from a point $\mathbf{w}^{[0]}$ and let $k = 0$.
- Repeat until convergence:
 - Draw m indices $1 \leq i_j \leq n$ uniformly at random
 - Set $\widehat{\nabla}F(\mathbf{w}^{[k]}) = \frac{1}{m} \sum_{j=1}^m \nabla F_{i_j}(\mathbf{w}^{[k]})$
 - $\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \alpha^{[k]} \left(\widehat{\nabla}F(\mathbf{w}^{[k]}) + \nabla R(\mathbf{w}^{[k]}) \right)$
 - $k \rightarrow k + 1$

with $\alpha^{[k]}$ a sequence of (small) steps.

- Each iteration has complexity $O(md)$ instead of $O(nd)$ for full gradient methods
- m is called the batch size.
- **Rk:** no gain when using a batch size larger than 1, except if the sum can be parallelized.
- Proximal variant if R is simple instead of smooth.

- Theoretical analysis requires some modifications:
 - at each step $\mathbf{w}^{[k+1]}$ should be projected into a ball $B(c, R)$ with $R > 0$.
 - G should be convex such that $\sup_i |\nabla F_i(\mathbf{w}) + \nabla R(\mathbf{w})| \leq b, \forall \mathbf{w} \in B(c, r)$
 - **Polyak-Ruppert** averaging: use SGD iterates w^t but return an average.

SGD Rate

- With $\alpha^{[k]} = 2R/(b\sqrt{k})$

$$\mathbb{E} \left[G \left(\frac{1}{k} \sum_{j=1}^k \mathbf{w}^{[j]} \right) \right] - G(\mathbf{w}^*) \leq \frac{3rb}{\sqrt{k}}$$

- If G is μ -strictly convex then with $\alpha^{[k]} = 2/(\mu(k+1))$,

$$\mathbb{E} \left[G \left(\frac{2}{k(k+1)} \sum_{j=1}^k j \mathbf{w}^{[j]} \right) \right] - G(\mathbf{w}^*) \leq \frac{2b^2}{\mu(k+1)}.$$

- Without averaging, logarithmic loss.
- Same rate for proximal algorithm if F is at least L -smooth and R simple.

Simplified Result for the (Averaged) SGD

- If G is convex and gradients are bounded ($\|\nabla F_i(\mathbf{w})\|_2 \leq b$) then the convergence rate is

$$O\left(\frac{1}{\sqrt{k}}\right) \quad \text{with} \quad \alpha^{[k]} = O\left(\frac{1}{\sqrt{k}}\right)$$

- G is μ -strongly convex, the rate is

$$O\left(\frac{1}{\mu k}\right) \quad \text{with} \quad \alpha^{[k]} = O\left(\frac{1}{\mu k}\right)$$

Comparison with GD

- Much slower rate:
 - for L smooth function: $O(1/\sqrt{k})$ vs $O(1/k)$
 - For μ convex function: $O(1/k)$ vs $O((1 - \mu/L)^k)$
 - Hidden factor σ in SGD measuring the variance of $\hat{\nabla}F$ which decreases when m increases.
- Much cheaper cost per iteration: $O(md)$ vs $O(nd)$

Stochastic Gradient Descent Algorithm

- Start from a point $\mathbf{w}^{[0]}$ and let $k = 0$.
- Repeat until convergence:
 - Draw m indices $1 \leq i_j \leq n$ uniformly at random
 - Set $\widehat{\nabla}F(\mathbf{w}^{[k]}) = \frac{1}{m} \sum_{k=1}^m \nabla F_{i_j}(\mathbf{w}^{[k]})$
 - $\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \alpha \left(\widehat{\nabla}F(\mathbf{w}^{[k]}) + \nabla R(\mathbf{w}^{[k]}) \right)$
 - $k \rightarrow k + 1$

with α constant.

- No convergence but for μ -convex function:

$$\mathbb{E} \left[\left\| \mathbf{w}^{[k]} - \mathbf{w}^* \right\|^2 \right] \rightarrow O(\alpha/\mu)$$

- SGD behavior in practice: very fast initially but much slower afterward.
- Justify the use of decreasing step size when the error decay becomes slow.

- Difference due to the variance of the stochastic gradient.

Variance reduction of the gradient

- In the iterations of SGD, replace $\nabla F_{i_k}(\mathbf{w}^{[k-1]})$ by

$$\nabla F_{i_k}(\mathbf{w}^{[k-1]}) - \nabla F_{i_k}(\tilde{\mathbf{w}}) + \nabla F(\tilde{\mathbf{w}})$$

where $\tilde{\mathbf{w}}$ is an **old** value of the iterate, namely use

$$\mathbf{w}^{[k]} \leftarrow \mathbf{w}^{[k-1]} - \alpha^{[k]} ((\nabla F_{i_k}(\mathbf{w}^{[k-1]}) - \nabla F_{i_k}(\tilde{\mathbf{w}})) + \nabla F(\tilde{\mathbf{w}}))$$

- Several instantiations
 - SVRG where $\tilde{\mathbf{w}}$ and $\nabla F(\tilde{\mathbf{w}})$ is computed every few step.
 - SAGA where, instead of $\nabla F_{i_k}(\tilde{\mathbf{w}})$, one use $F_i(\mathbf{w}^{[k_i]})$, where k_i is the last time the i th coordinate has been updated, even to compute the equivalent of $\nabla F(\tilde{\mathbf{w}})$.
- Lead to faster rate but requires more memory ($3d$ for SVRG and nd for SAGA).

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 **Optimization: Gradient Descent Algorithms**
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - **Gradient Descent Step**
 - Non-Convex Setting
 - References
 - 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
 - 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
 - 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
 - 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
 - 9 References

Stochastic Gradient Descent Algorithm

- Start from a point $\mathbf{w}^{(0)}$ and let $k = 0$.
- Repeat until convergence:
 - Draw m indices $1 \leq i_j \leq n$ uniformly at random
 - Set $\hat{\nabla} F(\mathbf{w}^{[k]}) = \frac{1}{m} \sum_{k=1}^m \nabla F_{i_j}(\mathbf{w}^{[k]})$
 - $\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \alpha^{[k]} \left(\hat{\nabla} F(\mathbf{w}^{[k]}) + \nabla R(\mathbf{w}^{[k]}) \right)$
 - $k \rightarrow k + 1$

with $\alpha^{[k]}$ a sequence of (small) steps.

- Step sizes $\alpha^{[k]}$ are crucial for convergence.
- How to choose them? Decaying? Fixed? **Adaptive?**
- Several algorithms available with no clear winner.

Determining a good learning rate becomes more of an art than science for many problems.

Gradient Descent and Reparameterization

- Gradient Descent: $\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \alpha \nabla G(\mathbf{w}^{[k]})$.
- Reparameterization: $\mathbf{w} = W\mathbf{z}$.
- Gradient Descent in \mathbf{z} : $\mathbf{z}^{[k+1]} = \mathbf{z}^{[k]} - \alpha W^\top \nabla G(W\mathbf{z}^{[k]})$.
- Implied Gradient Descent in \mathbf{w} : $\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \alpha WW^\top \nabla G(\mathbf{w}^{[k]})$.
- **Different dynamics if $WW^\top \neq \text{Id}$!**

Gradient Descent and Norm

- Descent lemma upperbound: $f(\mathbf{w}') \leq f(\mathbf{w}) + \nabla G(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) + \frac{1}{2\alpha} \|\mathbf{w}' - \mathbf{w}\|^2$.
- Uses the classical euclidean norm.
- Using $f(\mathbf{w}') \leq f(\mathbf{w}) + \nabla G(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) + \frac{1}{2\alpha} \|\mathbf{w}' - \mathbf{w}\|_\Sigma^2$ leads to $\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \alpha \Sigma^{-1} \nabla G(\mathbf{w}^{[k]})$.
- **Different dynamics if $\Sigma \neq \text{Id}$!**

Modified Gradient Descent Algorithm

- $\mathbf{w}^{[k+1]} = \mathbf{w}^{[k]} - \alpha P \nabla G(\mathbf{w}^{[k]})$ where P can be interpreted as WW^T or Σ of the previous slide.
- Convergence holds if corresponding descent lemma holds.
- P can even changes from one iteration to another.

Choices for P

- Newton method:
 - $P = (\nabla^2 G)^{-1}$
 - Hard to compute and no descent lemma property.
- Diagonal approximation:
 - $P = D$
 - Easy computation but no descent lemma in general.
 - Several strategies for the choice of D ...

- Coordinate Descent could be reframed in this framework!

$$\mathbf{w}_j^{[k+1]} = \mathbf{w}_j^{[k]} + \frac{\alpha_b^{[k]}}{d_j^{[k]}} m_j^{[k]}$$

Self Normalized Gradient Descent

- Gradient: $m^{[k]} \simeq \nabla G(\mathbf{w}^{[k]})$
- Renormalization factor: $d_j^{[k]} \simeq |\nabla G(\mathbf{w}^{[k]})|_j$
- Base stepsize: $\alpha_b^{[k]}$ could be constant or decaying. . .
- Makes the algorithm invariant to a diagonal rescaling.
- Amounts to use

$$\|\mathbf{w}' - \mathbf{w}\|_{\Sigma}^2 \simeq \sum_j \frac{|\mathbf{w}'_j - \mathbf{w}_j|^2}{|\nabla G(\mathbf{w})_i|^2}$$

$$\mathbf{w}_j^{[k+1]} = \mathbf{w}_j^{[k]} + \frac{\alpha_b^{[k]}}{d_j^{[k]}} m_j^{[k]}$$

ADaptive GRADient update

- $m^{[k]} = \nabla G(\mathbf{w}^{[k]})$
- $d_j^{[k]} = \sqrt{\sum_{k'=1}^k (\nabla G(\mathbf{w}^{[k']}))_j^2}$
- $\alpha_b^{[k]} = \alpha$

- Step sizes grow as the inverse of the gradient magnitude.
- Accumulation of the gradients acts as a decreasing learning rate.
- Sensitive to initial condition and may require large initial parameter or restarts.

$$\mathbf{w}_j^{[k+1]} = \mathbf{w}_j^{[k]} + \frac{\alpha_b^{[k]}}{d_j^{[k]}} m_j^{[k]}$$

ADaptive GRADient update reformulation

- $m^{[k]} = \nabla G(\mathbf{w}^{[k]})$
- $d_j^{[k]} = \sqrt{v_j^{[k]}}$ with $v_j^{[k]} = \frac{1}{k} \sum_{k'=1}^k (\nabla G(\mathbf{w}^{[k']}))_j^2$
- $\alpha_b^{[k]} = \alpha / \sqrt{k}$
- Explicit averaging and decay.

$$\mathbf{w}_j^{[k+1]} = \mathbf{w}_j^{[k]} + \frac{\alpha_b^{[k]}}{d_j^{[k]}} m_j^{[k]}$$

RMSprop update

- $m^{[k]} = \nabla G(\mathbf{w}^{[k]})$
 - $d_j^{[k]} = \sqrt{v_j^{[k]}}$ with $v_j^{[k]} = \rho v_j^{[k-1]} + (1 - \rho)(\nabla G(\mathbf{w}^{[k]})_j)^2$
 - $\alpha_b^{[k]} = \alpha$
-
- Unpublished method!
 - Exponential average instead of classical one.
 - No step size decay.
 - Often very efficient (especially in a non stationary setting).

$$\mathbf{w}_j^{[k+1]} = \mathbf{w}_j^{[k]} + \frac{\alpha_b^{[k]}}{d_j^{[k]}} m_j^{[k]}$$

ADAM

- Exponential averaging estimation: let $\tilde{m}_0 = \tilde{v}_0 = 0$,
 $\tilde{m}^{[k]} = \beta_m \tilde{m}^{[k-1]} + (1 - \beta_m) \nabla G(\mathbf{w}^{[k]})$ and $\tilde{v}^{[k]} = \beta_v \tilde{v}^{[k-1]} + (1 - \beta_v) (\nabla G(\mathbf{w}^{[k]}))^2$
 - Bias correction: $m^{[k]} = \frac{\tilde{m}^{[k]}}{1 - \beta_m^k}$ and $v^{[k]} = \frac{\tilde{v}^{[k]}}{1 - \beta_v^k}$
 - $\alpha_b^{[k]} = \alpha$.
-
- Use a smoothed estimation of the first two moments of the gradients.
 - Often efficient.
 - ADAMax: $v^{[k]} = \tilde{v}^{[k]} = \max(\beta_v v^{[k-1]}, |\nabla G(\mathbf{w}^{[k]})|)$.

$$\mathbf{w}_j^{[k+1]} = \mathbf{w}_j^{[k]} + \frac{\alpha_b^{[k]}}{d_j^{[k]}} m_j^{[k]}$$

AdaDelta

- $m^{[k]} = \nabla G(\mathbf{w}^{[k]})$
- $d_j^{[k]} = \frac{\sqrt{v_j^{[k]}}}{\sqrt{\tilde{v}^{[k-1]}}}$ with $v_j^{[k]} = \rho v_j^{[k-1]} + (1 - \rho)(\nabla G(\mathbf{w}^{[k]})_j)^2$ and
 $\tilde{v}_j^{[k]} = \rho \tilde{v}_j^{[k-1]} + (1 - \rho) \left(\frac{\alpha_j^{[k]}}{d_k^{[k]}} \nabla G(\mathbf{w}^{[k]})_j \right)^2$
- $\alpha_b^{[k]} = 1$

- Only one parameter!
- Based on a dimensional analysis. . . with a strange looking renormalization!

$$\mathbf{w}_j^{[k+1]} = \mathbf{w}_j^{[k]} + \frac{\alpha_b^{[k]}}{d_j^{[k]}} m_j^{[k]}$$

AdaDelta

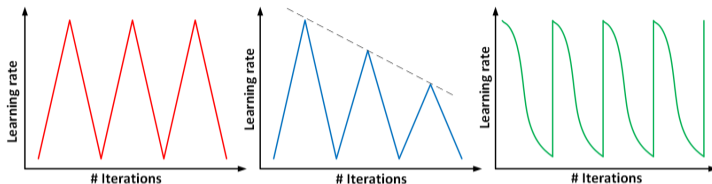
- $m^{[k]} = \nabla G(\mathbf{w}^{[k]})$ and $d_j^{[k]} = 1$

- $\alpha_j^{[0]} = \alpha$ and $\alpha_j^{[k]} = \frac{\bar{\alpha}_j^{[k-1]}}{\sqrt{\rho + (1 - \rho)(\nabla G(\mathbf{w}^{[k]})_j^2 / v_j^{[k-1]})}}$ with

$$v_j^{[k]} = \rho v_j^{[k-1]} + (1 - \rho)(\nabla G(\mathbf{w}^{[k]})_j^2), \text{ and}$$

$$v_j^{[k]} = \rho \tilde{v}_j^{[k-1]} + (1 - \rho) (\alpha_j^{[k]})^2 (\nabla G(\mathbf{w}^{[k]})_j^2) \text{ and } \bar{\alpha}_j^{[k-1]} = \sqrt{\tilde{v}_j^{[k]} / v_j^{[k]}}$$

- Better insight (and more flexible initial step)!
- Quite different from the previous schemes.



Cyclic Learning Rate

- Cycle between small and large learning rate with stopping when the learning rate is small.
- One cycle strategy seems possible!
- Plain SGD with a constant step is also making a comeback!

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 **Optimization: Gradient Descent Algorithms**
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - **Non-Convex Setting**
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- No global convergence result for non convex functions optimization (NP Hard problem)
- Nevertheless (Stochastic) Gradient Descent algorithm can be used.

Theoretical results?

- Typical weak convergence results:
 - convergence to a stationary point
 - convergence to a local minimum. . .
- Often in specific cases.
- In practice: convergence depends a lot on the step size α choice strategy.
- In general, more exploration than in the convex setting.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 **Optimization: Gradient Descent Algorithms**
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - **References**
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References



T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning.
Springer Series in Statistics, 2009



M. Mohri, A. Rostamizadeh, and A. Talwalkar.
Foundations of Machine Learning (2nd ed.)
MIT Press, 2018



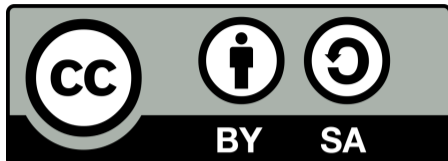
A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



Ch. Giraud.
Introduction to High-Dimensional Statistics.
CRC Press, 2014



S. Bubeck.
Convex Optimization: Algorithms and Complexity.
Now Publisher, 2015



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

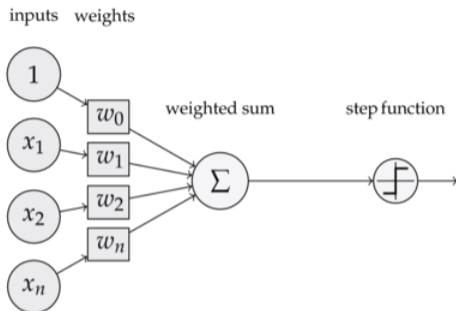
- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.

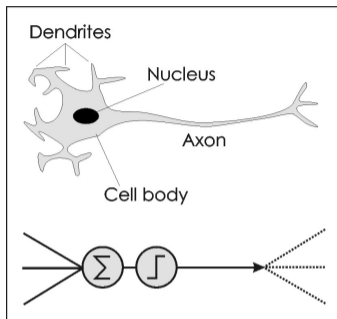
- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 8 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 9 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 8 Unsupervised Learning: Dimension Reduction and Clustering
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 9 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References



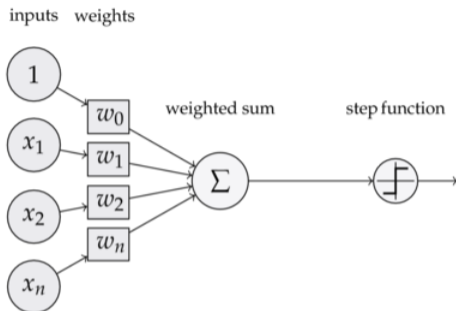
Perceptron (Rosenblatt 1957)

- Inspired from biology.
- Very simple (linear) model!
- Physical implementation and proof of concept.



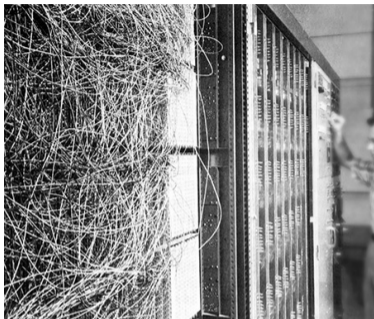
Perceptron (Rosenblatt 1957)

- Inspired from biology.
- Very simple (linear) model!
- Physical implementation and proof of concept.



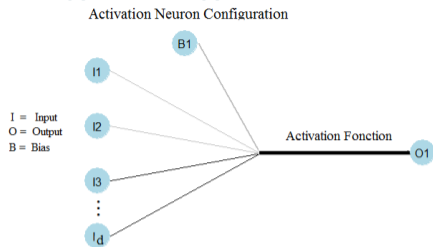
Perceptron (Rosenblatt 1957)

- Inspired from biology.
- Very simple (linear) model!
- Physical implementation and proof of concept.



Perceptron (Rosenblatt 1957)

- Inspired from biology.
- Very simple (linear) model!
- Physical implementation and proof of concept.



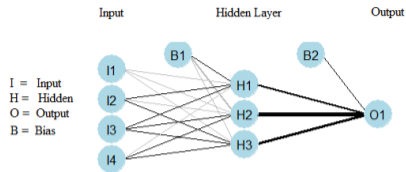
Artificial neuron

- Structure:
 - Mix inputs with a **weighted sum**,
 - Apply a (non linear) **activation function** to this sum,
 - Possibly threshold the result to make a decision.
- Weights learned by minimizing a loss function.

Logistic unit

- Structure:
 - Mix inputs with a **weighted sum**,
 - Apply the **logistic function**
 $\sigma(t) = e^t / (1 + e^t)$,
 - Threshold at 1/2 to make a decision!
- Logistic weights learned by minimizing the -log-likelihood.

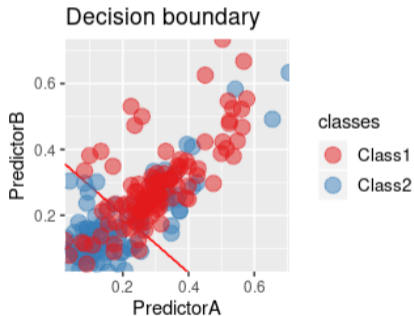
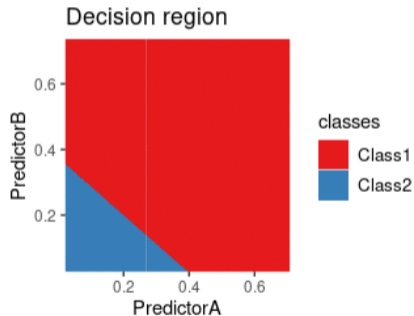
- Equivalent to linear regression when using a linear activation function!

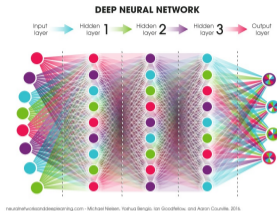


MLP (Rumelhart, McClelland, Hinton - 1986)

- Multilayer Perceptron: cascade of layers of artificial neuron units.
- Optimization through a gradient descent algorithm with a clever implementation (**Backprop**).
- Construction of a function by composing simple units.
- MLP corresponds to a specific direct acyclic graph structure.
- Non convex optimization problem!

Neural Network

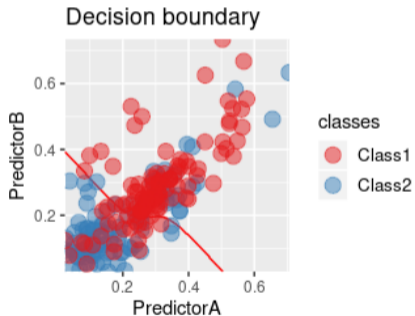
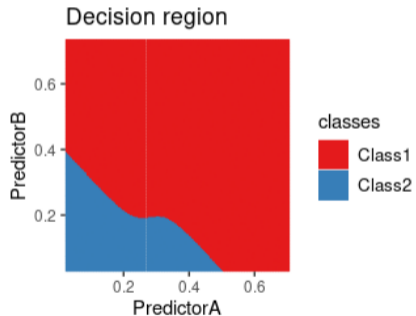


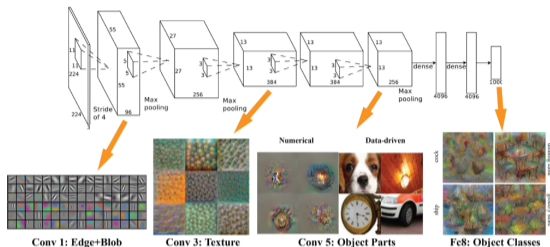


Deep Neural Network structure

- Deep cascade of layers!
- No conceptual novelty. . .
- But a **lot of tricks** allowing to obtain a good solution: clever initialization, better activation function, weight regularization, accelerated stochastic gradient descent, early stopping. . .
- Use of GPU and a lot of data. . .
- Very impressive results!

H2O NN

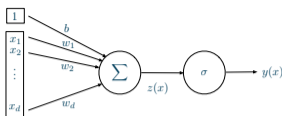




Family of Machine Learning algorithm combining:

- a (deep) multilayered structure,
 - a clever optimization including initialization and regularization.
-
- Examples: Deep NN, AutoEncoder, Recursive NN, GAN, Transformer...
 - Interpretation as a **Representation Learning**.
 - **Transfer learning**: use as initialization a pretrained net.
 - Very efficient and still evolving!

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 ML Methods: Neural Networks and Deep Learning
 - Gradient Descent Step
 - Non-Convex Setting
 - References
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 8 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 9 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 10 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 11 Applications to Text
 - References
- 12 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 13 References



Binary Logistic Regression

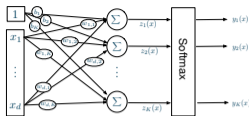
- Conditional probability model:

$$\mathbb{P}(Y = 1|\underline{X}) = \sigma(\mathbf{w}^\top \underline{X} + b) = \frac{e^{\mathbf{w}^\top \underline{X} + b}}{1 + e^{\mathbf{w}^\top \underline{X} + b}}$$

- Model **weights** $\mathbf{w} \in \mathbb{R}^d$ and intercept (or **bias**) $b \in \mathbb{R}$.

Artificial Neuron or Unit

- \underline{X} is the input (single feature vector).
- $z(\underline{X}) = \mathbf{w}^\top \underline{X} + b$ is called **pre-activation**.
- $y(\underline{X}) = \sigma(z(\underline{X}))$ is the output (in $[0, 1]$ in this case)
- \mathbf{w} is **weights** and b is **bias**
- σ is an **activation** function

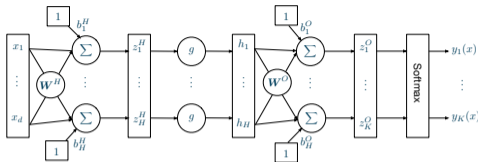


Multinomial Logistic Model

- Cond. prob. model: $\mathbb{P}(Y = k | \underline{X}) = \frac{e^{\mathbf{w}_k^\top \underline{X} + b_k}}{\sum_{k'=1}^K e^{\mathbf{w}_{k'}^\top \underline{X} + b_{k'}}$
- Model weights: k vector $\mathbf{w}_k \in \mathbb{R}^d$ or equivalently a model weight matrix \mathbf{W} with $\mathbf{W}_{\bullet,k} = \mathbf{w}_k$

Softmax Regression (Softmax Layer)

- $z_k(\underline{X}) = \mathbf{W}_{\bullet,k}^\top \underline{X} + b_k$ pre-activations or *logits*
- $y_k(\underline{X}) = \frac{e^{z_k(\underline{X})}}{\sum_{k'=1}^K e^{z_{k'}(\underline{X})}}$ coming out of the softmax activation
- Matricially:
$$y(\underline{X}) = \text{softmax}(z(\underline{X})) = \text{softmax}(\mathbf{W}^\top \underline{X} + b)$$

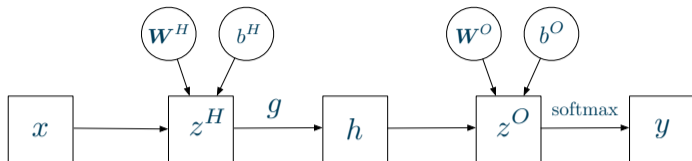


One-Hidden Layer Neural Network - Neural Representation

- Recursive definition:

$$\begin{aligned}y(\underline{X}) &= \text{softmax}(z^{(O)}) = \text{softmax}(\mathbf{W}^{(O)\top} h + b^{(O)}) \\ &= \text{softmax}(\mathbf{W}^{(O)\top} g(z^{(H)}) + b^{(O)}) \\ &= \text{softmax}(\mathbf{W}^{(O)\top} g(\mathbf{W}^{(H)\top} x + b^{(H)}) + b^{(O)})\end{aligned}$$

- g is an **activation function** applied entrywise on z_1^H, \dots, z_H^H
- This neural network has a width- H hidden layer

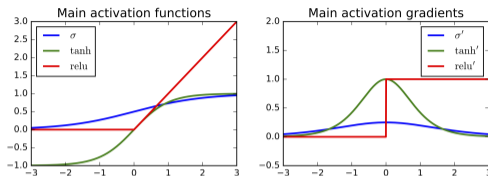


One-Hidden Layer Neural Network - Vector Representation

- Vector representation:

$$\begin{aligned}y(\underline{X}) &= \text{softmax}(z^{(O)}) = \text{softmax}(\mathbf{W}^{(O)\top} h + b^{(O)}) \\ &= \text{softmax}(\mathbf{W}^{(O)\top} g(z^{(H)}) + b^{(O)}) \\ &= \text{softmax}(\mathbf{W}^{(O)\top} g(\mathbf{W}^{(H)\top} \underline{X} + b^{(H)}) + b^{(O)})\end{aligned}$$

- g is an **activation function** applied entrywise
- Equivalent representation focusing on layers.



Entrywise Activation Functions

- Most classical: **sigmoid**, **tanh** and **ReLU** (rectified linear unit)

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \quad \tanh(z) = \frac{e^{2z} - 1}{e^{2z} + 1}, \quad \text{ReLU}(z) = \max(0, z)$$

- Their **derivatives** are given by

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

$$\tanh'(z) = 1 - \tanh(z)^2,$$

$$\text{ReLU}'(z) = \mathbf{1}_{z>0}$$

- Note that $\tanh(z) = 2\sigma(2z) - 1$

ReLU Activation

- ReLU: $z \mapsto \max(0, z)$
- Often the **best choice** for deep neural networks
- Easier to optimize, since their behavior is closer to linear
- Its gradient is not defined at $z = 0$: not a problem since, during training, it's unlikely that any input equals 0
- In contrast to ReLU, *Sigmoid* activations σ and \tanh (not recommended anymore) saturate for large positive (or large negative values).

Identity Activation

- Hidden layer:
 - Amount to use $h = g(\mathbf{V}\mathbf{U}^\top z + b)$ instead of $h = g(\mathbf{W}^\top z + b)$
 - \mathbf{W} is factorized as $\mathbf{U}\mathbf{V}^\top$
 - Reduces the dimension of the model
- Output layer: linear model on the last layer

Softmax

- Softmax is not an element-wise activation:

$$\text{softmax}(z) = \frac{1}{\sum_{k=1}^K e^{z_k}} \begin{bmatrix} e^{z_1} \\ \vdots \\ e^{z_K} \end{bmatrix}$$

- Gradient:

$$\frac{\text{softmax}(z)_k}{\partial x_{k'}} = \begin{cases} \text{softmax}(z)_k \times (1 - \text{softmax}(z)_k) & \text{if } k = k' \\ -\text{softmax}(z)_k \times \text{softmax}(z)_{k'} & \text{otherwise} \end{cases}$$

- Maps $z \in \mathbb{R}^K$ to the space of vector in $[0, 1]^K$ with entries that sum to 1 (**probabilities**)
- The inputs of the softmax are called the **logits** in deep learning

Feed-Forward Neural Network

- Heuristic: It's easy to construct a complex function by composing simple elements in some order.

- **Layerwise structure:**

$$h^{(1)} = g^{(1)}(\mathbf{W}^{(1)\top} x + b^{(1)})$$

$$h^{(2)} = g^{(2)}(\mathbf{W}^{(2)\top} h^{(1)} + b^{(2)})$$

⋮

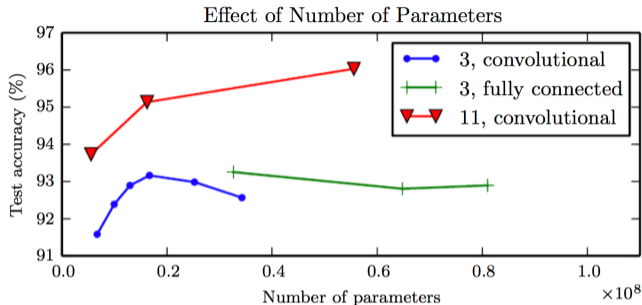
$$h^{(L)} = g^{(L)}(\mathbf{W}^{(L)\top} h^{(L-1)} + b^{(L)})$$

$$y = \text{softmax}(\mathbf{W}^{(O)\top} h^{(L)} + b^{(O)})$$

- Need to choose the width of each layer and the depth of the network.
- The name comes from the fact that information flows through the layers.
- For regression, replace the softmax by identity.

Universal Approximation Theorem (Hornik, 1991)

- A **single hidden layer neural network** with a linear output unit can **approximate** any continuous function **arbitrarily well**, given enough hidden units
- Valid for most activation functions.
- No bounds on the number of required units. . . (Asymptotic flavor)
- A single hidden layer is sufficient but more may require less units?

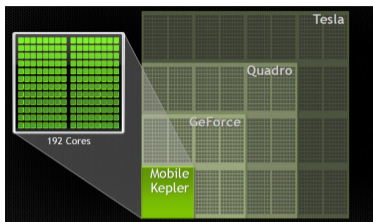


- What is best?
 - One hidden layer with large width?
 - ... or several layers with smaller width ?

A Recipe

- For the same number of parameters, several layers with a smaller width lead to better generalization. (More abstract layers)

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 **ML Methods: Neural Networks and Deep Learning**
 - Gradient Descent Step
 - Non-Convex Setting
 - References
 - Introduction
 - From Logistic Regression to NN
 - **NN Optimization**
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 8 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 9 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 10 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 11 Applications to Text
 - References
- 12 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 13 References



- Use a differentiable loss.

Loss Minimization

- Compute the prediction and the loss (**forward-propagation**)
- Compute the gradients (**back-propagation**)
- Use a Stochastic Gradient Descent algorithm
- In practice, use
 - A nice open source library (tensorflow, keras, pyTorch)
 - A GPU...since training a DNN can benefit from massive parallel computations

Feed-Forward Neural Network

- Network structure:

$$\begin{aligned}h^{(1)} &= g^{(1)}(\mathbf{W}^{(1)\top} \mathbf{x} + b^{(1)}) & h^{(2)} &= g^{(2)}(\mathbf{W}^{(2)\top} h^{(1)} + b^{(2)}) \\ \dots & & h^{(L)} &= g^{(L)}(\mathbf{W}^{(L)\top} h^{(L-1)} + b^{(L)})\end{aligned}$$

$$f(\mathbf{x}, \mathbf{w}) = \text{softmax}(\mathbf{W}^{(o)\top} h^{(L)} + b^{(o)})$$

with $\mathbf{w} = (\mathbf{W}^{(1)}, b^{(1)}, \dots, \mathbf{W}^{(L)}, b^{(L)}, \mathbf{W}^{(o)}, b^{(o)})$

Optimization Formulation

- Minimization of

$$G(\mathbf{w}) = \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(\underline{X}_i, \mathbf{w}))}_{F(\mathbf{w})} + \underbrace{\text{pen}(\mathbf{w})}_{R(\mathbf{w})}$$

with ℓ a loss function and pen a penalization.

- Most classical choices:
 - ℓ is the cross entropy (-log-likelihood)
 - pen is a ridge penalty on the \mathbf{W} .

Optimization Formulation

- Minimization of

$$G(\mathbf{w}) = \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(\underline{X}_i, \mathbf{w}))}_{F(\mathbf{w})} + \underbrace{\text{pen}(\mathbf{w})}_{R(\mathbf{w})}$$

with ℓ a loss function and pen a penalization.

Optimization Algorithm

- Non convex optimization problem.
- Stochastic Gradient Descent algorithm.
- Key issue: computing

$$\nabla F_i(\mathbf{w}) = \nabla \ell(Y_i, f(\underline{X}_i, \mathbf{w}))$$

- Clever algorithm to compute $\nabla F_i(\mathbf{w}) = \nabla \ell(Y_i, f(\underline{X}_i, \mathbf{w}))$.

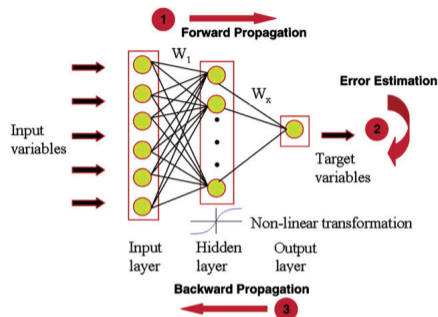
Heuristic

- Key observation: f is obtained through a direct acyclic graph of composition.
- Key lemma (Chain Rule): $\frac{\partial}{\partial \mathbf{w}} f(g(x, \mathbf{w})) = \frac{\partial}{\partial \mathbf{g}}(f(g(x, \mathbf{w}))) \frac{\partial \mathbf{g}}{\partial \mathbf{w}}(x, \mathbf{w})$.
- More generally: $\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$
- Back-propagation: use of this chain rule to compute the derivatives in $\nabla F_i(\mathbf{w})$ starting from the parameters of the last layer and going backward to the ones of the first one.

Set of Partial Derivative Linked $\nabla F_i(\mathbf{w}) = \nabla \ell(Y_i, f(\underline{X}_i, \mathbf{w}))$

- Loss: $\frac{\partial F_i}{\partial f_j}(\mathbf{w}) = \frac{\partial \ell}{\partial f_j}(Y_i, f(\underline{X}_i, \mathbf{w}))$.
- Output layer:
 - Softmax: $\frac{\partial f_j}{\partial z_k^{(o)}}(\mathbf{w})$ is known.
 - Pre-activation at the output layer:
 - Parameters: $\frac{\partial z_j^{(o)}}{\partial \mathbf{w}_{j,k}^{(o)}} = h_k^L(\mathbf{w})$ and $\frac{\partial z_j^{(o)}}{\partial b^{(o)}} = 1$
 - Previous Layer: $\frac{\partial z_j^{(o)}}{\partial h_k^{(l)}} = \mathbf{W}_{j,k}^{(o)}$
- Layer l :
 - Activation: $\frac{\partial h_j^{(l)}}{\partial z_j^{(l)}} = g^{(l)'}(z_j^{(l)})$
 - Pre-activation:
 - Parameters: $\frac{\partial z_j^{(l)}}{\partial \mathbf{w}_{j,k}^{(l)}} = h_k^{l-1}(\mathbf{w})$ and $\frac{\partial z_j^{(l)}}{\partial b^{(l)}} = 1$
 - Previous Layer: $\frac{\partial z_j^{(l)}}{\partial h_k^{(l-1)}} = \mathbf{W}_{j,k}^{(l)}$

- Recurs. comput. with the chain rule \Rightarrow efficient comput. of $\frac{\partial F_i}{\partial \mathbf{w}_{j,k}^{(\cdot)}}$ and $\frac{\partial F_i}{\partial b^{(\cdot)}}$.



Back-Propagation Algorithm

- This set is sufficient to compute $\nabla F_i(\mathbf{w}) = \nabla \ell(Y_i, f(\underline{X}_i, \mathbf{w}))$.
- Those partial derivatives can be computed in a backward-pass. . .
- provided all the intermediate values of h and z have been computed in a forward-pass

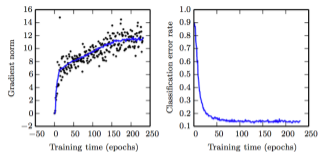
- Nearly all of deep learning is powered by **SGD**.

Stochastic Gradient Descent for DNN

- Regular Stochastic Gradient Algorithm with
 - adaptive step size (ADAM, RMSProp. . .) or fixed size with or without restarts,
 - batch size larger than 1 to use GPU parallelism but not too large.

Initialization

- SGD does not work without a good initialization scheme.
- Key properties:
 - Break symmetry by random initialization around 0.
 - Order of magnitude at layer l given by $1/\sqrt{H_l + H_{l-1}}$ provided the input is normalized.
 - Examples:
 - Normal with $\mathcal{N}(0, 2/(H_l + H_{l-1}))$
 - Uniform with $\mathcal{U}(-\sqrt{6/(H_l + H_{l-1})}, \sqrt{6/(H_l + H_{l-1})})$



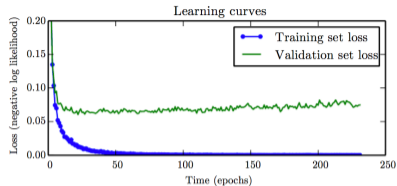
Challenges

- Non convex optimization in very high dimension.
 - SGD can be stuck around critical point and the gradient may not diminishes.
 - Several local minima due to invariance by permutation and scaling.
 - Some local minima may be bad.
-
- **Natural concerns but it seems that for sufficiently large network most of the local minima have a similar low cost value.**
 - Larger networks may be easier to train than smaller ones. . .

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 **ML Methods: Neural Networks and Deep Learning**
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - **NN Regularization**
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 8 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- Applications to Text
- References
- 9 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 10 References

Penalization or Projection

- Regularization on weights and not on bias.
- Mostly ℓ^1 or ℓ^2 .
- Two possibilities:
 - penalization: $\text{pen}(\mathbf{W}) = \lambda \|\mathbf{W}\|$
 - projection: $\|\mathbf{W}\| \leq r$
- Projection seems to work better as it only affects large weights.
- Avoid **dead** units.



Early Stopping

- SGD optimizes Train error but true goal is Test error.
- **Early Stopping:** stop SGD when test error rises.
- In practice, store the parameters each time the test error improves and use the best set of parameters in the end.
- Most widely used of regularization in DL.
- Number of iterations becomes a method parameter.
- For least-squares regression and gradient descent, early stopping \simeq ridge penalization on the weights

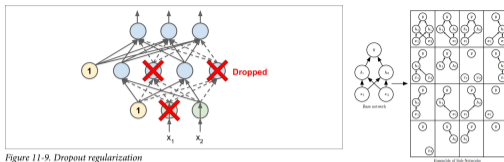


Figure 11-9. Dropout regularization

Dropout

- Idea: At each step, disable temporarily some connections by modifying temporarily the weights.
- Importance sampling idea:

$$\tilde{w}_{j,k}^{(\cdot)} = \begin{cases} 0 & \text{with probability } p \\ \frac{1}{(1-p)} w_{j,k}^{(\cdot)} & \text{with probability } 1 - p \end{cases}$$

- The average value of the weights are preserved.
- Interpretation in term of ensemble methods (bagging)
- No dropout after training!

- Idea: Maintain the scales of the neuron outputs during training.

Batch Normalization Layer

- During training, at each step and for each neuron
 - Compute the mean value on the current batch

$$\mu_k^{(l)} = \frac{1}{m} \sum_{i=1}^m h_k^{(l)}(\underline{X}_i, \mathbf{w})$$

- Compute the variance on the current batch

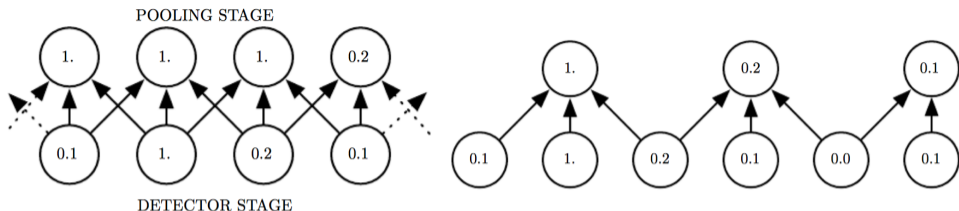
$$\sigma_k^{2,(l)} = \frac{1}{m} \sum_{i=1}^m (h_k^{(l)}(\underline{X}_i, \mathbf{w}) - \mu_k^{(l)})^2$$

- Renormalize the entries $\tilde{h}_k^{(l)}(\underline{X}_i, \mathbf{w}) = \frac{h_k^{(l)}(\underline{X}_i, \mathbf{w}) - \mu_k^{(l)}}{\sigma_k^{(l)}}$

- Train (update) the scale and shift parameters:

$$y_k^{(l+1)} = \gamma_k^{(l+1)} \tilde{h}_k^{(l)}(\underline{X}_i, \mathbf{w}) + \beta_k^{(l+1)}$$

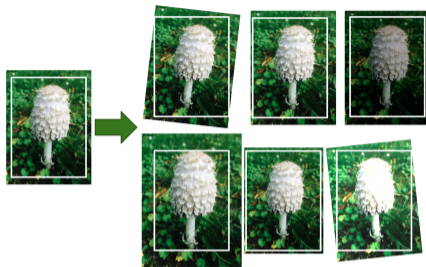
- Differentiability in γ and β : SGD to train them.
- Final freeze of the mean and the variance yield the final NN.
- Helps a lot!



- Idea: replace some outputs by a statistic.

Pooling

- Replace values by a local summary statistic.
 - Most classical summary choice: maximum.
 - (Sub)Gradient can still be computed.
 - Often combine with a subsampling to reduce the network size.
-
- Most useful when output position matters as in images.
 - Very useful to deal with pictures of various sizes.

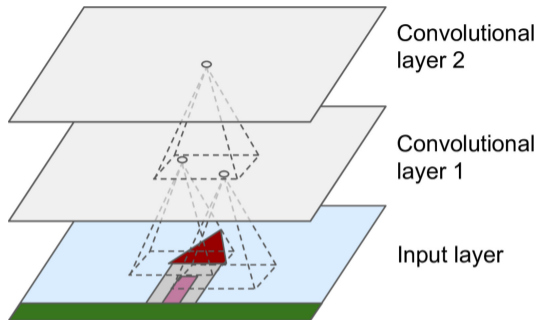


- Idea: Augment artificially the number of samples

Data Augmentation

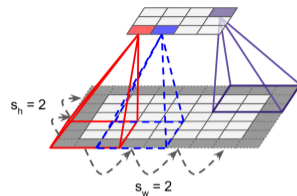
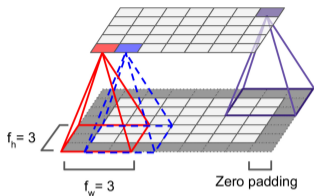
- Augment the dataset by generating new samples from original ones.
- Often mild variations of existing one: shift, rotation, zoom, noise. . .
- Often very efficient if the variations used correspond to specific problem invariance.
- Not so easy to generate completely new samples (GAN?)

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 **ML Methods: Neural Networks and Deep Learning**
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - **Image and CNN**
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 8 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- Applications to Text
- References
- 9 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 10 References



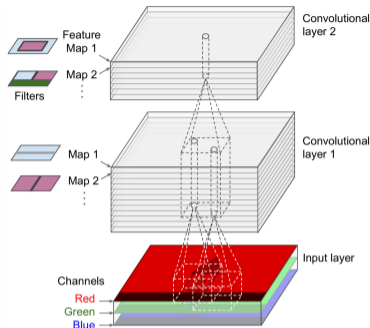
Convolutional Neural Networks (CNN)

- Reduce number of network weights:
 - Use only **local** computation (locality of information),
 - Use the **same weights** everywhere (translation invariance)
- Keep the hierarchical structure.



Convolution

- Filter structure:
 - Output is computed using only neighbors (locality of information)
 - Same computation with respect to the neighbors position everywhere (translation invariance)
- Translation invariance \Leftrightarrow convolution
- Subsampling of resulting output possible (stride) but issues with high frequency part of filters.
- Better subsampling stability after pooling.



Maps and Layers

- Each layer contains a **stack** of outputs (3D tensor).
- Filter : convolution in 2D but free along the stack direction.
- Further reduction of the weight number using 2D spatial convolution followed by arbitrary 1D filter along stack.

PROC. OF THE IEEE, NOVEMBER 1998

7

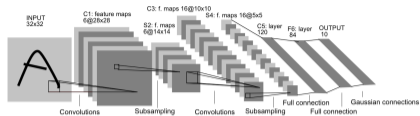
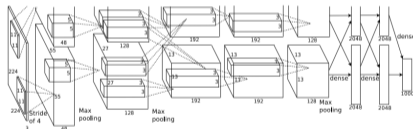


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

LeNet (LeCun, 89)



AlexNet (Krizhevsky et al., 12)

CNN

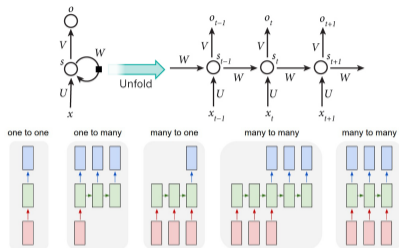
- Stack of convolutional layers followed by a fully connected one.
- Key is to be able to learn in less than a few days!

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 ML Methods: Neural Networks and Deep Learning
 - Gradient Descent Step
 - Non-Convex Setting
 - References
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - **Text, Recurrent Neural Networks and Transformers**
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 8 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 9 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 10 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 11 Applications to Text
 - References
- 12 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 13 References

A recurrent neural network (RNN) is a class of artificial neural network where connections between units form a directed cycle. This creates an internal state of the network which allows it to exhibit dynamic temporal behavior. Unlike feedforward neural networks, RNNs can use their internal memory to process arbitrary sequences of inputs. This makes them applicable to tasks such as unsegmented connected handwriting recognition or speech recognition.

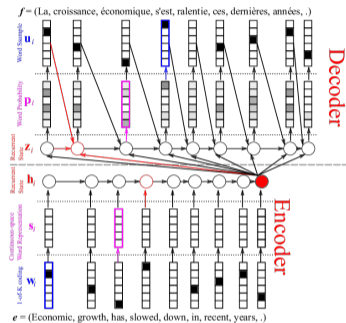
Sequences

- Word = sequence of letters.
- Text = sequence of letters/words.
- Capitalize on this structure.



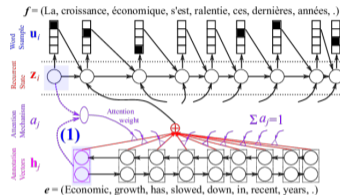
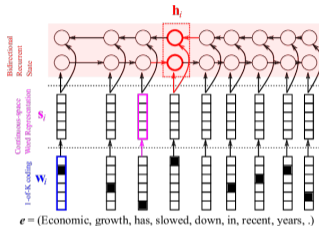
Recurrent Neural Network Unit

- Input seen as a sequence.
- **Simple** computational units with shared weights.
- Information transfer through a context!
- Several architectures!



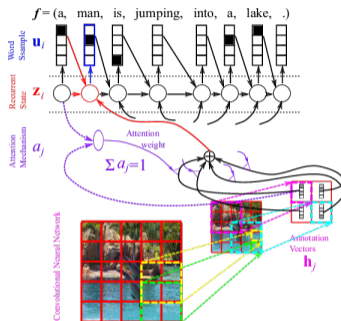
Encoder/Decoder structure

- Word vectors, RNN, stacked structure.



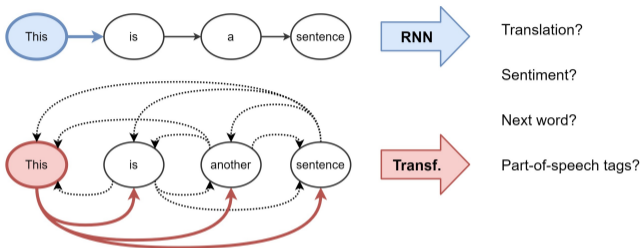
Encoder/Decoder structure

- Much more complex structure: asymmetric, attention order...



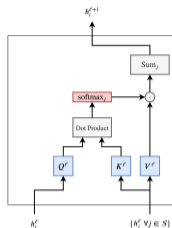
Encoder/Decoder structure

- Much more complex structure: asymmetric, attention order...



Text as Graph

- More than just sequential dependency.
- Each word is related to (all the) other words.
- Graph structure with words and directed relations between words.

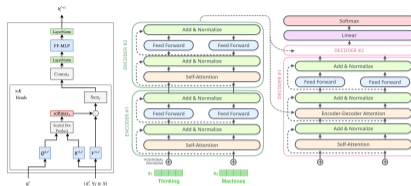


Attention between words

- Words encoded by h_i at layer l .
- Compute individual value for each word: $v_i = V^l h_i$
- Compute combined value for each word: $h'_i = \sum_j w_{i,j} v_j$
- **(Self) Attention:** weight $w_{i,j}$ defined by

$$w_{i,j} = \text{SoftMax} \left(\left\langle Q^l h_i, K^l h_j \right\rangle \right)$$

- $Q^l h_i$ is called a query and $K^l h_j$ a key.



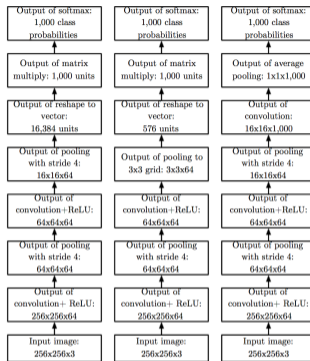
Transformer

- Block combining several attention heads and a classical MLP.

Encoder/Decoder Architecture

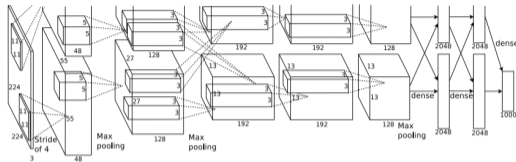
- Combine several transformers and more MLP in a task-adapted architecture.
- End-to-end training is not easy (initialization, optimization...).
- Initial embedding at token level rather than word level to cope with new words!

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 **ML Methods: Neural Networks and Deep Learning**
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - **NN Architecture**
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 8 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- Applications to Text
- References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References



Basic choices

- Nb of layers, size of layers,
- Activation function, pooling. . .
- Optimization algorithm!



¹Inception 5 (GoogLeNet)



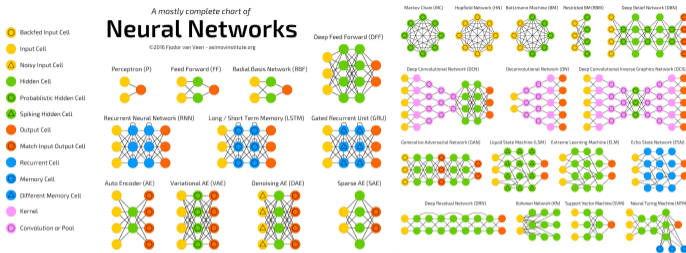
Inception 7a

¹Going Deeper with Convolutions, [C. Szegedy et al, CVPR 2015]



CNN Architecture

- Lot of freedom!



Other Architectures

- Several structures:
 - Different tasks. . .
 - Different inputs.
- **Representation Learning is everywhere**
- Key: Differentiability by composition of basic differentiable units.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 ML Methods: Neural Networks and Deep Learning
 - Gradient Descent Step
 - Non-Convex Setting
 - References
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - **References**
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 8 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 9 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 10 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 11 Applications to Text
 - References
- 12 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 13 References



T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning.
Springer Series in Statistics, 2009



M. Mohri, A. Rostamizadeh, and A. Talwalkar.
Foundations of Machine Learning (2nd ed.)
MIT Press, 2018



A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



Ch. Giraud.
Introduction to High-Dimensional Statistics.
CRC Press, 2014



S. Bubeck.
Convex Optimization: Algorithms and Complexity.
Now Publisher, 2015



I. Goodfellow, Y. Bengio, and A. Courville.
Deep Learning.
MIT Press, 2016



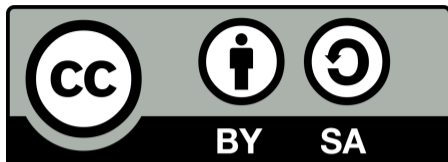
F. Chollet.
Deep Learning with Python (2nd ed.)
Manning, 2021



F. Chollet, T. Kalinowski, and J. J. Allaire.
Deep Learning with R (2nd ed.)
Manning, 2022



A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- Gradient Descent Step
- Non-Convex Setting
- References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5
 - Gradient Descent Step
 - Non-Convex Setting
 - References

5 ML Methods: Neural Networks and Deep Learning

 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 **ML Methods: Trees and Ensemble Methods**
 - **Trees**
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Guess Who?



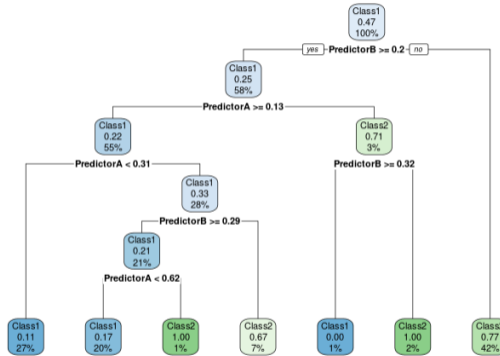
A game of questions

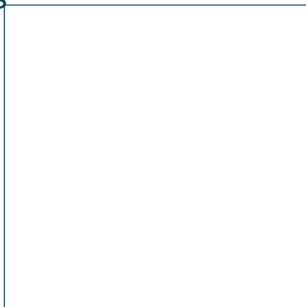
- Game invented in 1979 in the UK.
 - **Goal:** discover the character chosen by your opponent before he discovers yours.
 - **Optimal strategy:** choose at each step the question that splits the remaining characters in two groups with the least possible difference in size.
 - **Information Theory!**
-
- Adaptive construction of a tree of questions!
 - Optimal tree of questions can be constructed without knowing the answers... but during a game only a path of the tree is used...



Tree principle (CART by Breiman (85) / ID3 by Quinlan (86))

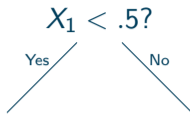
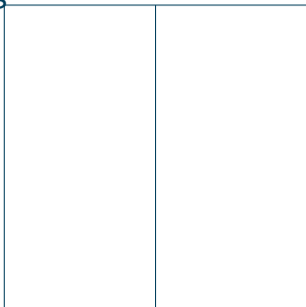
- Construction of a recursive partition through a tree structured set of questions (splits around a given value of a variable)
- For a given partition, probabilistic approach **and** optimization approach yield the same predictor!
- A simple majority vote/averaging in each leaf
- Quality of the prediction depends on the tree (the partition).
- **Intuitively:**
 - small leaves lead to low bias, but large variance
 - large leaves lead to large bias, but low variance. . .
- **Issue:** Minim. of the (penalized) empirical risk is NP hard!
- Practical tree construction are all based on two steps:
 - a top-down step in which branches are created (branching)
 - a bottom-up in which branches are removed (pruning)





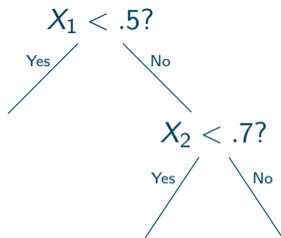
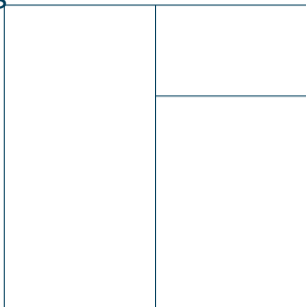
Greedy top-bottom approach

- Start from a single region containing all the data
- Recursively split those regions along a certain variable and a certain value
- **No regret strategy** on the choice of the splits!
- **Heuristic:** choose a split so that the two new regions are as *homogeneous* possible. . .



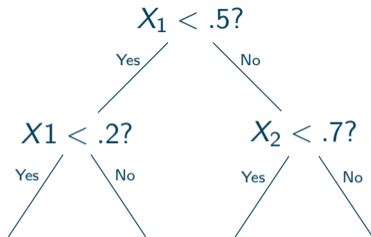
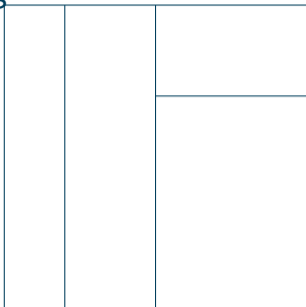
Greedy top-bottom approach

- Start from a single region containing all the data
- Recursively split those regions along a certain variable and a certain value
- **No regret strategy** on the choice of the splits!
- **Heuristic:** choose a split so that the two new regions are as *homogeneous* possible. . .



Greedy top-bottom approach

- Start from a single region containing all the data
- Recursively split those regions along a certain variable and a certain value
- **No regret strategy** on the choice of the splits!
- **Heuristic:** choose a split so that the two new regions are as *homogeneous* possible. . .



Greedy top-bottom approach

- Start from a single region containing all the data
- Recursively split those regions along a certain variable and a certain value
- **No regret strategy** on the choice of the splits!
- **Heuristic:** choose a split so that the two new regions are as *homogeneous* possible. . .

Various definition of *inhomogeneous*

- **CART:** empirical loss based criterion (least squares/prediction error)

$$C(R, \bar{R}) = \sum_{\underline{x}_i \in R} \bar{\ell}(y_i, y(R)) + \sum_{\underline{x}_i \in \bar{R}} \bar{\ell}(y_i, y(\bar{R}))$$

- **CART:** Gini index (Classification)

$$C(R, \bar{R}) = \sum_{\underline{x}_i \in R} p(R)(1 - p(R)) + \sum_{\underline{x}_i \in \bar{R}} p(\bar{R})(1 - p(\bar{R}))$$

- **C4.5:** entropy based criterion (Information Theory)

$$C(R, \bar{R}) = \sum_{\underline{x}_i \in R} H(R) + \sum_{\underline{x}_i \in \bar{R}} H(\bar{R})$$

- CART with Gini is probably the most used technique. . .
- Other criterion based on χ^2 homogeneity or based on different local predictors (generalized linear models. . .)

Choice of the split in a given region

- Compute the criterion for **all features and all possible splitting points** (necessarily among the data values in the region)
- Choose the split **minimizing** the criterion
- **Variations:** split at all categories of a categorical variable using a clever category ordering (ID3), split at a restricted set of points (quantiles or fixed grid)
- **Stopping rules:**
 - when a leaf/region contains less than a prescribed number of observations
 - when the region is sufficiently homogeneous. . .
- May lead to a quite complex tree: over-fitting possible!
- Additional pruning often use.



- **Model selection** within the (rooted) subtrees of previous tree!
- Number of subtrees can be quite large, but the tree structure allows to find the best model efficiently.

Key idea

- The predictor in a leaf depends only on the values in this leaf.
- **Efficient bottom-up (dynamic programming) algorithm** if the criterion used satisfies an additive property

$$C(\mathcal{T}) = \sum_{\mathcal{L} \in \mathcal{T}} c(\mathcal{L})$$

- Example: AIC / CV.

Examples of criterion satisfying this assumptions

- AIC type criterion:

$$\sum_{i=1}^n \bar{\ell}(y_i, f_{\mathcal{L}(\underline{x}_i)}(\underline{x}_i)) + \lambda |\mathcal{T}| = \sum_{\mathcal{L} \in \mathcal{T}} \left(\sum_{\underline{x}_i \in \mathcal{L}} \bar{\ell}(y_i, f_{\mathcal{L}}(\underline{x}_i)) + \lambda \right)$$

- Simple cross-Validation (with (\underline{x}'_i, y'_i) a different dataset):

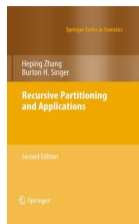
$$\sum_{i=1}^{n'} \bar{\ell}(y'_i, f_{\mathcal{L}}(\underline{x}'_i)) = \sum_{\mathcal{L} \in \mathcal{T}} \left(\sum_{\underline{x}'_i \in \mathcal{L}} \bar{\ell}(y'_i, f_{\mathcal{L}}(\underline{x}'_i)) \right)$$

- Limit over-fitting for a single tree.
- **Rk:** almost never used when combining several trees. . .

- **Key observation:** at a given node, the best subtree is either the current node or the union of the best subtrees of its child.

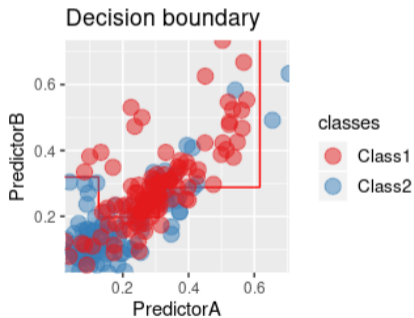
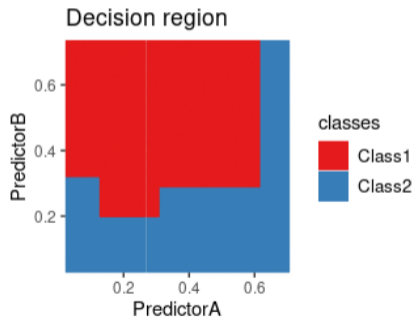
Dynamic programming algorithm

- Compute the individual cost $c(\mathcal{L})$ of each node (including the leaves)
 - Scan all the nodes in reverse order of depth:
 - If the node \mathcal{L} has no child, set its best subtree $\mathcal{T}(\mathcal{L})$ to $\{\mathcal{L}\}$ and its current best cost $c'(\mathcal{L})$ to $c(\mathcal{L})$
 - If the children \mathcal{L}_1 and \mathcal{L}_2 are such that $c'(\mathcal{L}_1) + c'(\mathcal{L}_2) \geq c(\mathcal{L})$, then prune the child by setting $\mathcal{T}(\mathcal{L}) = \{\mathcal{L}\}$ and $c'(\mathcal{L}) = c(\mathcal{L})$
 - Otherwise, set $\mathcal{T}(\mathcal{L}) = \mathcal{T}(\mathcal{L}_1) \cup \mathcal{T}(\mathcal{L}_2)$ and $c'(\mathcal{L}) = c'(\mathcal{L}_1) + c'(\mathcal{L}_2)$
 - The best subtree is the best subtree $\mathcal{T}(\mathcal{R})$ of the root \mathcal{R} .
-
- Optimization cost proportional to the **number of nodes** and not the number of subtrees!



- **Local estimation** of the proportions or of the conditional mean.
- **Recursive Partitioning methods:**
 - Recursive construction of a partition
 - Use of simple local model on each part of the partition
- **Examples:**
 - CART, ID3, C4.5, C5
 - MARS (local linear regression models)
 - Piecewise polynomial model with a dyadic partition. . .
- **Book:** *Recursive Partitioning and Applications* by Zhang and Singer

CART



Pros

- Leads to an easily interpretable model
- Fast computation of the prediction
- Easily deals with categorical features (and missing values)

Cons

- Greedy optimization
- Hard decision boundaries
- Lack of stability

- Lack of robustness for single trees.
- How to combine trees?

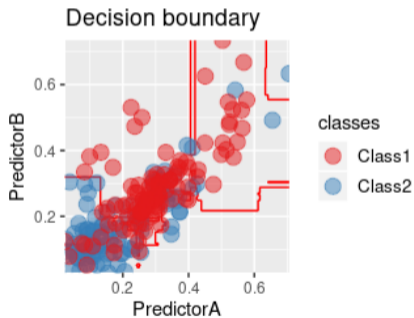
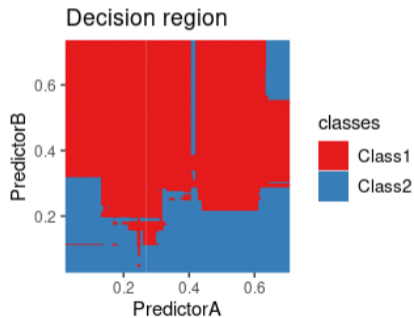
Parallel construction

- Construct several trees from bootstrapped samples and average the responses (**Bagging**)
- Add more randomness in the tree construction (**Random Forests**)

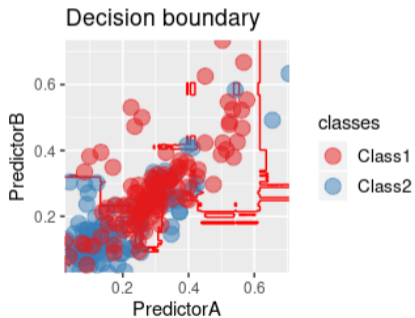
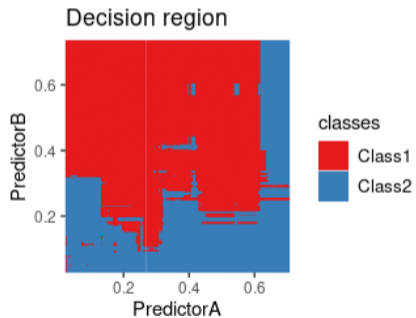
Sequential construction

- Construct a sequence of trees by reweighting sequentially the samples according to their difficulties (**AdaBoost**)
- Reinterpretation as a stagewise additive model (**Boosting**)

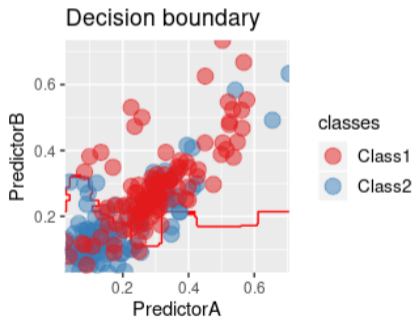
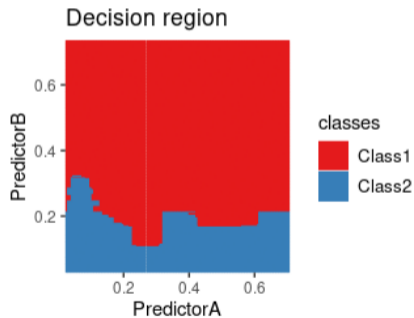
Bagging



Random Forest



AdaBoost



- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- Gradient Descent Step
- Non-Convex Setting
- References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - **Bagging and Random Forests**
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- Applications to Text
- References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Stability through averaging

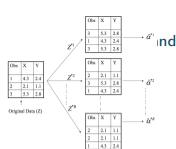
- Very simple idea to obtain a more stable estimator.
- **Vote/average** of B predictors f_1, \dots, f_B obtained with **independent datasets** of size n !

$$f_{\text{agr}} = \text{sign} \left(\frac{1}{B} \sum_{b=1}^B f_b \right) \quad \text{or} \quad f_{\text{agr}} = \frac{1}{B} \sum_{i=1}^B f_b$$

- **Regression:** $\mathbb{E}[f_{\text{agr}}(x)] = \mathbb{E}[f_b(x)]$ and $\text{Var}[f_{\text{agr}}(x)] = \frac{\text{Var}[f_b(x)]}{B}$
 - **Prediction:** slightly more complex analysis
 - Averaging leads to **variance reduction**, i.e. stability!
-
- **Issue:** cost of obtaining B independent datasets of size n !

Bagging and Bootstrap

- Strategy proposed by Breiman in 1994.



Stability through bootstrapping

- Instead of using B independent datasets of size n , draw B datasets from a single one using a **uniform with replacement** scheme (Bootstrap).
- **Rk:** On average, a fraction of $(1 - 1/e) \simeq .63$ examples are unique among each drawn dataset. . .

- The f_b are still identically distributed but **not independent** anymore.
- Price for the non independence: $\mathbb{E}[f_{\text{agr}}(x)] = \mathbb{E}[f_b(x)]$ and

$$\text{Var} [f_{\text{agr}}(x)] = \frac{\text{Var} [f_b(x)]}{B} + \left(1 - \frac{1}{B}\right) \rho(x)$$

with $\rho(x) = \text{Cov} [f_b(x), f_{b'}(x)] \leq \text{Var} [f_b(x)]$ with $b \neq b'$.

- **Bagging:** Bootstrap Aggregation

- Better aggregation scheme exists. . .

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

- Correlation leads to less variance reduction:

$$\text{Var} [f_{\text{agr}}(x)] = \frac{\text{Var} [f_b(x)]}{B} + \left(1 - \frac{1}{B}\right) \rho(x)$$

with $\rho(x) = \text{Cov} [f_b(x), f_{b'}(x)]$ with $b \neq b'$.

- **Idea:** Reduce the correlation by adding more randomness in the predictor.

Randomized Predictors

- Construct predictors that depend on a **randomness source** R that may be chosen independently for all bootstrap samples.
- This **reduces** the correlation between the estimates and thus the **variance**...
- But may **modify heavily the estimates** themselves!
- **Performance gain** not obvious from theory...

- Example of randomized predictors based on trees proposed by Breiman in 2001...

Random Forest

- Draw B resampled datasets from a single one using a uniform with replacement scheme (**Bootstrap**)
- For each resampled dataset, construct a tree using a different **randomly drawn subset of variables** at each split.
- Most important parameter is the **subset size**:
 - if it is too large then we are back to bagging
 - if it is too small the mean of the predictors is probably not a good predictor...
- **Recommendation**:
 - Classification: use a proportion of $1/\sqrt{p}$
 - Regression: use a proportion of $1/3$
- **Sloppier stopping rules** and pruning than in CART...

- Extremely randomized trees!

Extra Trees

- Variation of random forests.
- Instead of trying all possible cuts, try only K cuts at random for each variable.
- No bootstrap in the original article.
- Cuts are defined by a threshold drawn uniformly in the feature range.
- Much faster than the original forest and similar performance.
- Theoretical performance analysis very challenging!

Out Of the Box Estimate

- For each sample x_i , a prediction can be made using only the resampled datasets not containing x_i . . .
- The corresponding empirical prediction error is **not prone to overfitting** but does not correspond to the final estimate. . .
- Good proxy nevertheless.

Forests and Variable Ranking

- **Importance:** Number of use or criterion gain at each split can be used to rank the variables.
- **Permutation tests:** Difference between OOB estimate using the true value of the j th feature and a value drawn a random from the list of possible values.
- Up to OOB error, the permutation technique is not specific to trees.

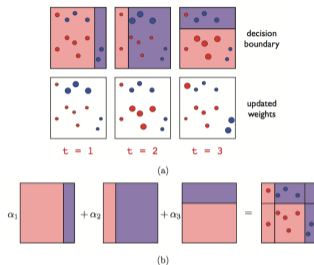
- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5
 - Gradient Descent Step
 - Non-Convex Setting
 - References
 - 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
 - 6 **ML Methods: Trees and Ensemble Methods**
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - **Boosting**
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
 - 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
 - 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
 - 9 References

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

- **Idea:** learn a sequence of predictors trained on weighted dataset with weights depending on the loss so far.

Iterative scheme proposed by Schapire and Freund

- Set $w_{1,i} = 1/n$; $t = 0$ and $f = 0$
- For $t = 1$ to $t = T$
 - $h_t = \operatorname{argmin}_{h \in \mathcal{H}} \sum_{i=1}^n w_{t,i} \ell^{0/1}(y_i, h(x_i))$
 - Set $\epsilon_t = \sum_{i=1}^n w_{t,i} \ell^{0/1}(y_i, h_t(x_i))$ and $\alpha_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$
 - let $w_{t+1,i} = \frac{w_{t,i} e^{-\alpha_t y_i h_t(x_i)}}{Z_{t+1}}$ where Z_{t+1} is a renormalization constant such that $\sum_{i=1}^n w_{t+1,i} = 1$
 - $f = f + \alpha_t h_t$
- Use $f = \sum_{i=1}^T \alpha_t h_t$
- **Intuition:** $w_{t,i}$ measures the difficulty of learning the sample i at step t ...



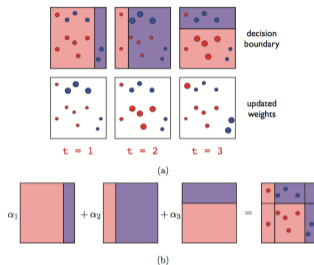
AdaBoost Intuition

- h_t obtained by minimizing a weighted loss

$$h_t = \operatorname{argmin}_{h \in \mathcal{H}} \sum_{i=1}^n w_{t,i} \ell^{0/1}(y_i, h(\underline{x}_i))$$

- Update the current estimate with

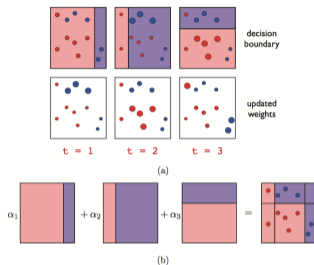
$$f_t = f_{t-1} + \alpha_t h_t$$



AdaBoost Intuition

- Weight $w_{t,i}$ should be large if \underline{x}_i is not well-fitted at step $t - 1$ and small otherwise.
- Use a weight proportional to $e^{-y_i f_{t-1}(\underline{x}_i)}$ so that it can be recursively updated by

$$w_{t+1,i} = w_{t,i} \times \frac{e^{-\alpha_t y_i h_t(\underline{x}_i)}}{Z_t}$$



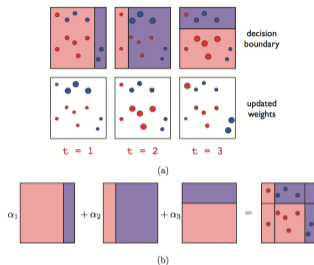
AdaBoost Intuition

- Set α_t such that

$$\sum_{y_i h_t(\underline{x}_i)=1} w_{t+1,i} = \sum_{y_i h_t(\underline{x}_i)=-1} w_{t+1,i}$$

or equivalently

$$\left(\sum_{y_i h_t(\underline{x}_i)=1} w_{t,i} \right) e^{-\alpha_t} = \left(\sum_{y_i h_t(\underline{x}_i)=-1} w_{t,i} \right) e^{\alpha_t}$$



AdaBoost Intuition

- Using

$$\epsilon_t = \sum_{y_i h_t(\underline{x}_i) = -1} w_{t,i}$$

leads to

$$\alpha_t = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t} \quad \text{and} \quad Z_t = 2\sqrt{\epsilon_t(1 - \epsilon_t)}$$

Exponential Stagewise Additive Modeling

- Set $t = 0$ and $f = 0$.
 - For $t = 1$ to T ,
 - $(h_t, \alpha_t) = \operatorname{argmin}_{h, \alpha} \sum_{i=1}^n e^{-y_i(f(\underline{x}_i) + \alpha h(\underline{x}_i))}$
 - $f = f + \alpha_t h_t$
 - Use $f = \sum_{t=1}^T \alpha_t h_t$
-
- **Greedy optimization** of a classifier as a linear combination of T classifier for the **exponential loss**.
 - Those two algorithms are **exactly the same!**

AdaBoost

- Set $t = 0$ and $f = 0$.
 - For $t = 1$ to T ,
 - $(h_t, \alpha_t) = \operatorname{argmin}_{h, \alpha} \sum_{i=1}^n e^{-y_i(f(x_i) + \alpha h(x_i))}$
 - $f = f + \alpha_t h_t$
 - Use $f = \sum_{t=1}^T \alpha_t h_t$
-
- **Greedy iterative scheme** with only two parameters: the class \mathcal{H} of *weak* classifiers and the number of step T .
 - In the literature, one can read that Adaboost does not overfit! This is not true and T should be chosen with care...

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Weak Learner

- Simple predictor belonging to a set \mathcal{H} .
- Easy to learn.
- Need to be only slightly better than a constant predictor.

Weak Learner Examples

- **Decision Tree** with few splits.
- **Stump** decision tree with one split.
- **(Generalized) Linear Regression** with few variables.

Boosting

- Sequential Linear Combination of Weak Learner
- Attempt to minimize a loss.
- Example of ensemble method.
- Link with Generalized Additive Modeling.

- **Greedy optim.** yielding a linear combination of *weak* learners.

Generic Boosting

- Algorithm:
 - Set $t = 0$ and $f = 0$.
 - For $t = 1$ to T ,
 - $(h_t, \alpha_t) = \operatorname{argmin}_{h, \alpha} \sum_{i=1}^n \bar{\ell}(y_i, f(x_i) + \alpha h(x_i))$
 - $f = f + \alpha_t h_t$
 - Use $f = \sum_{t=1}^T \alpha_t h_t$
- AKA as **Forward Stagewise Additive Modeling**
 - AdaBoost with $\bar{\ell}(y, h) = e^{-yh}$
 - LogitBoost with $\bar{\ell}(y, h) = \log_2(1 + e^{-yh})$
 - L_2 Boost with $\bar{\ell}(y, h) = (y - h)^2$ (Matching pursuit)
 - L_1 Boost with $\bar{\ell}(y, h) = |y - h|$
 - HuberBoost with $\bar{\ell}(y, h) = |y - h|^2 \mathbf{1}_{|y-h| < \epsilon} + (2\epsilon|y - h| - \epsilon^2) \mathbf{1}_{|y-h| \geq \epsilon}$
- Simple principle but **no easy numerical scheme** except for AdaBoost and L_2 Boost. . .

- **Issue:** At each boosting step, one need to solve

$$(h_t, \alpha_t) = \operatorname{argmin}_{h, \alpha} \sum_{i=1}^n \bar{\ell}(y_i, f(x_i) + \alpha h(x_i)) = L(y, f + \alpha h)$$

- **Idea:** Replace the function by a **first order approximation**

$$L(y, f + \alpha h) \sim L(y, f) + \alpha \langle \nabla L(y, f), h \rangle$$

Gradient Boosting

- Replace the minimization step by a **gradient descent** step:
 - Choose h_t as the best possible descent direction in \mathcal{H} according to the approximation
 - Choose α_t that minimizes $L(y, f + \alpha h_t)$ (line search)
- **Rk:** Exact gradient direction often not possible!
- Need to find efficiently this best possible direction. . .

- Gradient direction:

$$\begin{aligned}\nabla L(y, f) \quad \text{with} \quad \nabla_i L(y, f) &= \frac{\partial}{\partial f(x_i)} \left(\sum_{i'=1}^n \bar{\ell}(y_{i'}, f(x_{i'})) \right) \\ &= \frac{\partial}{\partial f(x_i)} \bar{\ell}(y_i, f(x_i))\end{aligned}$$

Best Direction within \mathcal{H}

- Direct formulation:

$$h_t \in \operatorname{argmin}_{h \in \mathcal{H}} \frac{\sum_{i=1}^n \nabla_i L(y, f) h(x_i)}{\sqrt{\sum_{i=1}^n |h(x_i)|^2}} \quad \left(= \frac{\langle \nabla L(y, f), h \rangle}{\|h\|} \right)$$

- Equivalent (least-squares) formulation: $h_t = -\beta_t h'_t$ with

$$(\beta_t, h'_t) \in \operatorname{argmin}_{(\beta, h) \in \mathbb{R} \times \mathcal{H}} \sum_{i=1}^n |\nabla_i L(y, f) - \beta h(x_i)|^2 \quad \left(= \|\nabla L - \beta h\|^2 \right)$$

- Choice of the formulation will depend on \mathcal{H} ...

- **Assumptions:**

- h is a binary classifier, $h(x) = \pm 1$ and thus $\|h\|^2 = n$.
- $\bar{\ell}(y, f(x)) = l(yf(x))$ so that $\nabla_i L(y, f) = y_i l'(y_i f(x_i))$.

- Best direction h_t in \mathcal{H} using the first formulation

$$h_t = \operatorname{argmin}_{h \in \mathcal{H}} \sum_i \nabla_i L(y, f) h(x_i)$$

AdaBoost Type Minimization

- Best direction rewriting

$$\begin{aligned} h_t &= \operatorname{argmin}_{h \in \mathcal{H}} \sum_i l'(y_i f(x_i)) y_i h(x_i) \\ &= \operatorname{argmin}_{h \in \mathcal{H}} \sum_i (-l')(y_i f(x_i)) (2\ell^{0/1}(y_i, h(x_i)) - 1) \end{aligned}$$

- **AdaBoost type weighted loss minimization** as soon as $(-l')(y_i f(x_i)) \geq 0$:

$$h_t = \operatorname{argmin}_i \sum_i (-l')(y_i f(x_i)) \ell^{0/1}(y_i, h(x_i))$$

Gradient Boosting

- **(Gradient) AdaBoost:** $\bar{\ell}(y, f) = \exp(-yf)$
 - $l(x) = \exp(-x)$ and thus $(-l')(y_i f(x_i)) = e^{-y_i f(x_i)} \geq 0$
 - h_t is the same as in AdaBoost
 - α_t also... (explicit computation)
- **LogitBoost:** $\bar{\ell}(y, f) = \log_2(1 + e^{-yf})$
 - $l(x) = \log_2(1 + e^{-x})$ and thus $(-l')(y_i f(x_i)) = \frac{e^{-y_i f(x_i)}}{\log(2)(1 + e^{-y_i f(x_i)})} \geq 0$
 - Less weight on misclassified samples than in AdaBoost...
 - No explicit formula for α_t (line search)
 - Different path than with the (non-computable) classical boosting!
- **SoftBoost:** $\bar{\ell}(y, f) = \max(1 - yf, 0)$
 - $l(x) = \max(1 - x, 0)$ and $(-l')(y_i f(x_i)) = \mathbf{1}_{y_i f(x_i) \leq 1} \geq 0$
 - Do not use the samples that are sufficiently well classified!

- Least squares formulation is often preferred when $|h| \neq 1$.

Least Squares Gradient Boosting

- Find $h_t = -\beta_t h'_t$ with

$$(\beta_t, h'_t) \in \operatorname{argmin}_{(\beta, h) \in \mathbb{R} \times \mathcal{H}} \sum_{i=1}^n |\nabla_i L(y, f) - \beta h(x_i)|^2$$

- Classical least squares if \mathcal{H} is a finite dimensional vector space!
- Not a usual least squares in general but a classical regression problem!
- Numerical scheme depends on the loss...

Examples

- **Gradient L_2 Boost:**

- $\ell(y, f) = |y - f|^2$ and $\nabla_i L(y_i, f(x_i)) = -2(y_i - f(x_i))$:

$$(\beta_t, h'_t) \in \underset{(\beta, h) \in \mathbb{R} \times \mathcal{H}}{\operatorname{argmin}} \sum_{i=1}^n |2y_i - 2(f(x_i) - \beta/2h(x_i))|^2$$

- $\alpha_t = -\beta_t/2$
- Equivalent to classical L_2 -Boosting

- **Gradient L_1 Boost:**

- $\ell(y, f) = |y - f|$ and $\nabla_i L(y_i, f(x_i)) = -\operatorname{sign}(y_i - f(x_i))$:

$$(\beta_t, h'_t) \in \underset{(\beta, h) \in \mathbb{R} \times \mathcal{H}}{\operatorname{argmin}} \sum_{i=1}^n |-\operatorname{sign}(y_i - f(x_i)) - \beta h(x_i)|^2$$

- Robust to outliers. . .

- Classical choice for \mathcal{H} : Linear Model in which each h depends on a small subset of variables.

- Least squares formulation can also be used in classification!
- Assumption:
 - $\ell(y, f(x)) = l(yf(x))$ so that $\nabla_i L(y_i, f(x_i)) = y_i l'(y_i f(x_i))$

Least Squares Gradient Boosting for Classifiers

- Least Squares formulation:

$$(\beta_t, h'_t) \in \operatorname{argmin}_{(\beta, h) \in \mathbb{R} \times \mathcal{H}} \sum_{i=1}^n |y_i l'(y_i f(x_i)) - \beta h(x_i)|^2$$

- Equivalent formulation:

$$(\beta_t, h'_t) \in \operatorname{argmin}_{(\beta, h) \in \mathbb{R} \times \mathcal{H}} \sum_{i=1}^n |(-l')(y_i f(x_i)) - (-\beta) y_i h(x_i)|^2$$

- **Intuition:** Modify misclassified examples without modifying too much the well-classified ones. . .

Stochastic Boosting

- **Idea:** change the learning set at each step.
- Two possible reasons:
 - Optimization over all examples too costly
 - Add variability to use an averaged solution
- Two different samplings:
 - Use sub-sampling, if you need to reduce the complexity
 - Use re-sampling, if you add variability...
- Stochastic Gradient name mainly used for the first case...

Second Order Boosting

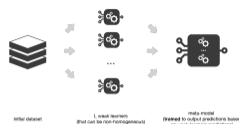
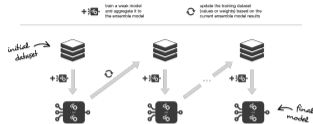
- Replace the first order approximation by a second order one and avoid the line search...

- Very efficient boosting algorithm proposed by Chen and Guestrin in 2014.

eXtreme Gradient Boosting

- Gradient boosting for a (penalized) smooth loss using a second order approximation and the least squares approximation.
 - Reduced stepsize with a shrinkage of the *optimal* parameter.
 - Feature subsampling.
 - Weak learners:
 - Trees: limited depth, penalized size and parameters, fast approximate best split.
 - Linear model: elastic-net penalization.
-
- Excellent baseline!
 - Lightgbm and CatBoost are also excellent similar choices!

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5
 - Gradient Descent Step
 - Non-Convex Setting
 - References
 - 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
 - 6 **ML Methods: Trees and Ensemble Methods**
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - **Ensemble Methods**
 - References
 - 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
 - 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
 - 9 References



Ensemble Methods

- **Averaging:** combine several models by averaging (bagging, random forests, . . .)
 - **Boosting:** construct a sequence of (weak) classifiers (XGBoost, LightGBM, CatBoost)
 - **Stacking:** use the outputs of several models as features (tpot. . .)
-
- Loss of interpretability but gain in performance
 - Beware of overfitting with stacking: the second learning step should be done with fresh data.
 - No end to end optimization as in deep learning!

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5
 - Gradient Descent Step
 - Non-Convex Setting
 - References
 - 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
 - 6 **ML Methods: Trees and Ensemble Methods**
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - **References**
 - 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - 8 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - 9 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - 10 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - 11 Applications to Text
 - References
 - 12 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
 - 13 References



T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning.
Springer Series in Statistics, 2009



M. Mohri, A. Rostamizadeh, and A. Talwalkar.
Foundations of Machine Learning (2nd ed.)
MIT Press, 2018



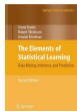
A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



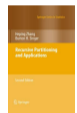
Ch. Giraud.
Introduction to High-Dimensional Statistics.
CRC Press, 2014



S. Bubeck.
Convex Optimization: Algorithms and Complexity.
Now Publisher, 2015



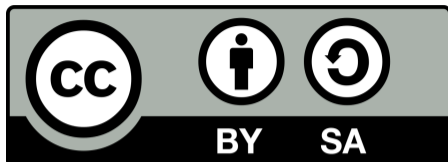
T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning.
Springer Series in Statistics, 2009



H. Zhang and B. Singer.
Recursive Partitioning and Applications.
Springer, 2010



A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (2nd ed.)
O'Reilly, 2019



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

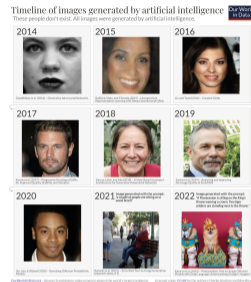
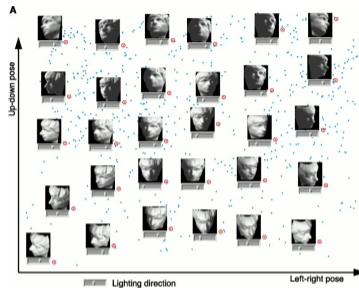
- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References



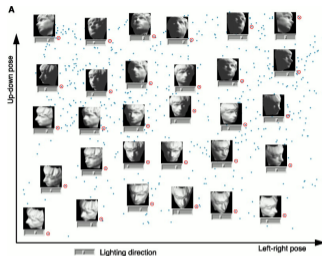
What is possible with data without labels?

- To group them?
- To visualize them in a 2 dimensional space?
- To generate more data?



To group them?

- **Data:** Base of customer data containing their properties and past buying records
- **Goal:** Use the customers *similarities* to find groups.
- **Clustering:** propose an explicit *grouping* of the customers
- **Visualization:** propose a representation of the customers so that the groups are *visible*. (Bonus)



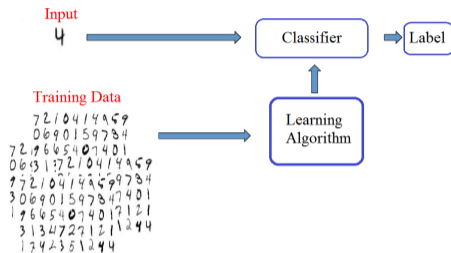
To visualize them?

- **Data:** Images of a single object
- **Goal:** Visualize the *similarities* between images.
- **Visualization:** propose a representation of the images so that similar images are *close*.
- **Clustering:** use this representation to cluster the images. (Bonus)



To generate more data?

- **Data:** Images.
- **Goal:** Generate images similar to the ones in the dataset.
- **Generative Modeling:** propose (and train) a generator.



A definition by Tom Mitchell

A computer program is said to learn from **experience E** with respect to some **class of tasks T** and **performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E.

Experience, Task and Performance measure

- **Training data** : $\mathcal{D} = \{(\underline{X}_1, Y_1), \dots, (\underline{X}_n, Y_n)\}$ (i.i.d. $\sim \mathbb{P}$)
- **Predictor**: $f : \mathcal{X} \rightarrow \mathcal{Y}$ measurable
- **Cost/Loss function**: $\ell(f(\underline{X}), Y)$ measure how well $f(\underline{X})$ predicts Y
- **Risk**:

$$\mathcal{R}(f) = \mathbb{E}[\ell(Y, f(\underline{X}))] = \mathbb{E}_{\underline{X}} \left[\mathbb{E}_{Y|\underline{X}}[\ell(Y, f(\underline{X}))] \right]$$

- Often $\ell(f(\underline{X}), Y) = \|f(\underline{X}) - Y\|^2$ or $\ell(f(\underline{X}), Y) = \mathbf{1}_{Y \neq f(\underline{X})}$

Goal

- Learn a rule to construct a **predictor** $\hat{f} \in \mathcal{F}$ from the training data \mathcal{D}_n s.t. **the risk** $\mathcal{R}(\hat{f})$ is **small on average** or with high probability with respect to \mathcal{D}_n .

Experience, Task and Performance measure

- **Training data** : $\mathcal{D} = \{\underline{X}_1, \dots, \underline{X}_n\}$ (i.i.d. $\sim \mathbb{P}$)
- **Task**: ???
- **Performance measure**: ???

- No obvious task definition!

Tasks for this lecture

- **Dimension reduction**: construct a map of the data in a **low dimensional** space without **distorting** it too much.
- **Clustering (or unsupervised classification)**: construct a **grouping** of the data in **homogeneous** classes.

- **Training data** : $\mathcal{D} = \{\underline{X}_1, \dots, \underline{X}_n\} \in \mathcal{X}^n$ (i.i.d. $\sim \mathbb{P}$)
- Space \mathcal{X} of possibly high dimension.

Dimension Reduction Map

- Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of **smaller dimension**:

$$\Phi : \mathcal{X} \rightarrow \mathcal{X}'$$

$$\underline{X} \mapsto \Phi(\underline{X})$$

- Map can be defined only on the dataset.

Motivations

- Visualization of the data
- Dimension reduction (or embedding) before further processing

- Need to control the **distortion** between \mathcal{D} and $\Phi(\mathcal{D}) = \{\Phi(\underline{X}_1), \dots, \Phi(\underline{X}_n)\}$

Distortion(s)

- Reconstruction error:
 - Construct $\tilde{\Phi}$ from \mathcal{X}' to \mathcal{X}
 - Control the error between \underline{X} and its reconstruction $\tilde{\Phi}(\Phi(\underline{X}))$
 - Relationship preservation:
 - Compute a *relation* \underline{X}_i and \underline{X}_j and a *relation* between $\Phi(\underline{X}_i)$ and $\Phi(\underline{X}_j)$
 - Control the difference between those two *relations*.
-
- Lead to different constructions. . . .

- **Training data** : $\mathcal{D} = \{\underline{X}_1, \dots, \underline{X}_n\} \in \mathcal{X}^n$ (i.i.d. $\sim \mathbb{P}$)
- Latent groups?

Clustering

- Construct a map f from \mathcal{D} to $\{1, \dots, K\}$ where K is a number of classes to be fixed:

$$f : \underline{X}_i \mapsto k_i$$

- Similar to classification except:
 - no ground truth (no given labels)
 - label only elements of the dataset!

Motivations

- Interpretation of the groups
- Use of the groups in further processing

- Need to define the **quality** of the cluster.
- No obvious measure!

Clustering quality

- Inner homogeneity: samples in the same group should be similar.
- Outer inhomogeneity: samples in two different groups should be different.

- Several possible definitions of similar and different.
- Often based on the distance between the samples.
- Example based on the Euclidean distance:
 - Inner homogeneity = intra-class variance,
 - Outer inhomogeneity = inter-class variance.
- **Beware:** choice of the number of clusters K often complex!

Generative Modeling

- **Training data** : $\mathcal{D} = \{(\underline{X}_1, \underline{Y}_1), \dots, (\underline{X}_n, \underline{Y}_n)\} \in (\mathcal{X} \times \mathcal{Y})^n$ (i.i.d. $\sim \mathbb{P}$)
- Same kind of data than for supervised learning if $\mathcal{Y} \neq \emptyset$.

Generative Modeling

- Construct a map G from the product of \mathcal{Y} and a randomness source Ω to \mathcal{X}

$$G : \mathcal{Y} \times \Omega \rightarrow \mathcal{X}$$

$$(Y, \omega) \mapsto X$$

- Unconditional model if $\mathcal{Y} = \emptyset$...

Motivation

- Generate plausible novel conditional samples based on a given dataset.

Sample Quality

- Related to the proximity between the law of $G(Y, \omega)$ and the law of $X|Y$.
- Most classical choice is the Kullback-Leibler divergence.

Ingredients

- Generator $F_\theta(Y, \omega)$ and cond. density prob. $p_\theta(X|Y)$ (Explicit vs implicit link)
- Simple / Complex / Approximate estimation...

Some Possible Choices

	Probabilistic model	Generator	Estimation
Base	Simple (parametric)	Explicit	Simple (ML)
Flow	Image of simple model	Explicit	Simple (ML)
Factorization	Factorization of simple model	Explicit	Simple (ML)
VAE	Simple model with latent var.	Explicit	Approximate (ML)
EBM	Arbitrary	Implicit (MCMC)	Complex (ML/score/discrim.)
Diffusion	Continuous noise	Implicit (MCMC)	Complex (score)
	Discrete Noise with latent var.	Explicit	Approximate (ML)
GAN	Implicit	Explicit	Complex (Discrimination)

- SOTA: Diffusion based approach!

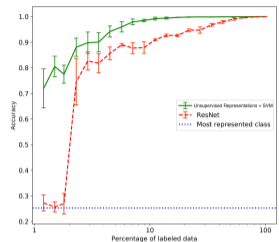
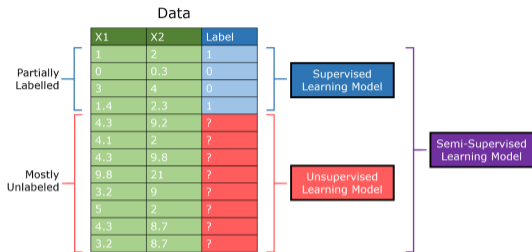
- **General observation:** most data do not have a label !
- **Example:** The number of images on which someone has described the content of the image is a *tiny fraction* of the images online.
- Labeling is very expensive and time consuming
- A lot of information can be extracted from the structure of the data, before seeing any label.

How can we leverage the large quantity of un-labeled data?

- Learn relevant features (= **representations**) in an unsupervised fashion
 - Use those features to solve a supervised task with a fraction of labeled data.
-
- **Semi-supervised framework**
 - ↗ Very useful in practice, for images, time series, text.

Semi-supervised Framework

Unsupervised Learning:
Dimension Reduction and
Clustering



Semu-Supervised Framework

- With representation learned in an unsupervised fashion + a simple linear model, one can achieve the same performance with 10% of data labeled than with a fully annotated dataset.
- Complementary regularization based approaches also exist.

Except for generative modeling, the learner is always right

- A subjective measure of performance
- Subjective choices for the algorithmic constraints (e.g., the type of transformation of the data we allow for low-dimensional representation, type of groups in clustering)
- \Rightarrow Very difficult or impossible to tell which is the *best* method.
- Yet:
 - Extremely important in practice:
 - 90-99% of the data is un-labeled!
 - the tasks themselves are fundamental
 - Huge success in various fields (NLP, images. . .)

Lecture goals for the three main tasks

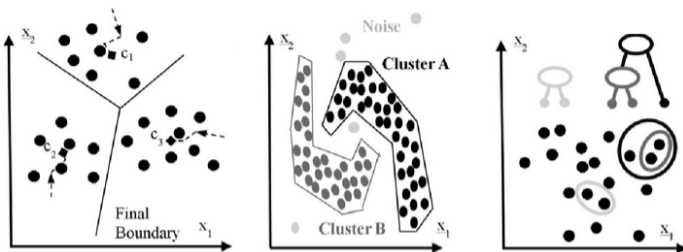
- Discussing possible choices of measures of performance and algorithmic constraints
- Understand the correspondences between those choices and a variety of classical algorithms
- For the simplest algorithms (PCA, k-means), get a precise mathematical understanding of the learning process.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 **Unsupervised Learning: Dimension Reduction and Clustering**
 - Unsupervised Learning?
 - **A First Glimpse**
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

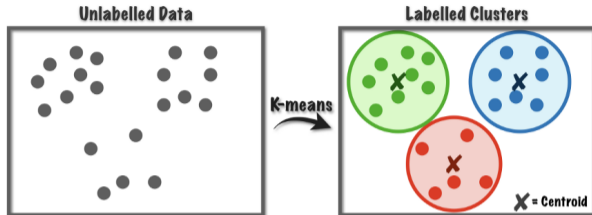
What's a group?



- No simple or unanimous definition!
- Require a notion of similarity/difference. . .

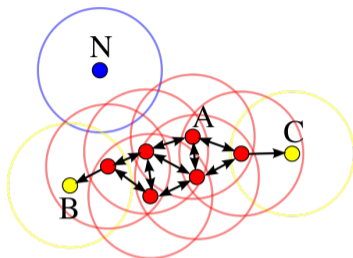
Three main approaches

- A group is a set of samples similar to a prototype.
- A group is a set of samples that can be linked by contiguity.
- A group can be obtained by fusing some smaller groups. . .



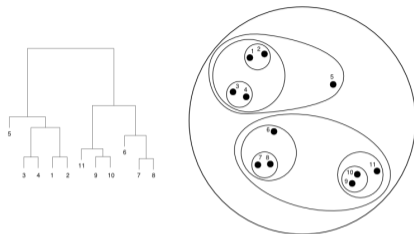
Prototype Approach

- A group is a set of samples similar to a prototype.
- Most classical instance: *k*-means algorithm.
- Principle: alternate prototype choice for the current groups and group update based on those prototypes.
- Number of groups fixed at the beginning
- No need to compare the samples between them!



Contiguity Approach

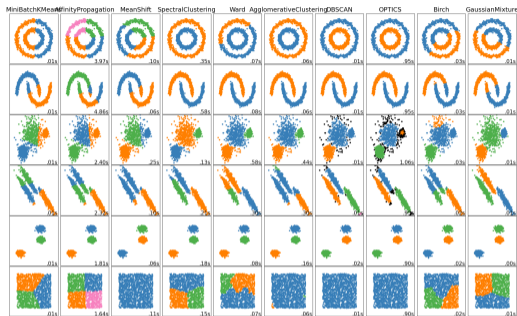
- A group is the set of samples that can be linked by contiguity.
- Most classical instance: DBScan
- Principle: group samples by contiguity if possible (proximity and density)
- Some samples may remain isolated.
- Number of groups controlled by the scale parameter.



Agglomerative Approach

- A group can be obtained by fusing some smaller groups. . .
- Hierarchical clustering principle: sequential merging of groups according to a *best merge* criterion
- Numerous variations on the merging criterion. . .
- Number of groups chosen afterward.

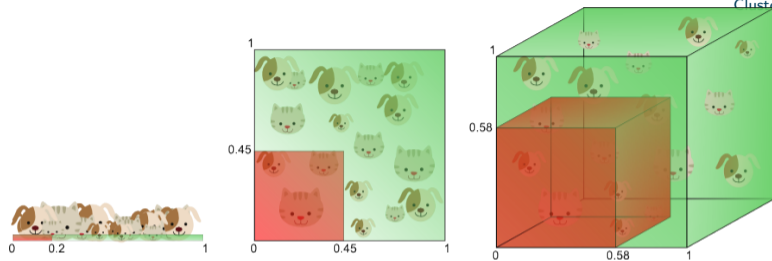
Choice of the method and of the number of groups



No method or number of groups is better than the others...

- Criterion not necessarily explicit!
- No cross validation possible
- Choice of the number of groups: a priori, heuristic, *based on the final usage...*

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View



- **DISCLAIMER: Even if they are used everywhere beware of the usual distances in high dimension!**

Dimensionality Curse

- Previous approaches based on distances.
- Surprising behavior in high dimension: everything is ((often) as) far away.
- Beware of categories. . .

- **DISCLAIMER: Even if they are used everywhere beware of the usual distances in high dimension!**

High Dimensional Geometry Curse

- Folks theorem: In high dimension, everyone is alone.
- Theorem: If $\underline{X}_1, \dots, \underline{X}_n$ in the hypercube of dimension d such that their coordinates are i.i.d then

$$d^{-1/p} \left(\max \|\underline{X}_i - \underline{X}_j\|_p - \min \|\underline{X}_i - \underline{X}_j\|_p \right) = 0 + O_P \left(\sqrt{\frac{\log n}{d}} \right)$$
$$\frac{\min \|\underline{X}_i - \underline{X}_j\|_p}{\max \|\underline{X}_i - \underline{X}_j\|_p} = 1 + O_P \left(\sqrt{\frac{\log n}{d}} \right).$$

- When d is large, all the points are almost equidistant. . .
- Nearest neighbors are meaningless!

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

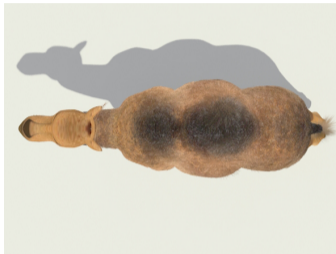
Visualization and Dimension Reduction

- How to view a dataset in high dimension !
- High dimension: dimension larger than 2!
- **Projection** onto a 2D space.



Visualization and Dimension Reduction

- How to view a dataset in high dimension !
- High dimension: dimension larger than 2!
- **Projection** onto a 2D space.



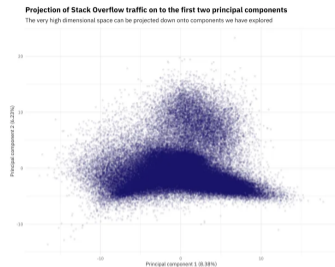
Visualization and Dimension Reduction

- How to view a dataset in high dimension !
- High dimension: dimension larger than 2!
- **Projection** onto a 2D space.



Visualization and Dimension Reduction

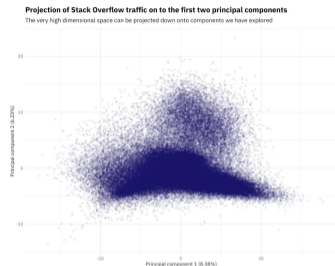
- How to view a dataset in high dimension !
- High dimension: dimension larger than 2!
- **Projection** onto a 2D space.



- Simple formula: $\tilde{X} = P(X - m)$

How to chose P ?

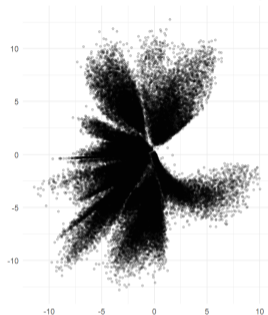
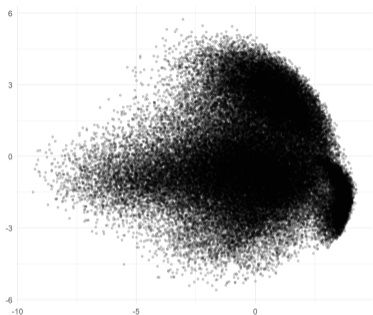
- Maximising the dispersion of the points?
- Allowing to well reconstruct X from \tilde{X} ?
- Preserving the relationship between the X through those between the \tilde{X} ?



- Simple formula: $\tilde{X} = P(X - m)$

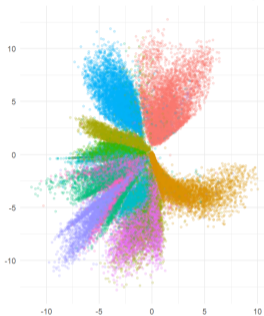
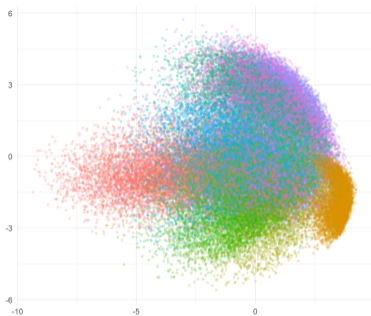
How to chose P ?

- Maximising the dispersion of the points?
- Allowing to well reconstruct X from \tilde{X} ?
- Preserving the relationship between the X through those between the \tilde{X} ?
- The 3 approaches yield the same solution!



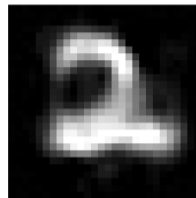
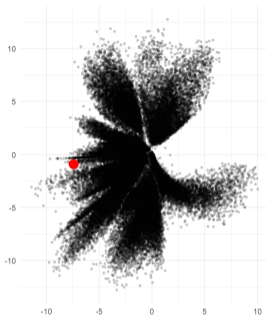
Reconstruction Approaches

- Learn a formula to encode and one formula to decode.
- Auto-encoder structure
- Yields a formula for new points.



Reconstruction Approaches

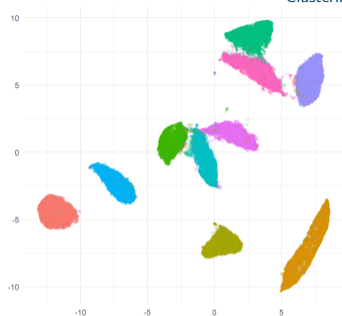
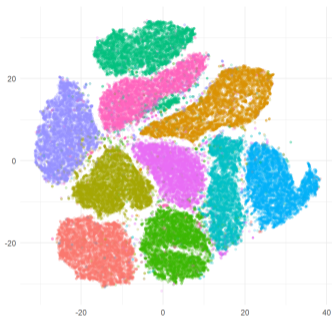
- Learn a formula to encode and one formula to decode.
- Auto-encoder structure
- Yields a formula for new points.



Reconstruction Approaches

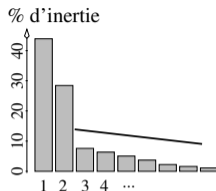
- Learn a formula to encode and one formula to decode.
- Auto-encoder structure
- Yields a formula for new points.

Relationship Preservation Approaches



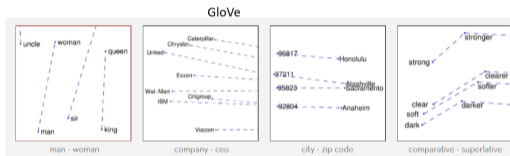
Relationship Preservation Approaches

- Based on the definition of the relationship notion (in both worlds).
- Huge flexibility
- Not always yields a formula for new points.



No Better Choice?

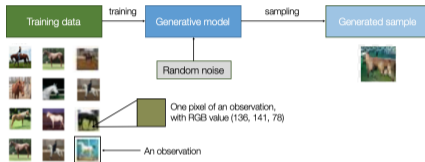
- Different criterion for different methods: impossible to use cross-validation.
 - The larger the dimension the easier is to be faithful!
 - In visualization, dimension 2 is the only choice.
 - Heuristic criterion for the dimension choice: elbow criterion (no more gain), stability...
-
- Dimension Reduction is rarely used standalone but rather as a step in a predictive/prescriptive method.
 - The dimension becomes an hyper-parameter of this method.



Representation Learning

- How to transform arbitrary objects into numerical vectors?
- Objects: Categorical variables, Words, Images/Sounds. . .
- The two previous dimension reduction approaches can be used (given possibly a first simple high dimensional representation)

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View



Timeline of images generated by artificial intelligence

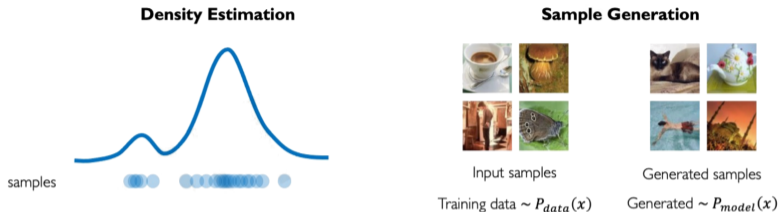
These people don't exist. All images were generated by artificial intelligence.

Our World
in Data



Generative Modeling

- Generate new samples similar to the ones in an original dataset.
- Generation may be conditioned by an input.
- Key for image generation. . . and chatbot!

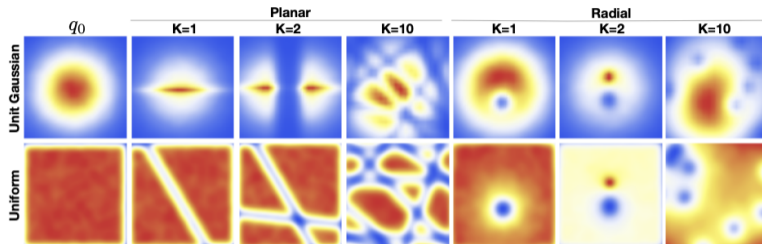


How can we learn $P_{model}(x)$ similar to $P_{data}(x)$?

- **Heuristic:** If we can estimate the (conditional) law \mathbb{P} of the data and can simulate it, we can obtain new samples similar to the input ones.

Estimation and Simulation

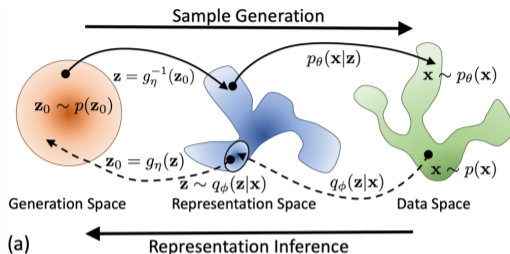
- How to estimate the density?
- How to simulate the estimate density?
- Other possibilities?



Parametric Model, Image and Factorization

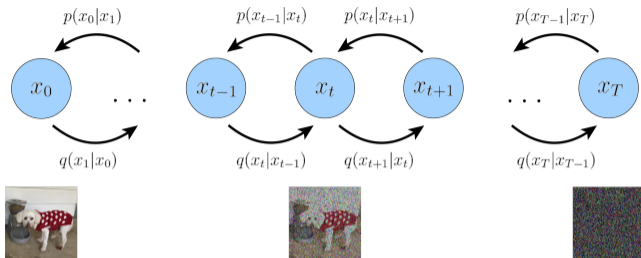
- Use
 - a simple parametric model, . . .
 - or the image of a parametric model (flow), . . .
 - or a factorization of a parametric model (recurrent model)as they are *simple* to estimate and to simulate.

- Estimation by Maximum Likelihood principle.
- Recurrent models are used in Large Language Models!



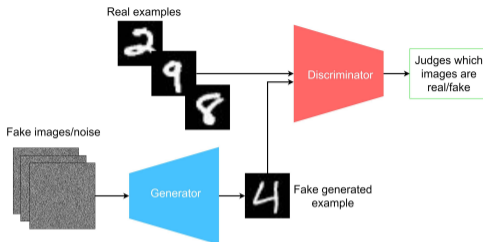
Latent Variable

- Introduce a latent variable Z from which X is easy to sample.
- Estimation based on approximate Maximum Likelihood (VAE/ELBO)
- The latent variable Z often lives in a smaller dimensional space.
- It can be generated by a simple method (or a more complex one...).



Monte Carlo Markov Chain

- Rely on much more complex probability model. . .
- which can only be simulated numerically.
- Often combined with noise injection to stabilize the numerical scheme (Diffusion).
- Much more expensive to simulate than with Latent Variable approaches.



Generative Adversarial Network

- Bypass the density estimation problem, by transforming the problem into a competition between the generator and a discriminator.
- The better the generator, the harder it is for the generator to distinguish true samples from synthetic ones.
- No explicit density!
- Fast simulator but unstable training. . .

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- Gradient Descent Step
- Non-Convex Setting
- References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 **Unsupervised Learning: Dimension Reduction and Clustering**
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- **Training data** : $\mathcal{D} = \{\underline{X}_1, \dots, \underline{X}_n\} \in \mathcal{X}^n$ (i.i.d. $\sim \mathbb{P}$)
- Space \mathcal{X} of possibly high dimension.

Dimension Reduction Map

- Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of **smaller dimension**:

$$\Phi : \mathcal{X} \rightarrow \mathcal{X}'$$

$$\underline{X} \mapsto \Phi(\underline{X})$$

Criterion

- Reconstruction error
- Relationship preservation

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

A Projection Based Approach

- Observations: $\underline{X}_1, \dots, \underline{X}_n \in \mathbf{R}^d$
- Simplified version: $\Phi(\underline{X}_1), \dots, \Phi(\underline{X}_n) \in \mathbf{R}^d$ with Φ an affine projection preserving the mean $\Phi(\underline{X}) = P(\underline{X} - m) + m$ with $P^\top = P = P^2$ and $m = \frac{1}{n} \sum_i \underline{X}_i$.

How to choose P ?

- **Inertia criterion:**

$$\max_P \sum_{i,j} \|\Phi(\underline{X}_i) - \Phi(\underline{X}_j)\|^2?$$

- **Reconstruction criterion:**

$$\min_P \sum_i \|\underline{X}_i - \Phi(\underline{X}_i)\|^2?$$

- **Relationship criterion:**

$$\min_P \sum_{i,j} |(\underline{X}_i - m)^\top (\underline{X}_j - m) - (\Phi(\underline{X}_i) - m)^\top (\Phi(\underline{X}_j) - m)|^2?$$

- **Rk:** Best solution is $P = I$! Need to reduce the rank of the projection to $d' < d \dots$

- **Heuristic:** a good representation is such that the projected points are far apart.

Two views on inertia

- Inertia:

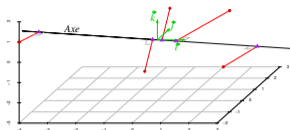
$$I = \frac{1}{2n^2} \sum_{i,j} \|\underline{X}_i - \underline{X}_j\|^2 = \frac{1}{n} \sum_{i=1}^n \|\underline{X}_i - m\|^2$$

- 2 times the mean squared distance to the mean = Mean squared distance between individual

Inertia criterion (Principal Component Analysis)

- Criterion: $\max_P \sum_{i,j} \frac{1}{2n^2} \|P\underline{X}_i - P\underline{X}_j\|^2 = \max_P \frac{1}{n} \sum_i \|P\underline{X}_i - m\|^2$

- **Solution:** Choose P as a projection matrix on the space spanned by the d' first eigenvectors of $\Sigma = \frac{1}{n} \sum_i (\underline{X}_i - m)(\underline{X}_i - m)^\top$



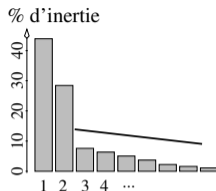
- $\tilde{X} = m + a^\top (X - m)a$ with $\|a\| = 1$
- Inertia: $\frac{1}{n} \sum_{i=1}^n a^\top (X_i - m)(X_i - m)^\top a$

Principal Component Analysis: optimization of the projection

- Maximization of $\tilde{l} = \frac{1}{n} \sum_{i=1}^n a^\top (X_i - m)(X_i - m)^\top a = a^\top \Sigma a$ with

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (X_i - m)(X_i - m)^\top \text{ the empirical covariance matrix.}$$

- Explicit optimal choice given by the eigenvector of the largest eigenvalue of Σ .



Principal Component Analysis : sequential optimization of the projection

- Explicit optimal solution obtain by the projection on the eigenvectors of the largest eigenvalues of Σ .
- Projected inertia given by the sum of those eigenvalues.
- Often fast decay of the eigenvalues: some dimensions are much more important than others.
- Not exactly the curse of dimensionality setting. . .
- Yet a lot of *small* dimension can drive the distance!

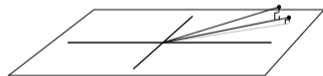
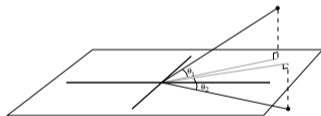
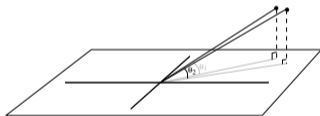
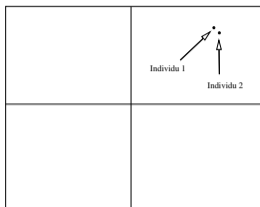
- **Heuristic:** a good representation is such that the projected points are close to the original ones.

Reconstruction Criterion

- Criterion: $\min_P \sum_i \frac{1}{n} \|\underline{X}_i - (P(\underline{X}_i - m) + m)\|^2 = \min_P \frac{1}{n} \sum_i \|(I - P)(\underline{X}_i - m)\|^2$
- **Solution:** Choose P as a projection matrix on the space spanned by the d' first eigenvectors of $\Sigma = \frac{1}{n} \sum_i (\underline{X}_i - m)(\underline{X}_i - m)^\top$

- Same solution with a different heuristic!
- Proof (Pythagora):

$$\sum_i \|\underline{X}_i - m\|^2 = \sum_i \left(\|P(\underline{X}_i - m)\|^2 + \|(I - P)(\underline{X}_i - m)\|^2 \right)$$



Close projection doesn't mean close individuals!

- Same projections but different situations.
- Quality of the reconstruction measured by the angle with the projection space!

- **Heuristic:** a good representation is such that the projected points scalar products are similar to the original ones.

Relationship Criterion (Multi Dimensional Scaling)

- Criterion: $\min_P \sum_{i,j} |(\underline{X}_i - m)^\top (\underline{X}_j - m) - (\Phi(\underline{X}_i) - m)^\top (\Phi(\underline{X}_j) - m)|^2$
- **Solution:** Choose P as a projection matrix on the space spanned by the d' first eigenvectors of $\Sigma = \frac{1}{n} \sum_i (\underline{X}_i - m)(\underline{X}_i - m)^\top$
- Same solution with a different heuristic!
- Much more involved justification!

- PCA model: $\underline{X} - m \simeq P(\underline{X} - m)$
- **Prop:** $P = VV^T$ with V an orthormal family in dimension d' of size d' .
- PCA model with V : $\underline{X} - m \simeq VV^T(\underline{X} - m)$ where $\tilde{\underline{X}} = V^T(\underline{X} - m) \in \mathbb{R}^{d'}$
- Row vector rewriting: $\underline{X}^T - m^T \simeq \tilde{\underline{X}}^T V^T$

Matrix Rewriting and Low Rank Factorization

- Matrix rewriting

$$\begin{array}{ccc}
 \begin{array}{|c|} \hline \underline{X}_1^T - m^T \\ \hline \vdots \\ \hline \underline{X}_n^T - m^T \\ \hline \end{array} & \simeq & \begin{array}{|c|} \hline \tilde{\underline{X}}_1^T \\ \hline \vdots \\ \hline \tilde{\underline{X}}_n^T \\ \hline \end{array} \begin{array}{|c|} \hline \mathbf{V}^T \\ \hline \end{array} \\
 (n \times d) & & (n \times d') \quad (d' \times d)
 \end{array}$$

- Low rank matrix factorization! (Truncated SVD solution...)

SVD Decomposition

- Any matrix $n \times d$ matrix A can be decomposed as

$$\begin{array}{c} \boxed{\mathbf{A}} \\ (n \times d) \end{array} = \begin{array}{c} \boxed{\mathbf{U}} \\ (n \times n) \end{array} \begin{array}{c} \boxed{\mathbf{D}} \\ (n \times d) \end{array} \begin{array}{c} \boxed{\mathbf{W}^T} \\ (d \times d) \end{array}$$

with U and W two orthonormal matrices and D a *diagonal* matrix with decreasing values.

Low Rank Approximation

- The best low rank approximation or rank r is obtained by restriction of the matrices to the first r dimensions:

$$\begin{array}{c}
 \boxed{\mathbf{A}} \\
 (n \times d)
 \end{array}
 \approx
 \begin{array}{c}
 \boxed{\mathbf{U}_r} \\
 (n \times r)
 \end{array}
 \begin{array}{c}
 \boxed{D_{r,r}} \\
 (r \times r)
 \end{array}
 \begin{array}{c}
 \boxed{\mathbf{W}_r^\top} \\
 (r \times d)
 \end{array}$$

for both the operator norm and the Frobenius norm!

- PCA: Low rank approximation with Frobenius norm, $d' = r$ and

$$\begin{pmatrix} \underline{X}_1^\top - m^\top \\ \vdots \\ \underline{X}_n^\top - m^\top \end{pmatrix} \leftrightarrow A, \quad \begin{pmatrix} \tilde{X}_1^\top \\ \vdots \\ \tilde{X}_n^\top \end{pmatrix} \leftrightarrow \mathbf{U}_r D_{r,r}, \quad \mathbf{V}^\top \leftrightarrow \mathbf{W}_r^\top$$

SVD Decompositions

- Recentered data:

$$\mathbf{R} = \begin{pmatrix} \underline{X}_1^\top - m^\top \\ \vdots \\ \underline{X}_n^\top - m^\top \end{pmatrix} = \mathbf{U}\mathbf{D}\mathbf{W}^\top$$

- Covariance matrix:

$$\Sigma = \mathbf{R}^\top \mathbf{R} = \mathbf{W}\mathbf{D}^\top \mathbf{D}\mathbf{W}$$

with $\mathbf{D}^\top \mathbf{D}$ diagonal.

- Gram matrix (matrix of scalar products):

$$\mathbf{G} = \mathbf{R}\mathbf{R}^\top = \mathbf{U}\mathbf{D}\mathbf{D}^\top \mathbf{U}$$

with $\mathbf{D}\mathbf{D}^\top$ diagonal.

- Those are the same \mathbf{U} , \mathbf{W} and \mathbf{D} , hence the link between all the approaches.

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Goal

- Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of **smaller dimension**:

$$\Phi : \mathcal{X} \rightarrow \mathcal{X}'$$

$$\underline{X} \mapsto \Phi(\underline{X})$$

- Construct $\tilde{\Phi}$ from \mathcal{X}' to \mathcal{X}
- Control the error between \underline{X} and its reconstruction $\tilde{\Phi}(\Phi(\underline{X}))$
- Canonical example for $\underline{X} \in \mathbb{R}^d$: find Φ and $\tilde{\Phi}$ in a parametric family that minimize

$$\frac{1}{n} \sum_{i=1}^n \|\underline{X}_i - \tilde{\Phi}(\Phi(\underline{X}_i))\|^2$$

- $\mathcal{X} \in \mathbb{R}^d$ and $\mathcal{X}' = \mathbb{R}^{d'}$
- Affine model $\underline{X} \sim m + \sum_{l=1}^{d'} \underline{X}'^{(l)} V^{(l)}$ with $(V^{(l)})$ an orthonormal family.
- Equivalent to:

$$\Phi(\underline{X}) = V^\top (\underline{X} - m) \quad \text{and} \quad \tilde{\Phi}(\underline{X}') = m + V \underline{X}'$$

- Reconstruction error criterion:

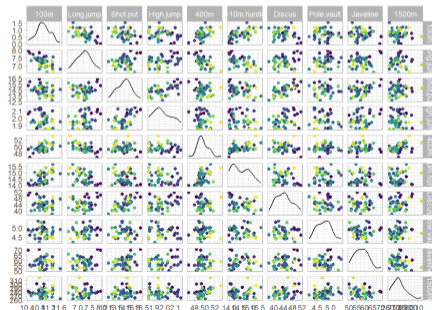
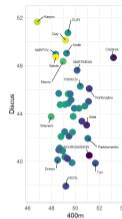
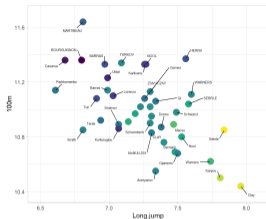
$$\frac{1}{n} \sum_{i=1}^n \|\underline{X}_i - (m + VV^\top (\underline{X}_i - m))\|^2$$

- **Explicit solution:** m is the empirical mean and V is any orthonormal basis of the space spanned by the d' first eigenvectors (the one with largest eigenvalues) of the empirical covariance matrix $\frac{1}{n} \sum_{i=1}^n (\underline{X}_i - m)(\underline{X}_i - m)^\top$.

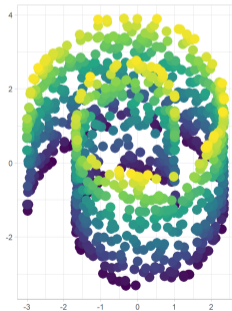
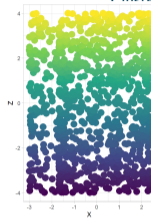
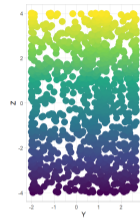
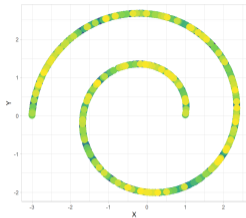
PCA Algorithm

- Compute the empirical mean $m = \frac{1}{n} \sum_{i=1}^n \underline{X}_i$
 - Compute the empirical covariance matrix $\frac{1}{n} \sum_{i=1}^n (\underline{X}_i - m)(\underline{X}_i - m)^\top$.
 - Compute the d' first eigenvectors of this matrix: $V^{(1)}, \dots, V^{(d')}$
 - Set $\Phi(\underline{X}) = V^\top (\underline{X} - m)$
-
- Complexity: $O(n(d + d^2) + d'd^2)$
 - Interpretation:
 - $\Phi(\underline{X}) = V^\top (\underline{X} - m)$: coordinates in the restricted space.
 - $V^{(i)}$: influence of each original coordinates in the i th new one.
 - **Scaling:** This method is not invariant to a scaling of the variables! It is custom to normalize the variables (at least within groups) before applying PCA.

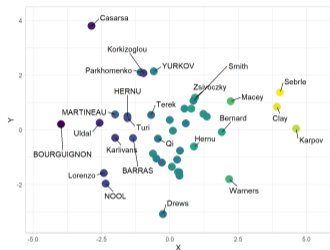
Decathlon



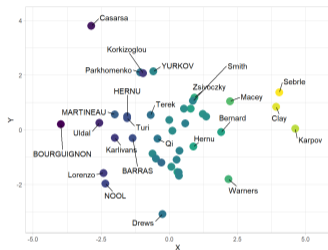
Swiss Roll



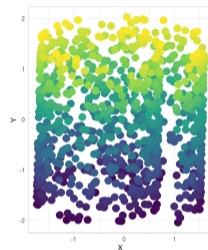
Principal Component Analysis



Decathlon



Decathlon
Renormalized



Swiss Roll

Multiple Factor Analysis

- PCA assumes $\mathcal{X} = \mathbb{R}^d$!
- How to deal with categorical values?
- MFA = PCA with clever coding strategy for categorical values.

Categorical value code for a single variable

- Classical redundant dummy coding:

$$\underline{X} \in \{1, \dots, V\} \mapsto P(\underline{X}) = (\mathbf{1}_{\underline{X}=1}, \dots, \mathbf{1}_{\underline{X}=V})^\top$$

- Compute the mean (i.e. the empirical proportions): $\bar{P} = \frac{1}{n} \sum_{i=1}^n P(\underline{X}_i)$

- Renormalize $P(\underline{X})$ by $1/\sqrt{(V-1)\bar{P}}$:

$$P(\underline{X}) = (\mathbf{1}_{\underline{X}=1}, \dots, \mathbf{1}_{\underline{X}=V}) \mapsto \left(\frac{\mathbf{1}_{\underline{X}=1}}{\sqrt{(V-1)\bar{P}_1}}, \dots, \frac{\mathbf{1}_{\underline{X}=V}}{\sqrt{(V-1)\bar{P}_V}} = P^r(\underline{X}) \right)$$

- χ^2 type distance!

- PCA becomes the minimization of

$$\frac{1}{n} \sum_{i=1}^n \|P^r(\underline{X}_i) - (m + VV^\top(P^r(\underline{X}_i) - m))\|^2$$
$$= \frac{1}{n} \sum_{i=1}^n \sum_{v=1}^V \frac{|\mathbf{1}_{\underline{X}_i=v} - (m' + \sum_{l=1}^{d'} V^{(l)\top}(P(\underline{X}_i) - m')V^{(l,v)})|^2}{(V-1)\bar{P}_v}$$

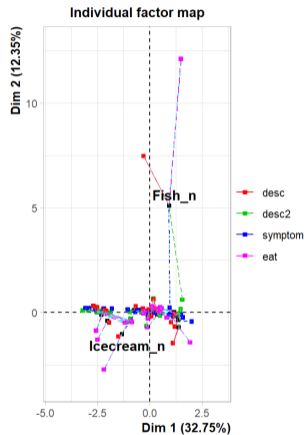
- Interpretation:

- $m' = \bar{P}$
- $\Phi(\underline{X}) = V^\top(P^r(\underline{X}) - m)$: coordinates in the restricted space.
- $V^{(l)}$ can be interpreted as a probability profile.
- Complexity: $O(n(V + V^2) + d'V^2)$
- Link with Correspondence Analysis (CA)

MFA Algorithm

- Redundant dummy coding of each categorical variable.
 - Renormalization of each block of dummy variable.
 - Classical PCA algorithm on the resulting variables
-
- Interpretation as a reconstruction error with a rescaled/ χ^2 metric.
 - Interpretation:
 - $\Phi(\underline{X}) = V^T (P^r(\underline{X}) - m)$: coordinates in the restricted space.
 - $V^{(l)}$: influence of each modality/variable in the l th new coordinates.
 - **Scaling:** This method is not invariant to a scaling of the continuous variables! It is custom to normalize the variables (at least within groups) before applying PCA.

Multiple Factor Analysis



PCA Model

- PCA: Linear model assumption

$$\underline{X} \simeq m + \sum_{l=1}^{d'} \underline{X}'^{(l)} V^{(l)} = m + V \underline{X}'$$

- with
 - $V^{(l)}$ orthonormal
 - $\underline{X}'^{(l)}$ without constraints.
- Two directions of extension:
 - Other constraints on V (or the coordinates in the restricted space): ICA, NMF, Dictionary approach
 - PCA on a non-linear image of \underline{X} : kernel-PCA
- Much more complex algorithm!

ICA (Independent Component Analysis)

- Linear model assumption

$$\underline{X} \simeq m + \sum_{l=1}^{d'} \underline{X}^{',(l)} V^{(l)} = m + V \underline{X}'$$

- with

- $V^{(l)}$ without constraints.
- $\underline{X}^{',(l)}$ independent

NMF (Non Negative Matrix Factorization)

- (Linear) Model assumption

$$\underline{X} \simeq \sum_{l=1}^{d'} \underline{X}^{',(l)} V^{(l)} = V \underline{X}'$$

- with

- $V^{(l)}$ non-negative
- $\underline{X}^{',(l)}$ non-negative.

Dictionary

- (Linear) Model assumption

$$\underline{X} \simeq m + \sum_{l=1}^{d'} \underline{X}'^{(l)} V^{(l)} = m + V \underline{X}'$$

- with
 - $V^{(l)}$ without constraints
 - \underline{X}' sparse (with a lot of 0)

kernel PCA

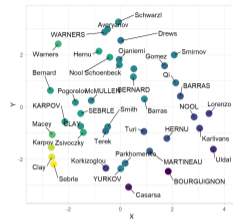
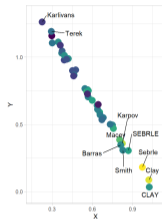
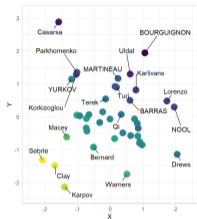
- Linear model assumption

$$\Psi(\underline{X} - m) \simeq \sum_{l=1}^{d'} \underline{X}'^{(l)} V^{(l)} = V \underline{X}'$$

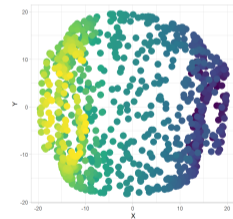
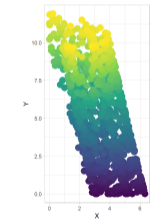
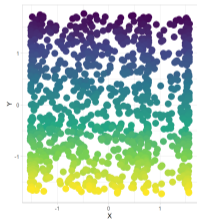
- with
 - $V^{(l)}$ orthonormal
 - \underline{X}'_l without constraints.

Non Linear PCA

Decathlon



Swiss Roll



ICA

NMF

Kernel PCA

Deep Auto Encoder

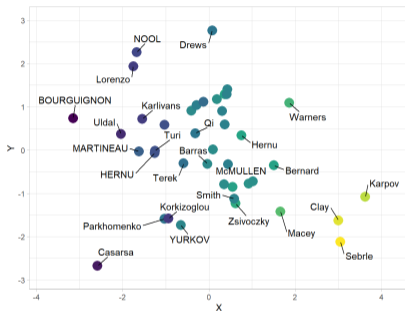
- Construct a map Φ with a **NN** from the space \mathcal{X} into a space \mathcal{X}' of smaller dimension:

$$\begin{aligned}\Phi : \mathcal{X} &\rightarrow \mathcal{X}' \\ \underline{X} &\mapsto \Phi(\underline{X})\end{aligned}$$

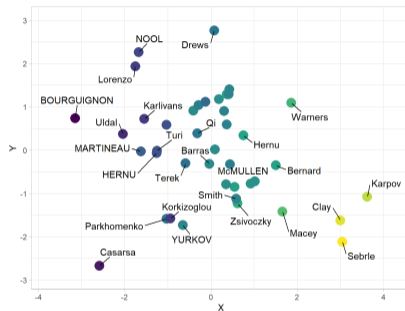
- Construct $\tilde{\Phi}$ with a **NN** from \mathcal{X}' to \mathcal{X}
- Control the error between \underline{X} and its reconstruction $\tilde{\Phi}(\Phi(\underline{X}))$:

$$\frac{1}{n} \sum_{i=1}^n \|\underline{X}_i - \tilde{\Phi}(\Phi(\underline{X}_i))\|^2$$

- Optimization by gradient descent.
- NN can be replaced by another parametric function...



Shallow Auto Encoder
(PCA)



Deep Auto Encoder

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

- Different point of view!
- Focus on pairwise relation $\mathcal{R}(\underline{X}_i, \underline{X}_j)$.

Distance Preservation

- Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of **smaller dimension**:

$$\Phi: \mathcal{X} \rightarrow \mathcal{X}'$$

$$\underline{X} \mapsto \Phi(\underline{X}) = \underline{X}'$$

- such that

$$\mathcal{R}(\underline{X}_i, \underline{X}_j) \sim \mathcal{R}'(\underline{X}'_i, \underline{X}'_j)$$

- Most classical version (MDS):

- Scalar product relation: $\mathcal{R}(\underline{X}_i, \underline{X}_j) = (\underline{X}_i - m)^\top (\underline{X}_j - m)$
- Linear mapping $\underline{X}' = \Phi(\underline{X}) = V^\top (\underline{X} - m)$.
- Euclidean scalar product matching:

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left| (\underline{X}_i - m)^\top (\underline{X}_j - m) - (\underline{X}'_i)^\top \underline{X}'_j \right|^2$$

- Φ often defined only on $\mathcal{D} \dots$

MDS Heuristic

- Match the *scalar* products:

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left| (\underline{X}_i - m)^\top (\underline{X}_j - m) - \underline{X}_i'^\top \underline{X}_j' \right|^2$$

- Linear method: $\underline{X}' = U^\top (\underline{X} - m)$ with U orthonormal

- **Beware:** \underline{X} can be unknown, only the scalar products are required!
- Resulting criterion: minimization in $U^\top (\underline{X}_i - m)$ of

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left| (\underline{X}_i - m)^\top (\underline{X}_j - m) - (\underline{X}_i - m)^\top U U^\top (\underline{X}_j - m) \right|^2$$

without using explicitly \underline{X} in the algorithm...

- Explicit solution obtained through the eigendecomposition of the know Gram matrix $(\underline{X}_i - m)^\top (\underline{X}_j - m)$ by keeping only the d' largest eigenvalues.

- In this case, MDS yields the same result as the PCA (but with different inputs, distance between observation vs correlations)!
- **Explanation:** Same SVD problem up to a transposition:

- MDS

$$\underline{\bar{X}}_{(n)}^\top \underline{\bar{X}}_{(n)} \sim \underline{\bar{X}}_{(n)}^\top U U^\top \underline{\bar{X}}_{(n)}$$

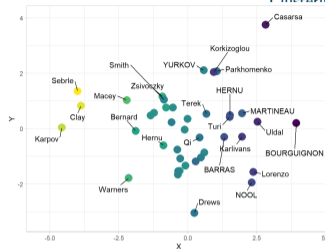
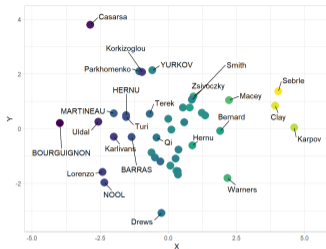
- PCA

$$\underline{\bar{X}}_{(n)} \underline{\bar{X}}_{(n)}^\top \sim U^\top \underline{\bar{X}}_{(n)} \underline{\bar{X}}_{(n)}^\top U$$

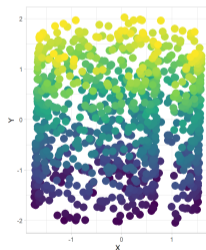
- Complexity: PCA $O((n + d')d^2)$ vs MDS $O((d + d')n^2)$...

MultiDimensional Scaling

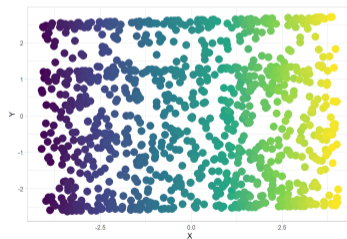
Decathlon



Swiss Roll



PCA



MDS

- Preserving the scalar products amounts to preserve the Euclidean distance.
- Easier **generalization** if we work in terms of distance!

Generalized MDS

- Generalized MDS:
 - Distance relation: $\mathcal{R}(\underline{X}_i, \underline{X}_j) = d(\underline{X}_i, \underline{X}_j)$
 - Linear mapping $\underline{X}' = \Phi(\underline{X}) = V^T(\underline{X} - m)$.
 - Euclidean matching:

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n |d(\underline{X}_i, \underline{X}_j) - d'(\underline{X}'_i, \underline{X}'_j)|^2$$

- Strong connection (but no equivalence) with MDS when $d(x, y) = \|x - y\|^2$!
- **Minimization:** Simple gradient descent can be used (can be stuck in local minima).

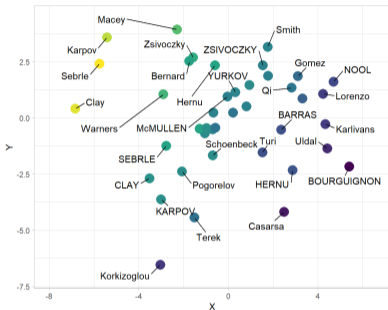
- MDS: equivalent to PCA (but more expensive) if $d(x, y) = \|x - y\|^2$!
- ISOMAP: use a *localized* distance instead to limit the influence of very far point.

ISOMAP

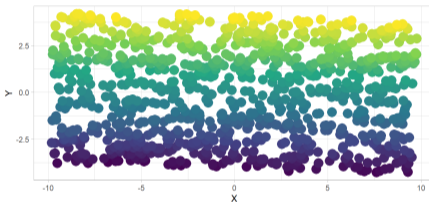
- For each point \underline{X}_i , define a neighborhood \mathcal{N}_i (either by a distance or a number of points) and let

$$d_0(\underline{X}_i, \underline{X}_j) = \begin{cases} +\infty & \text{if } \underline{X}_j \notin \mathcal{N}_i \\ \|\underline{X}_i - \underline{X}_j\|^2 & \text{otherwise} \end{cases}$$

- Compute the shortest path distance for each pair.
- Use the MDS algorithm with this distance



Decathlon

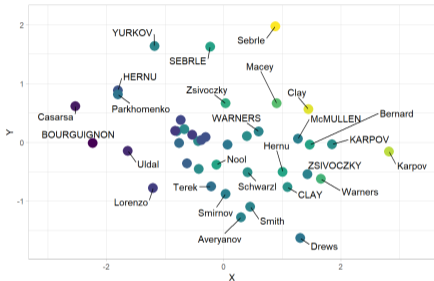


Swiss Roll

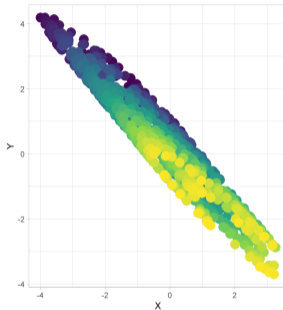
Random Projection Heuristic

- Draw at random d' unit vector (direction) U_j .
- Use $\underline{X}' = U^\top (\underline{X} - m)$ with $m = \frac{1}{n} \sum_{i=1}^n \underline{X}_i$
- **Property:** If \underline{X} lives in a space of dimension d'' , then, as soon as, $d' \sim d'' \log(d'')$,
$$\|\underline{X}_i - \underline{X}_j\|^2 \sim \frac{d}{d'} \|\underline{X}'_i - \underline{X}'_j\|^2$$
- Do not really use the data!

Random Projection



Decathlon



Swiss Roll

SNE heuristic

- From $\underline{X}_i \in \mathcal{X}$, construct a set of conditional probability:

$$P_{j|i} = \frac{e^{-\|\underline{X}_i - \underline{X}_j\|^2 / 2\sigma_i^2}}{\sum_{k \neq i} e^{-\|\underline{X}_i - \underline{X}_k\|^2 / 2\sigma_i^2}} \quad P_{i|i} = 0$$

- Find \underline{X}'_i in $\mathbb{R}^{d'}$ such that the set of conditional probability:

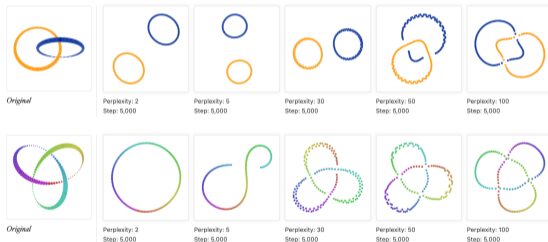
$$Q_{j|i} = \frac{e^{-\|\underline{X}'_i - \underline{X}'_j\|^2 / 2\sigma_i^2}}{\sum_{k \neq i} e^{-\|\underline{X}'_i - \underline{X}'_k\|^2 / 2\sigma_i^2}} \quad Q_{i|i} = 0$$

is close from P .

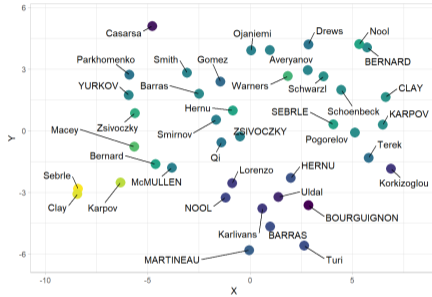
- t-SNE:** use a Student-t term $(1 + \|\underline{X}'_i - \underline{X}'_j\|^2)^{-1}$ for \underline{X}'_i
- Minimize the Kullback-Leibler divergence $(\sum_{i,j} P_{j|i} \log \frac{P_{j|i}}{Q_{j|i}})$ by a simple gradient descent (can be stuck in local minima).
- Parameters σ_i such that $H(P_i) = -\sum_{j=1}^n P_{j|i} \log P_{j|i} = \text{cst.}$

t-Stochastic Neighbor Embedding

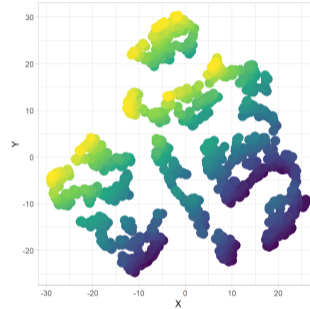
- Very successful/ powerful technique in practice
- Convergence may be long, unstable, or strongly depending on parameters.
- See this [distill post](#) for many impressive examples



Representation depending on t-SNE parameters



Decathlon



Swiss Roll

- Topological Data Analysis inspired.

Uniform Manifold Approximation and Projection

- Define a notion of asymmetric scaled local proximity between neighbors:
 - Compute the k -neighborhood of \underline{X}_i , its diameter σ_i and the distance ρ_i between \underline{X}_i and its nearest neighbor.
 - Define

$$w_i(\underline{X}_i, \underline{X}_j) = \begin{cases} e^{-(d(\underline{X}_i, \underline{X}_j) - \rho_i) / \sigma_i} & \text{for } \underline{X}_j \text{ in the } k\text{-neighborhood} \\ 0 & \text{otherwise} \end{cases}$$

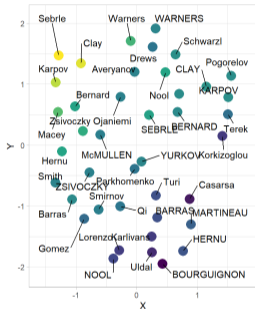
- Symmetrize into a *fuzzy* nearest neighbor criterion

$$w(\underline{X}_i, \underline{X}_j) = w_i(\underline{X}_i, \underline{X}_j) + w_j(\underline{X}_j, \underline{X}_i) - w_i(\underline{X}_i, \underline{X}_j)w_j(\underline{X}_j, \underline{X}_i)$$

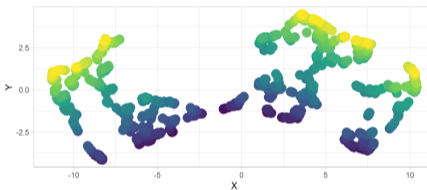
- Determine the points \underline{X}'_i in a low dimensional space such that

$$\sum_{i \neq j} w(\underline{X}_i, \underline{X}_j) \log \left(\frac{w(\underline{X}_i, \underline{X}_j)}{w'(\underline{X}'_i, \underline{X}'_j)} \right) + (1 - w(\underline{X}_i, \underline{X}_j)) \log \left(\frac{(1 - w(\underline{X}_i, \underline{X}_j))}{(1 - w'(\underline{X}'_i, \underline{X}'_j))} \right)$$

- Can be performed by local gradient descent.



Decathlon



Swiss Roll

Graph heuristic

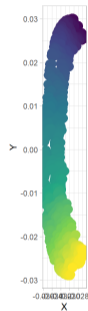
- Construct a graph with weighted edges $w_{i,j}$ measuring the *proximity* of \underline{X}_i and \underline{X}_j ($w_{i,j}$ large if close and 0 if there is no information).
- Find the points $\underline{X}'_j \in \mathbb{R}^{d'}$ minimizing

$$\frac{1}{n} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n w_{i,j} \|\underline{X}'_i - \underline{X}'_j\|^2$$

- Need of a constraint on the size of \underline{X}'_j ...
- Explicit solution through linear algebra: d' eigenvectors with smallest eigenvalues of the Laplacian of the graph $D - W$, where D is a diagonal matrix with $D_{i,i} = \sum_j w_{i,j}$.
- Variation on the definition of the Laplacian...



Decathlon



Swiss Roll

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

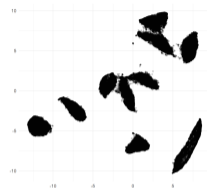
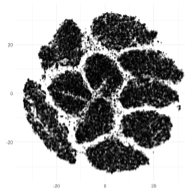
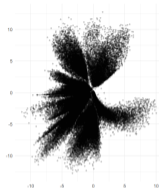
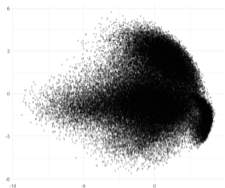
How to Compare Different Dimensionality Reduction Methods ?

- **Difficult!** Once again, the metric is very subjective.

However, a few possible attempts

- Did we preserve a lot of inertia with only a few directions?
- Do those directions *make sense* from an expert point of view?
- Do the low dimension representation *preserve* some important information?
- Are we better on **subsequent task**?

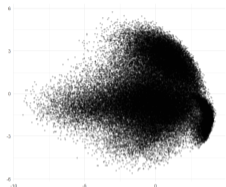
A Challenging Example: MNIST



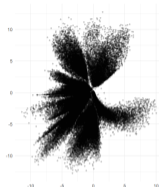
MNIST Dataset

- Images of 28×28 pixels.
- No label used!
- 4 different embeddings.

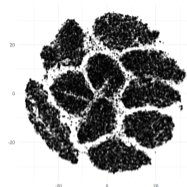
A Challenging Example: MNIST



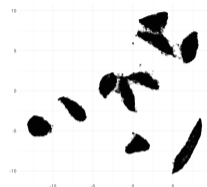
PCA



autoencoder



t-SNE

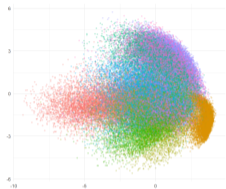


UMAP

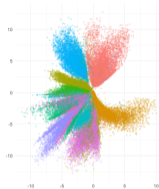
MNIST Dataset

- Images of 28×28 pixels.
- No label used!
- 4 different embeddings.

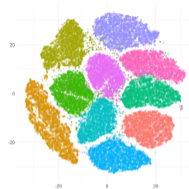
A Challenging Example: MNIST



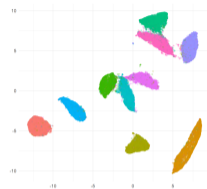
PCA



autoencoder



t-SNE

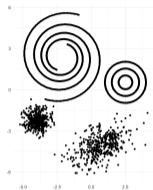


UMAP

MNIST Dataset

- Images of 28×28 pixels.
- No label used!
- 4 different embeddings.
- Quality evaluated by visualizing the true labels **not used to obtain the embeddings.**
- Only a few labels could have been used.

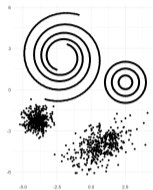
A Simpler Example: A 2D Set



Cluster Dataset

- Set of points in 2D.
- No label used!
- 3 different embeddings.

A Simpler Example: A 2D Set



PCA



t-SNE

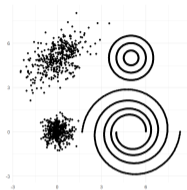


UMAP

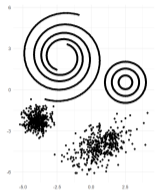
Cluster Dataset

- Set of points in 2D.
- No label used!
- 3 different embeddings.

A Simpler Example: A 2D Set



Original



PCA



t-SNE



UMAP

Cluster Dataset

- Set of points in 2D.
- No label used!
- 3 different embeddings.
- Quality evaluated by stability. . .

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- **Training data** : $\mathcal{D} = \{\underline{X}_1, \dots, \underline{X}_n\} \in \mathcal{X}^n$ (i.i.d. $\sim \mathbb{P}$)
- Latent groups?

Clustering

- Construct a map f from \mathcal{D} to $\{1, \dots, K\}$ where K is a number of classes to be fixed:

$$f : \underline{X}_j \mapsto k_j$$

Motivations

- Interpretation of the groups
- Use of the groups in further processing

- Several strategies possible!
- Can use dimension reduction as a preprocessing.

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Partition Heuristic

- Clustering is defined by a partition in K classes. . .
- that minimizes a homogeneity criterion.

K- Means

- Cluster k defined by a *center* μ_k .
- Each sample is associated to the closest center.
- Centers defined as the minimizer of
$$\sum_{i=1}^n \min_k \|\underline{X}_i - \mu_k\|^2$$
- Iterative scheme (Lloyd):
 - Start by a (pseudo) random choice for the centers μ_k
 - Assign each samples to its nearby center
 - Replace the center of a cluster by the mean of its assigned samples.
 - Repeat the last two steps until convergence.

Partition Based

Unsupervised Learning:
Dimension Reduction and
Clustering



- Other schemes:
 - McQueen: modify the mean each time a sample is assigned to a new cluster.
 - Hartigan: modify the mean by removing the considered sample, assign it to the nearby center and recompute the new mean after assignment.

A good initialization is crucial!

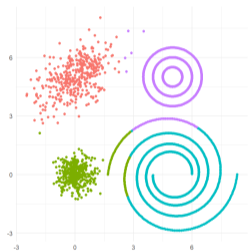
- Initialize by samples.
 - k-Mean++: try to take them as separated as possible.
 - No guarantee to converge to a global optimum: repeat and keep the best result!
-
- Complexity : $O(n \times K \times T)$ where T is the number of steps in the algorithm.

- k-Medoid: use a sample as a center
 - PAM: for a given cluster, use the sample that minimizes the intra distance (sum of the squared distance to the other points)
 - Approximate medoid: for a given cluster, assign the point that is the closest to the mean.

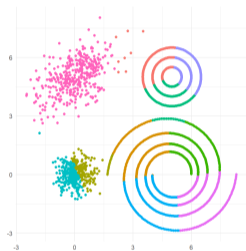
Complexity

- PAM: $O(n^2 \times T)$ in the worst case!
- Approximate medoid: $O(n \times K \times T)$ where T is the number of steps in the algorithm.
- **Remark:** Any distance can be used... but the complexity of computing the centers can be very different.

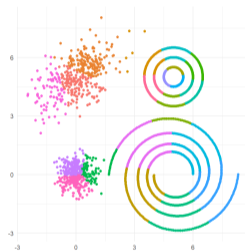
K-Means



$k = 4$



$k = 10$



$k = 10$

Model Heuristic

- Use a generative model of the data:

$$\mathbb{P}(\underline{X}) = \sum_{k=1}^K \pi_k \mathbb{P}_{\theta_k}(\underline{X}|k)$$

where π_k are proportions and $\mathbb{P}_{\theta}(\underline{X}|k)$ are parametric probability models.

- Estimate those parameters (often by a ML principle).
- Assign each observation to the class maximizing the a posteriori probability (obtained by Bayes formula)

$$\frac{\widehat{\pi}_k \mathbb{P}_{\widehat{\theta}_k}(\underline{X}|k)}{\sum_{k'=1}^K \widehat{\pi}_{k'} \mathbb{P}_{\widehat{\theta}_{k'}}(\underline{X}|k')}$$

- Link with Generative model in supervised classification!

- Large choice of parametric models.

Gaussian Mixture Model

- Use

$$\mathbb{P}_{\theta_k}(\vec{X}|k) \sim \mathcal{N}(\mu_k, \Sigma_k)$$

with $\mathcal{N}(\mu, \Sigma)$ the Gaussian law of mean μ and covariance matrix Σ .

- Efficient optimization algorithm available (EM)
- Often some constraint on the covariance matrices: identical, with a similar structure. . .
- Strong connection with K -means when the covariance matrices are assumed to be the same multiple of the identity.

Probabilistic latent semantic analysis (PLSA)

- Documents described by their word counts w
- Model:

$$\mathbb{P}(w) = \sum_{k=1}^K \pi_k \mathbb{P}_{\theta_k}(w|k)$$

with k the (hidden) topic, π_k a topic probability and $\mathbb{P}_{\theta_k}(w|k)$ a multinomial law for a given topic.

- Clustering according to

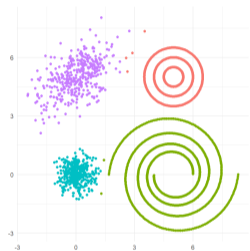
$$\mathbb{P}(k|w) = \frac{\hat{\pi}_k \mathbb{P}_{\hat{\theta}_k}(w|k)}{\sum_{k'} \hat{\pi}_{k'} \mathbb{P}_{\hat{\theta}_{k'}}(w|k')}$$

- Same idea than GMM!
- Bayesian variant called LDA.

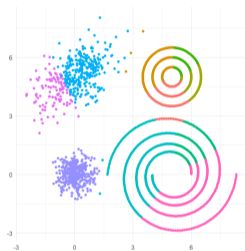
Parametric Density Estimation Principle

- Assign a probability of membership.
- Lots of theoretical studies. . .
- Model selection principle can be used to select K the number of classes (or rather to avoid using a nonsensical K . . .):
 - AIC / BIC / MDL penalization
 - Cross Validation is also possible!
- Complexity: $O(n \times K \times T)$

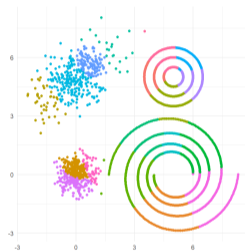
Gaussian Mixture Models



$k = 4$



$k = 10$



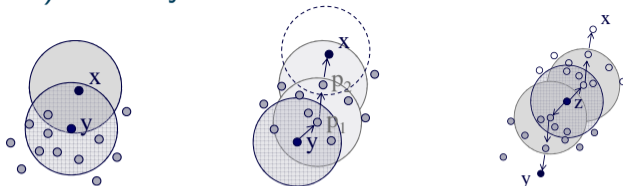
$k = 10$

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Density Heuristic

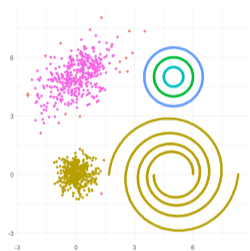
- Cluster are connected dense zone separated by low density zone.
- Not all points belong to a cluster.
- Basic bricks:
 - Estimate the density.
 - Find points with high densities.
 - Gather those points according to the density.
- Density estimation:
 - Classical kernel density estimators. . .
- Gathering:
 - Link points of high density and use the resulted component.
 - Move them toward top of density *hill* by following the gradient and gather all the points arriving at the same *summit*.

(Non Parametric) Density Based

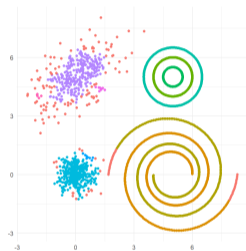


Examples

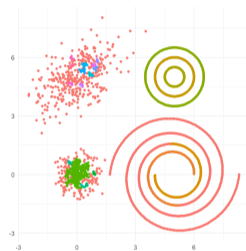
- DBSCAN: link point of high densities using a very simple kernel.
- PdfCLuster: find connected zone of high density.
- Mean-shift: move points toward top of density *hill* following an evolving kernel density estimate.
- Complexity: $O(n^2 \times T)$ in the worst case.
- Can be reduced to $O(n \log(n) T)$ if samples can be encoded in a tree structure (n-body problem type approximation).



$\epsilon = .45$



$\epsilon = .2$



$\epsilon = .1$

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

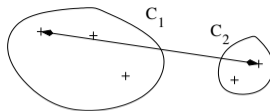
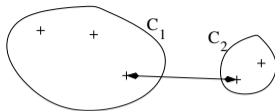
Agglomerative Clustering Heuristic

- Start with very small clusters (a sample by cluster?)
 - Sequential merging of the most similar clusters. . .
 - according to some *greedy* criterion Δ .
-
- Generates a hierarchy of clustering instead of a single one.
 - Need to select the number of cluster afterwards.
 - Several choices for the merging criterion. . .
 - Examples:
 - Minimum Linkage: merge the closest cluster in term of the usual distance
 - Ward's criterion: merge the two clusters yielding the less inner inertia loss (k-means criterion)

Algorithm

- Start with $(\mathcal{C}_i^{(0)}) = (\{\underline{X}_i\})$ the collection of all singletons.
- At step s , we have $n - s$ clusters $(\mathcal{C}_i^{(s)})$:
 - Find the two most similar clusters according to a criterion Δ :
$$(i, i') = \underset{(j, j')}{\operatorname{argmin}} \Delta(\mathcal{C}_j^{(s)}, \mathcal{C}_{j'}^{(s)})$$
 - Merge $\mathcal{C}_i^{(s)}$ and $\mathcal{C}_{i'}^{(s)}$ into $\mathcal{C}_i^{(s+1)}$
 - Keep the $n - s - 2$ other clusters $\mathcal{C}_{i''}^{(s+1)} = \mathcal{C}_{i''}^{(s)}$
- Repeat until there is only one cluster.
- Complexity: $O(n^3)$ in general.
- Can be reduced to $O(n^2)$
 - if only a bounded number of merging is possible for a given cluster,
 - for the most classical distances by maintaining a nearest neighbors list.

Agglomerative Clustering



Merging criterion based on the distance between points

- Minimum linkage:

$$\Delta(C_i, C_j) = \min_{\underline{X}_i \in C_i} \min_{\underline{X}_j \in C_j} d(\underline{X}_i, \underline{X}_j)$$

- Maximum linkage:

$$\Delta(C_i, C_j) = \max_{\underline{X}_i \in C_i} \max_{\underline{X}_j \in C_j} d(\underline{X}_i, \underline{X}_j)$$

- Average linkage:

$$\Delta(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{\underline{X}_i \in C_i} \sum_{\underline{X}_j \in C_j} d(\underline{X}_i, \underline{X}_j)$$

- Clustering based on the proximity...

Merging criterion based on the inertia (distance to the mean)

- Ward's criterion:

$$\begin{aligned}\Delta(C_i, C_j) &= \sum_{\underline{X}_i \in C_i} \left(d^2(\underline{X}_i, \mu_{C_i \cup C_j}) - d^2(\underline{X}_i, \mu_{C_i}) \right) \\ &\quad + \sum_{\underline{X}_j \in C_j} \left(d^2(\underline{X}_j, \mu_{C_i \cup C_j}) - d^2(\underline{X}_j, \mu_{C_j}) \right)\end{aligned}$$

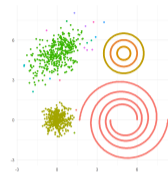
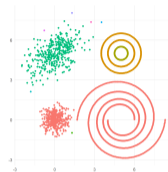
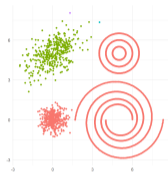
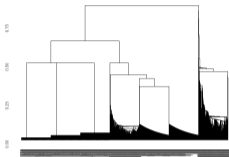
- If d is the Euclidean distance:

$$\Delta(C_i, C_j) = \frac{2|C_i||C_j|}{|C_i| + |C_j|} d^2(\mu_{C_i}, \mu_{C_j})$$

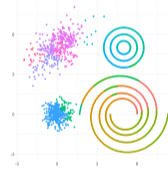
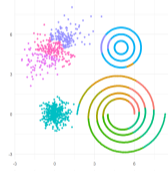
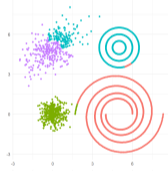
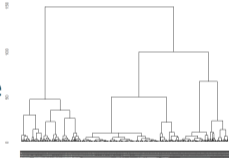
- Same criterion than in the k -means algorithm but greedy optimization.

Agglomerative Clustering

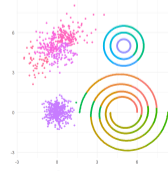
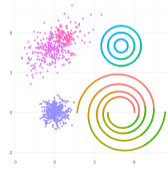
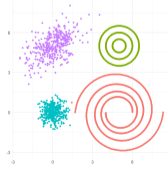
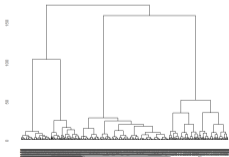
Single



Complete



Ward



Dendrogram

$k = 4$

$k = 10$

$k = 20$

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Grid heuristic

- Split the space in pieces
- Group those of high density according to their proximity
- Similar to density based estimate (with partition based initial clustering)
- Space splitting can be fixed or adaptive to the data.
- Examples:
 - STING (Statistical Information Grid): Hierarchical tree construction plus DBSCAN type algorithm
 - AMR (Adaptive Mesh Refinement): Adaptive tree refinement plus k -means type assignment from high density leaves.
 - CLIQUE: Tensorial grid and 1D detection.
- Linked to Divisive clustering (DIANA)

Graph based

- Spectral clustering: dimension reduction + k-means.
- Message passing: iterative local algorithm.
- Graph cut: min/max flow.

- Kohonen Map,
- ...

Outline



- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Large dataset issue

- When n is large, a $O(n^\alpha \log n)$ with $\alpha > 1$ is not acceptable!
- How to deal with such a situation?
- **Beware:** Computing all the pairwise distance requires $O(n^2)$ operations!

Ideas

- Sampling
- Online processing
- Simplification
- Parallelization

Sampling heuristic

- Use only a subsample to construct the clustering.
 - Assign the other points to the constructed clusters afterwards.
-
- Requires a clustering method that can assign new points (partition, model. . .)
 - Often repetition and choice of the best clustering
 - Example:
 - CLARA: K-medoid with sampling and repetition
 - Two-steps algorithm:
 - Generate a large number n' of clusters using a fast algorithm (with $n' \ll n$)
 - Cluster the clusters with a more accurate algorithm.

Online heuristic

- Modify the current clusters according to the value of a single observation.
- Requires compactly described clusters.
- Examples:
 - Add to an existing cluster (and modify it) if it is close enough and create a new cluster otherwise (k -means without reassignment)
 - Stochastic descent gradient (GMM)
- May leads to far from optimal clustering.

Simplification heuristic

- Simplify the algorithm to be more efficient at the cost of some precision.
- Algorithm dependent!
- Examples:
 - Replace groups of observation (preliminary cluster) by the (approximate) statistics.
 - Approximate the distances by cheaper ones.
 - Use n-body type techniques.

Parallelization heuristic

- Split the computation on several computers.
- Algorithm dependent!
- Examples:
 - Distance computation in k -means, parameter gradient in model based clustering
 - Grid density estimation, Space splitting strategies
- Classical batch sampling not easy to perform as partitions are not easily merged. . .

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Generative Modeling

- **Training data** : $\mathcal{D} = \{(\underline{X}_1, \underline{Y}_1), \dots, (\underline{X}_n, \underline{Y}_n)\} \in (\mathcal{X} \times \mathcal{Y})^n$ (i.i.d. $\sim \mathbb{P}$)
- Same kind of data than for supervised learning if $\mathcal{Y} \neq \emptyset$.

Generative Modeling

- Construct a map G from the product of \mathcal{Y} and a randomness source Ω to \mathcal{X}

$$G : \mathcal{Y} \times \Omega \rightarrow \mathcal{X}$$

$$(Y, \omega) \mapsto X$$

- Unconditional model if $\mathcal{Y} = \emptyset$...

Motivation

- Generate plausible novel conditional samples based on a given dataset.

Sample Quality

- Related to the proximity between the law of $G(Y, \omega)$ and the law of $X|Y$.
- Most classical choice is the Kullback-Leibler divergence.

Ingredients

- Generator $F_\theta(Y, \omega)$ and cond. density prob. $p_\theta(X|Y)$ (Explicit vs implicit link)
- Simple / Complex / Approximate estimation...

Some Possible Choices

	Probabilistic model	Generator	Estimation
Base	Simple (parametric)	Explicit	Simple (ML)
Flow	Image of simple model	Explicit	Simple (ML)
Factorization	Factorization of simple model	Explicit	Simple (ML)
VAE	Simple model with latent var.	Explicit	Approximate (ML)
EBM	Arbitrary	Implicit (MCMC)	Complex (ML/score/discrim.)
Diffusion	Continuous noise	Implicit (MCMC)	Complex (score)
	Discrete Noise with latent var.	Explicit	Approximate (ML)
GAN	Implicit	Explicit	Complex (Discrimination)

- SOTA: Diffusion based approach!

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Density Estimation

Unsupervised Learning:
Dimension Reduction and
Clustering



Factorization

Unsupervised Learning:
Dimension Reduction and
Clustering



Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

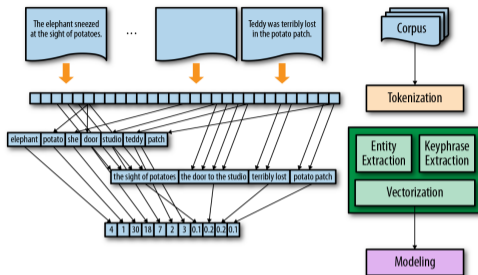
Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References



Text and Representation

- Need to transform a text into a numerical vector to reuse the previous algorithms!
- Art still in progress.
- Important steps:
 - Token extraction
 - Token vectorization
 - Learning algorithm

Stemming

adjustable → adjust
formality → formaliti
formaliti → formal
airliner → airlin ⚠

Lemmatization

was → (to) be
better → good
meeting → meeting

Token Extraction

- From a text to a sequence of *tokens* (words, characters, subwords. . .).
- Need of cleaning or pre-processing: spelling checker, stemming, lemmatization. . .
- Often with a further reduction of the number of possible tokens.
- Beware to not oversimplify!

The Bag of Words Representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

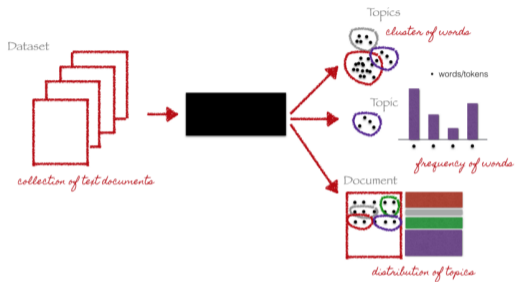
15



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

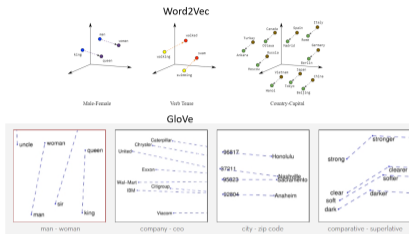
Bag of Words

- Most simple approach to transform a text into a vector.
- Simple count of the words belonging to a predefined vocabulary.
- Counts preferably replaced by frequencies (or tf-idf. . .)
- Often combined with dimension reduction:
 - restriction to an interesting vocabulary
 - use of principal component analysis (latent semantic analysis)



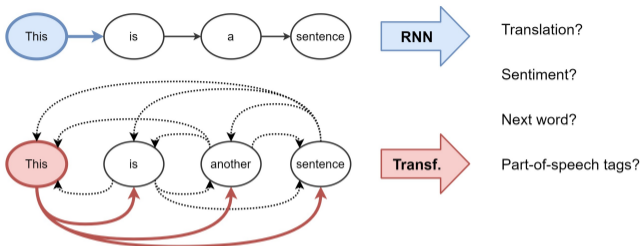
Article Clustering

- Clustering algorithms directly on the bag-of-words representation
- Most used algorithm is a variation around the k -means algorithm.



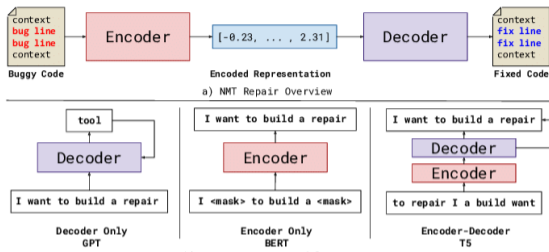
Word Representation

- More accuracy by working at the token (word) scale.
- Two approaches:
 - Associate to a word the frequency of the other words in its neighborhood and performing dimension reduction on this first representation.
 - Learn for each word a vector allowing to predict by a simple formula (scalar product) whether one word appears in the neighborhood of the other one.
- Similar results but the second approach is more flexible.



Deep Learning

- Propose a formula allowing to do computations on the word starting by associating vectors to each word.
- Learning the best possible vectors for a given task: auto-prediction (self-supervised) or prediction (supervised).
- Tremendous progress in the last years thanks to deep neural net architectures (RNN, Transformer...).



Large Language Models

- Huge neural networks relying on transformers and (pre)trained on huge corpus with self-supervised tasks.
- Three architectures:
 - Decoder: prediction of next word (online).
 - Encoder: prediction of inner word(s) (offline).
 - Encoder/Decoder: prediction of a sentence from another (offline).
- Can be used as a basis for further specialized training or directly.



How to associate a sentiment to a text?

- Four possible approaches:
 - Simple approach (without learning) that averages the *sentiments* of the words used in a text using a fixed table.
 - Simple approach (supervised and linear) where this table is learned from examples.
 - Direct approach (supervised) where one predicts directly the sentiment from examples.
 - Zero-Shot approach (without learning?) where one uses directly a Large Language Model trained on a huge corpus (unrelated to the application).
- Direct approach more efficient provided one has sufficient data and one starts from a pretrained model.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - **References**
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References



T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning.
Springer Series in Statistics, 2009



M. Mohri, A. Rostamizadeh, and A. Talwalkar.
Foundations of Machine Learning (2nd ed.)
MIT Press, 2018



A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



Ch. Giraud.
Introduction to High-Dimensional Statistics.
CRC Press, 2014



S. Bubeck.
Convex Optimization: Algorithms and Complexity.
Now Publisher, 2015

References



F. Husson, S. Le, and J. Pagès.
Exploratory Multivariate Analysis by Example Using R (2nd ed.)
Chapman and Hall/CRC, 2017



B. Ghojogh, M. Crowley, F. Karray, and
A. Ghodsi.
Elements of Dimensionality Reduction and Manifold Learning.
Springer, 2023



Ch. Aggarwal and Ch. Reddy.
Data Clustering: Algorithms and Applications.
Chapman and Hall/CRC, 2013



Ch. Hennig, M. Meila, F. Murtagh, and
R. Rocci.
Handbook of Cluster Analysis.
Chapman and Hall/CRC, 2015



Ch. Bouveyron, G. Celeux, B. Murphy,
and A. Raftery.
Model-Based Clustering and Classification for Data Science.
Cambridge University Press, 2019



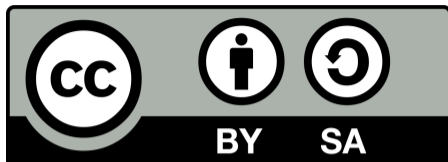
J. Tomczak.
Deep Generative Modeling.
Springer, 2021



D. Foster.
Generative Deep Learning (2nd ed.)
O'Reilly, 2023



A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 6 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 6 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 6 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 6 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 6 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 6 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 Applications to Text
 - References
- 8 **Statistical Learning: PAC-Bayesian Approach and Complexity Theory**
 - **Supervised Learning**
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Supervised Learning Framework

- Input measurement $\underline{X} \in \mathcal{X}$
- Output measurement $Y \in \mathcal{Y}$.
- $(\underline{X}, Y) \sim \mathbb{P}$ with \mathbb{P} unknown.
- **Training data** : $\mathcal{D}_n = \{(\underline{X}_1, Y_1), \dots, (\underline{X}_n, Y_n)\}$ (i.i.d. $\sim \mathbb{P}$)
- Often
 - $\underline{X} \in \mathbb{R}^d$ and $Y \in \{-1, 1\}$ (classification)
 - or $\underline{X} \in \mathbb{R}^d$ and $Y \in \mathbb{R}$ (regression).
- A **predictor** is a function in $\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \text{ meas.}\}$

Goal

- Construct a **good** predictor \hat{f} from the training data.
- Need to specify the meaning of good.
- Classification and regression are almost the **same** problem!

Loss function for a generic predictor

- **Loss function:** $\ell(Y, f(\underline{X}))$ measures the goodness of the prediction of Y by $f(\underline{X})$
- Examples:
 - 0/1 loss: $\ell(Y, f(\underline{X})) = \mathbf{1}_{Y \neq f(\underline{X})}$
 - Quadratic loss: $\ell(Y, f(\underline{X})) = |Y - f(\underline{X})|^2$

Risk function

- Risk measured as the average loss for a new couple:

$$\mathcal{R}(f) = \mathbb{E}_{(X, Y) \sim \mathbb{P}}[\ell(Y, f(\underline{X}))]$$

- Examples:
 - 0/1 loss: $\mathbb{E}[\ell(Y, f(\underline{X}))] = \mathbb{P}(Y \neq f(\underline{X}))$
 - Quadratic loss: $\mathbb{E}[\ell(Y, f(\underline{X}))] = \mathbb{E}[|Y - f(\underline{X})|^2]$

- **Beware:** As \hat{f} depends on \mathcal{D}_n , $\mathcal{R}(\hat{f})$ is a random variable!

- The best solution f^* (which is independent of \mathcal{D}_n) is

$$f^* = \arg \min_{f \in \mathcal{F}} \mathcal{R}(f) = \arg \min_{f \in \mathcal{F}} \mathbb{E}[\ell(Y, f(\underline{X}))] = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{\underline{X}} \left[\mathbb{E}_{Y|\underline{X}}[\ell(Y, f(\underline{X}))] \right]$$

Bayes Predictor (explicit solution)

- In binary classification with 0 – 1 loss:

$$f^*(\underline{X}) = \begin{cases} +1 & \text{if } \mathbb{P}(Y = +1|\underline{X}) \geq \mathbb{P}(Y = -1|\underline{X}) \\ & \Leftrightarrow \mathbb{P}(Y = +1|\underline{X}) \geq 1/2 \\ -1 & \text{otherwise} \end{cases}$$

- In regression with the quadratic loss

$$f^*(\underline{X}) = \mathbb{E}[Y|\underline{X}]$$

Issue: Solution requires to **know** $\mathbb{E}[Y|\underline{X}]$ for all values of \underline{X} !

Machine Learning

- Learn a rule to construct a **predictor** $\hat{f} \in \mathcal{F}$ from the training data \mathcal{D}_n s.t. **the risk** $\mathcal{R}(\hat{f})$ is **small on average** or with high probability with respect to \mathcal{D}_n .
- In practice, the rule should be an algorithm!

Canonical example: Empirical Risk Minimizer

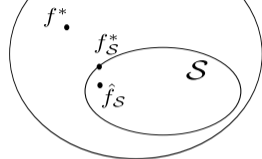
- One restricts f to a subset of functions $\mathcal{S} = \{f_\theta, \theta \in \Theta\}$
- One replaces the minimization of the average loss by the minimization of the empirical loss

$$\hat{f} = f_{\hat{\theta}} = \operatorname{argmin}_{f_\theta, \theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_\theta(\underline{X}_i))$$

- Examples:
 - Linear regression
 - Linear classification with

$$\mathcal{S} = \{\underline{x} \mapsto \operatorname{sign}\{\underline{x}^\top \beta + \beta^{(0)}\} / \beta \in \mathbb{R}^d, \beta^{(0)} \in \mathbb{R}\}$$

- General setting:
 - $\mathcal{F} = \{\text{measurable functions } \mathcal{X} \rightarrow \mathcal{Y}\}$
 - Best solution: $f^* = \operatorname{argmin}_{f \in \mathcal{F}} \mathcal{R}(f)$
 - Class $\mathcal{S} \subset \mathcal{F}$ of functions
 - Ideal target in \mathcal{S} : $f_S^* = \operatorname{argmin}_{f \in \mathcal{S}} \mathcal{R}(f)$
 - Estimate in \mathcal{S} : \hat{f}_S obtained with some procedure

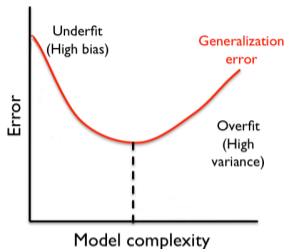


Approximation error and estimation error (Bias/Variance)

$$\mathcal{R}(\hat{f}_S) - \mathcal{R}(f^*) = \underbrace{\mathcal{R}(f_S^*) - \mathcal{R}(f^*)}_{\text{Approximation error}} + \underbrace{\mathcal{R}(\hat{f}_S) - \mathcal{R}(f_S^*)}_{\text{Estimation error}}$$

- Approx. error can be large if the model \mathcal{S} is not suitable.
- Estimation error can be large if the model is complex.

Under-fitting / Over-fitting Issue



- Different behavior for different model complexity
- **Low complexity model** are easily learned but the approximation error (**bias**) may be large (**Under-fit**).
- **High complexity model** may contain a good ideal target but the estimation error (**variance**) can be large (**Over-fit**)

Bias-variance trade-off \iff avoid **overfitting** and **underfitting**

- **Rk:** Better to think in term of method (including feature engineering and specific algorithm) rather than only of model.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 8
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Empirical Risk Minimizer (ERM)

- For any loss ℓ and function class \mathcal{S} ,

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(\underline{X}_i)) = \operatorname{argmin}_{f \in \mathcal{S}} \mathcal{R}_n(f)$$

- Key property:

$$\mathcal{R}_n(\hat{f}) \leq \mathcal{R}_n(f), \forall f \in \mathcal{S}$$

- **Minimization not always tractable in practice!**
- Focus on the $\ell^{0/1}$ case:
 - only algorithm is to try all the functions,
 - not feasible if there are many functions
 - but interesting hindsight!

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

- Theoretical control of the random (error estimation) term:

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*)$$

Probably Almost Correct Analysis

- **Theoretical guarantee** that

$$\mathbb{P}\left(\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \epsilon_S(\delta)\right) \geq 1 - \delta$$

for a suitable $\epsilon_S(\delta) \geq 0$.

- Implies:

- $\mathbb{P}\left(\mathcal{R}(\hat{f}) - \mathcal{R}(f^*) \leq \mathcal{R}(f_S^*) - \mathcal{R}(f^*) + \epsilon_S(\delta)\right) \geq 1 - \delta$

- $\mathbb{E}\left[\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*)\right] \leq \int_0^{+\infty} \delta_S(\epsilon) d\epsilon$

- The result should hold without any assumption on the law ***P***!

- By construction:

$$\begin{aligned}\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) &= \mathcal{R}(\hat{f}) - \mathcal{R}_n(\hat{f}) + \mathcal{R}_n(\hat{f}) - \mathcal{R}_n(f_S^*) + \mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*) \\ &\leq \mathcal{R}(\hat{f}) - \mathcal{R}_n(\hat{f}) + \mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*) \\ &\leq \left(\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \right) - \left(\mathcal{R}_n(\hat{f}) - \mathcal{R}_n(f_S^*) \right)\end{aligned}$$

Four possible upperbounds

- $\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sup_{f \in \mathcal{S}} ((\mathcal{R}(f) - \mathcal{R}(f_S^*)) - (\mathcal{R}_n(f) - \mathcal{R}_n(f_S^*)))$
- $\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) + (\mathcal{R}_n(f_S^*) - \mathcal{R}(f_S^*))$
- $\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) + \sup_{f \in \mathcal{S}} (\mathcal{R}_n(f) - \mathcal{R}(f))$
- $\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq 2 \sup_{f \in \mathcal{S}} |\mathcal{R}(f) - \mathcal{R}_n(f)|$

- Supremum of centered random variables!
- **Key:** Concentration of each variable...

- By construction, for any $f' \in \mathcal{S}$,

$$\mathcal{R}(f') = \mathcal{R}_n(f') + (\mathcal{R}(f') - \mathcal{R}_n(f'))$$

A uniform upper bound for the risk

- Simultaneously $\forall f' \in \mathcal{S}$,

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f))$$

- Supremum of centered random variables!
- **Key:** Concentration of each variable. . .
- Can be interpreted as a justification of the ERM!

Outline



- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

- Empirical loss:

$$\mathcal{R}_n(f) = \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f(\underline{X}_i))$$

Properties

- $\ell^{0/1}(Y_i, f(\underline{X}_i))$ are i.i.d. random variables in $[0, 1]$.

Concentration

$$\mathbb{P}(\mathcal{R}(f) - \mathcal{R}_n(f) \leq \epsilon) \geq 1 - e^{-2n\epsilon^2}$$

$$\mathbb{P}(\mathcal{R}_n(f) - \mathcal{R}(f) \leq \epsilon) \geq 1 - e^{-2n\epsilon^2}$$

$$\mathbb{P}(|\mathcal{R}_n(f) - \mathcal{R}(f)| \leq \epsilon) \geq 1 - 2e^{-2n\epsilon^2}$$

- Concentration of sum of bounded independent variables!
- Hoeffding theorem.
- Equiv. to $\mathbb{P}\left(\mathcal{R}(f) - \mathcal{R}_n(f) \leq \sqrt{\log(1/\delta)/(2n)}\right) \geq 1 - \delta$

Theorem

- Let Z_i be a sequence of ind. centered r.v. supported in $[a_i, b_i]$ then

$$\mathbb{P}\left(\sum_{i=1}^n Z_i \geq \epsilon\right) \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}}$$

- Proof ingredients:

- Chernov bounds:

$$\mathbb{P}\left(\sum_{i=1}^n Z_i \geq \epsilon\right) \leq \frac{\mathbb{E}\left[e^{\lambda \sum_{i=1}^n Z_i}\right]}{e^{\lambda \epsilon}} \leq \frac{\prod_{i=1}^n \mathbb{E}\left[e^{\lambda Z_i}\right]}{e^{\lambda \epsilon}}$$

- Exponential moment bounds: $\mathbb{E}\left[e^{\lambda Z_i}\right] \leq e^{\frac{\lambda^2 (b_i - a_i)^2}{8}}$
- Optimization in λ

- Prop:**

$$\mathbb{E}\left[e^{\lambda \sum_{i=1}^n Z_i}\right] \leq e^{\frac{\lambda^2 \sum_{i=1}^n (b_i - a_i)^2}{8}}.$$

Theorem

- Let Z_i be a sequence of independent centered random variables supported in $[a_i, b_i]$ then

$$\mathbb{P}\left(\sum_{i=1}^n Z_i \geq \epsilon\right) \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}}$$

- $Z_i = \frac{1}{n} \left(\mathbb{E}[\ell^{0/1}(Y, f(\underline{X}))] - \ell^{0/1}(Y_i, f(\underline{X}_i)) \right)$
- $\mathbb{E}[Z_i] = 0$ and $Z_i \in \left[\frac{1}{n} \left(\mathbb{E}[\ell^{0/1}(Y, f(\underline{X}))] - 1 \right), \frac{1}{n} \mathbb{E}[\ell^{0/1}(Y, f(\underline{X}))] \right]$
- Concentration:

$$\mathbb{P}(\mathcal{R}(f) - \mathcal{R}_n(f) \geq \epsilon) \leq e^{-2n\epsilon^2}$$

- By symmetry,

$$\mathbb{P}(\mathcal{R}_n(f) - \mathcal{R}(f) \geq \epsilon) \leq e^{-2n\epsilon^2}$$

- Combining the two yields

$$\mathbb{P}(|\mathcal{R}_n(f) - \mathcal{R}(f)| \geq \epsilon) \leq 2e^{-2n\epsilon^2}$$

Concentration

- If \mathcal{S} is finite of cardinality $|\mathcal{S}|$,

$$\mathbb{P} \left(\sup_f (\mathcal{R}(f) - \mathcal{R}_n(f)) \leq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} \right) \geq 1 - \delta$$

$$\mathbb{P} \left(\sup_f |\mathcal{R}_n(f) - \mathcal{R}(f)| \leq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} \right) \geq 1 - 2\delta$$

- Control of the supremum by a quantity depending on the cardinality and the probability parameter δ .
- Simple combination of Hoeffding and a union bound.

PAC Bounds

- If \mathcal{S} is finite of cardinality $|\mathcal{S}|$, with proba greater than $1 - 2\delta$

$$\begin{aligned}\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) &\leq \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}} \\ &\leq 2\sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}}\end{aligned}$$

- If \mathcal{S} is finite of cardinality $|\mathcal{S}|$, with proba greater than $1 - \delta$, simultaneously $\forall f' \in \mathcal{S}$,

$$\begin{aligned}\mathcal{R}(f') &\leq \mathcal{R}_n(f') + \sqrt{\frac{\log |\mathcal{S}| + \log(1/\delta)}{2n}} \\ &\leq \mathcal{R}_n(f') + \sqrt{\frac{\log |\mathcal{S}|}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}\end{aligned}$$

PAC Bounds

- If \mathcal{S} is finite of cardinality $|\mathcal{S}|$, with proba greater than $1 - 2\delta$

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_{\mathcal{S}}^*) \leq \sqrt{\frac{\log |\mathcal{S}|}{2n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- If \mathcal{S} is finite of cardinality $|\mathcal{S}|$, with proba greater than $1 - \delta$, simultaneously $\forall f' \in \mathcal{S}$,

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{\log |\mathcal{S}|}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- Risk increases with the cardinality of \mathcal{S} .
- Similar issue in cross-validation!
- No direct extension for an infinite \mathcal{S} ...

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

- Supremum of Empirical losses:

$$\begin{aligned}\Delta_n(\mathcal{S})(\underline{X}_1, \dots, \underline{X}_n) &= \sup_{f \in \mathcal{S}} \mathcal{R}(f) - \mathcal{R}_n(f) \\ &= \sup_{f \in \mathcal{S}} \left(\mathbb{E} \left[\ell^{0/1}(Y, f(\underline{X})) \right] - \frac{1}{n} \sum_{i=1}^n \ell^{0/1}(Y_i, f(\underline{X}_i)) \right)\end{aligned}$$

Properties

- Bounded difference:

$$|\Delta_n(\mathcal{S})(\underline{X}_1, \dots, \underline{X}_i, \dots, \underline{X}_n) - \Delta_n(\mathcal{S})(\underline{X}_1, \dots, \underline{X}'_i, \dots, \underline{X}_n)| \leq 1/n$$

Concentration

$$\mathbb{P}(\Delta_n(\mathcal{S}) - \mathbb{E}[\Delta_n(\mathcal{S})] \leq \epsilon) \geq 1 - e^{-2n\epsilon^2}$$

- Concentration of bounded difference function.
- Generalization of Hoeffding theorem: McDiarmid Theorem.

Bounded difference function

- $g : \mathcal{X}^n \rightarrow \mathbb{R}$ is a bounded difference function if it exist c_i such that

$$\forall (\underline{X}_i)_{i=1}^n, (\underline{X}'_i)_{i=1}^n \in \mathbb{R},$$

$$|g(\underline{X}_1, \dots, \underline{X}_i, \dots, \underline{X}_n) - g(\underline{X}_1, \dots, \underline{X}'_i, \dots, \underline{X}_n)| \leq c_i$$

Theorem

- If g is a bounded difference function and \underline{X}_i are independent random variables then

$$\mathbb{P}(g(\underline{X}_1, \dots, \underline{X}_n) - \mathbb{E}[g(\underline{X}_1, \dots, \underline{X}_n)] \geq \epsilon) \leq e^{\sum_{i=1}^n \frac{-2\epsilon^2}{c_i^2}}$$

$$\mathbb{P}(\mathbb{E}[g(\underline{X}_1, \dots, \underline{X}_n)] - g(\underline{X}_1, \dots, \underline{X}_n) \geq \epsilon) \leq e^{\sum_{i=1}^n \frac{-2\epsilon^2}{c_i^2}}$$

- Proof ingredients:
 - Chernov bounds
 - Martingale decomposition...

Theorem

- If g is a bounded difference function and \underline{X}_i are independent random variables then

$$\mathbb{P}(g(\underline{X}_1, \dots, \underline{X}_n) - \mathbb{E}[g(\underline{X}_1, \dots, \underline{X}_n)] \geq \epsilon) \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n c_i^2}}$$

- Using $g = \Delta_n(\mathcal{S})$ for which $c_i = 1/n$ yields immediately

$$\mathbb{P}(\Delta_n(\mathcal{S}) - \mathbb{E}[\Delta_n(\mathcal{S})] \geq \epsilon) \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n c_i^2}} = e^{-2n\epsilon^2}$$

- We derive then

$$\mathbb{P}(\Delta_n(\mathcal{S}) \geq \mathbb{E}[\Delta_n(\mathcal{S})] + \epsilon) \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n c_i^2}} = e^{-2n\epsilon^2}$$

- It remains to upperbound

$$\mathbb{E}[\Delta_n] = \mathbb{E} \left[\sup_{f \in \mathcal{S}} \mathcal{R}(f) - \mathcal{R}_n(f) \right]$$

Theorem

- Let σ_i be a sequence of i.i.d. random symmetric Bernoulli variables (Rademacher variables):

$$\mathbb{E} \left[\sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) \right] \leq 2 \mathbb{E} \left[\sup_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \sigma_i \ell^{0/1}(Y_i, f(\underline{X}_i)) \right]$$

Rademacher complexity

- Let $B \subset \mathbf{R}^n$, the Rademacher complexity of B is defined as

$$R_n(B) = \mathbb{E} \left[\sup_{b \in B} \frac{1}{n} \sum_{i=1}^n \sigma_i b_i \right]$$

- Theorem gives an upper bound of the expectation in terms of the **average Rademacher complexity of the random set**
 $B_n(\mathcal{S}) = \{(\ell^{0/1}(Y_i, f(\underline{X}_i)))_{i=1}^n, f \in \mathcal{S}\}$.
- Back to finite setting:** This set is at most of cardinality 2^n .

Theorem

- If B is finite and such that $\forall b \in B, \frac{1}{n} \|b\|_2^2 \leq M^2$, then

$$R_n(B) = \mathbb{E} \left[\sup_{b \in B} \frac{1}{n} \sum_{i=1}^n \sigma_i b_i \right] \leq \sqrt{\frac{2M^2 \log |B|}{n}}$$

- If $B = B_n(\mathcal{S}) = \{(\ell^{0/1}(Y_i, f(\underline{X}_i)))_{i=1}^n, f \in \mathcal{S}\}$, we have $M = 1$ and thus

$$R_n(B) \leq \sqrt{\frac{2 \log |B_n(\mathcal{S})|}{n}}$$

- We obtain immediately

$$\mathbb{E} \left[\sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) \right] \leq \mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right].$$

Theorem

- With probability greater than $1 - 2\delta$,

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right] + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- With probability greater than $1 - \delta$, simultaneously $\forall f' \in \mathcal{S}$

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right] + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- This is a direct consequence of the previous bound.

Corollary

- If \mathcal{S} is finite then with probability greater than $1 - 2\delta$

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_{\mathcal{S}}^*) \leq \sqrt{\frac{8 \log |\mathcal{S}|}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- If \mathcal{S} is finite then with probability greater than $1 - \delta$, simultaneously $\forall f' \in \mathcal{S}$

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{8 \log |\mathcal{S}|}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- It suffices to notice that

$$|B_n(\mathcal{S})| = |\{(\ell^{0/1}(Y_i, f(\underline{X}_i)))_{i=1}^n, f \in \mathcal{S}\}| \leq |\mathcal{S}|$$

- Same result with Hoeffding but with **better** constants!

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sqrt{\frac{\log |\mathcal{S}|}{2n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{\log |\mathcal{S}|}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- Difference due to the *crude* upperbound of

$$\mathbb{E} \left[\sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) \right]$$

- **Why bother?:** We do not have to assume that \mathcal{S} is finite!

$$|B_n(\mathcal{S})| \leq 2^n$$

Outline



- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

Theorem

$$\mathbb{E} \left[\sup_{f \in \mathcal{S}} (\mathcal{R}(f) - \mathcal{R}_n(f)) \right] \leq \mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right]$$

- Key quantity: $\mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right]$
- Hard to control due to its structure!

A first data dependent upperbound

$$\mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right] \leq \sqrt{\frac{8 \log \mathbb{E}[|B_n(\mathcal{S})|]}{n}} \quad (\text{Jensen})$$

- Depends on the unknown P !

Shattering Coefficient (or Growth Function)

- The shattering coefficient of the class \mathcal{S} , $s(\mathcal{S}, n)$, is defined as

$$s(\mathcal{S}, n) = \sup_{((\underline{X}_1, Y_1), \dots, (\underline{X}_n, Y_n)) \in (\mathcal{X} \times \{-1, 1\})^n} |\{(\ell^{0/1}(Y_i, f(\underline{X}_i)))_{i=1}^n, f \in \mathcal{S}\}|$$

- By construction, $|B_n(\mathcal{S})| \leq s(\mathcal{S}, n) \leq \min(2^n, |\mathcal{S}|)$.

A data independent upperbound

$$\mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S})|}{n}} \right] \leq \sqrt{\frac{8 \log s(\mathcal{S}, n)}{n}}$$

Theorem

- With probability greater than $1 - 2\delta$,

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sqrt{\frac{8 \log s(\mathcal{S}, n)}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- With probability greater than $1 - \delta$, simultaneously $\forall f' \in \mathcal{S}$,

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{8 \log s(\mathcal{S}, n)}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- Depends only on the class \mathcal{S} !

VC Dimension

- The VC dimension d_{VC} of \mathcal{S} is defined as the largest integer d such that
$$s(\mathcal{S}, d) = 2^d$$

- The VC dimension can be infinite!

VC Dimension and Dimension

- **Prop:** If $\text{span}(\mathcal{S})$ corresponds to the sign of functions in a linear space of dimension d then $d_{VC} \leq d$.
- VC dimension similar to the usual dimension.

Sauer's Lemma

- If the VC dimension d_{VC} of \mathcal{S} is finite

$$s(\mathcal{S}, n) \leq \begin{cases} 2^n & \text{if } n \leq d_{VC} \\ \left(\frac{en}{d_{VC}}\right)^{d_{VC}} & \text{if } n > d_{VC} \end{cases}$$

- **Cor.:** $\log s(\mathcal{S}, n) \leq d_{VC} \log \left(\frac{en}{d_{VC}}\right)$ if $n > d_{VC}$.

PAC Bounds

- If \mathcal{S} is of VC dimension d_{VC} then if $n > d_{VC}$
- With probability greater than $1 - 2\delta$,

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sqrt{\frac{8d_{VC} \log\left(\frac{en}{d_{VC}}\right)}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- With probability greater than $1 - \delta$, simultaneously $\forall f' \in \mathcal{S}$,

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{8d_{VC} \log\left(\frac{en}{d_{VC}}\right)}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- **Rk:** If $d_{VC} = +\infty$ no uniform PAC bounds exists!

Outline

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View

PAC Bounds

- Let $\pi_f > 0$ such that $\sum_{f \in \mathcal{S}} \pi_f = 1$
- With proba greater than $1 - 2\delta$,

$$\mathcal{R}(\hat{f}) - \mathcal{R}(f_S^*) \leq \sqrt{\frac{\log(1/\pi_f)}{2n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- With proba greater than $1 - \delta$, simultaneously $\forall f' \in \mathcal{S}$,

$$\mathcal{R}(f') \leq \mathcal{R}_n(f') + \sqrt{\frac{\log(1/\pi_f)}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

- Very similar proof than the uniform one!
- Much more interesting idea when combined with several models...

- Assume we have a countable collection of set $(\mathcal{S}_m)_{m \in \mathcal{M}}$ and let π_m be such that $\sum_{m \in \mathcal{M}} \pi_m = 1$.

Non Uniform Risk Bound

- With probability $1 - \delta$, simultaneously for all $m \in \mathcal{M}$ and all $f \in \mathcal{S}_m$,

$$\mathcal{R}(f) \leq \mathcal{R}_n(f) + \mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

Structural Risk Minimization

- Choose \hat{f} as the minimizer over $m \in \mathcal{M}$ and $f \in \mathcal{S}_m$ of

$$\mathcal{R}_n(f) + \mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}}$$

- Mimics the minimization of the integrated risk!

PAC Bound

- If \hat{f} is the SRM minimizer then with probability $1 - 2\delta$,

$$\mathcal{R}(\hat{f}) \leq \inf_{m \in \mathcal{M}} \inf_{f \in \mathcal{S}_m} \left(\mathcal{R}(f) + \mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}} \right) + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

- The SRM minimizer balances the risk $\mathcal{R}(f)$ and the upper bound on the estimation error $\mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right] + \sqrt{\frac{\log(1/\pi_m)}{2n}}$.
- $\mathbb{E} \left[\sqrt{\frac{8 \log |B_n(\mathcal{S}_m)|}{n}} \right]$ can be replaced by an upper bound (for instance a VC based one)...

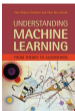
- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5
 - Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 8
 - Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
 - Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
 - Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
 - Applications to Text
 - References
- 8 **Statistical Learning: PAC-Bayesian Approach and Complexity Theory**
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - **References**
- 9 References



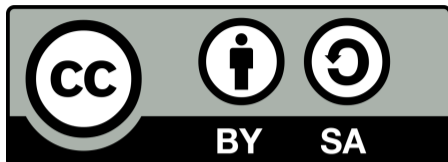
Ch. Giraud.
Introduction to High-Dimensional Statistics.
CRC Press, 2014



M. Mohri, A. Rostamizadeh, and A. Talwalkar.
Foundations of Machine Learning (2nd ed.)
MIT Press, 2018



S. Shalev-Shwartz and S. Ben-David.
Understanding Machine Learning.
Cambridge University Press, 2014



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.

- 1 Statistical Learning: Introduction, Setting and Risk Estimation
 - Introduction
 - Machine Learning
 - Supervised Learning
 - Risk Estimation and Model Selection
 - Cross Validation and Test
 - References
- 2 ML Methods: Probabilistic Point of View
 - Motivation
 - Supervised Learning
 - A Probabilistic Point of View
 - Parametric Conditional Density Modeling
 - Non Parametric Conditional Density Modeling
 - Generative Modeling
 - Model Selection
 - Penalization
- 3 ML Methods: Optimization Point of View
 - Supervised Learning
 - Optimization Point of View
 - SVM
 - Penalization
 - Cross Validation and Weights
- 4 Optimization: Gradient Descent Algorithms
 - Introduction
 - Gradient Descent
 - Proximal Descent
 - Coordinate Descent
 - Gradient Descent Acceleration
 - Stochastic Gradient Descent
- 5 Gradient Descent Step
 - Non-Convex Setting
 - References
- 5 ML Methods: Neural Networks and Deep Learning
 - Introduction
 - From Logistic Regression to NN
 - NN Optimization
 - NN Regularization
 - Image and CNN
 - Text, Recurrent Neural Networks and Transformers
 - NN Architecture
 - References
- 6 ML Methods: Trees and Ensemble Methods
 - Trees
 - Bagging and Random Forests
 - Bootstrap and Bagging
 - Randomized Rules and Random Forests
 - Boosting
 - AdaBoost as a Greedy Scheme
 - Boosting
 - Ensemble Methods
 - References
- 7 Unsupervised Learning: Dimension Reduction and Clustering
 - Unsupervised Learning?
 - A First Glimpse
 - Clustering
 - Dimensionality Curse
 - Dimension Reduction
 - Generative Modeling
- 6 Dimension Reduction
 - Simplification
 - Reconstruction Error
 - Relationship Preservation
 - Comparing Methods?
- 6 Clustering
 - Prototype Approaches
 - Contiguity Approaches
 - Agglomerative Approaches
 - Other Approaches
 - Scalability
- 6 Generative Modeling
 - (Plain) Parametric Density Estimation
 - Latent Variables
 - Approximate Simulation
 - Diffusion Model
 - Generative Adversarial Network
- 6 Applications to Text
 - References
- 8 Statistical Learning: PAC-Bayesian Approach and Complexity Theory
 - Supervised Learning
 - Empirical Risk Minimization
 - Empirical Risk Minimization
 - ERM and PAC Analysis
 - Hoeffding and Finite Class
 - McDiarmid and Rademacher Complexity
 - VC Dimension
 - Structural Risk Minimization
 - References
- 9 References



T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning.
Springer Series in Statistics, 2009



M. Mohri, A. Rostamizadeh, and A. Talwalkar.
Foundations of Machine Learning (2nd ed.)
MIT Press, 2018



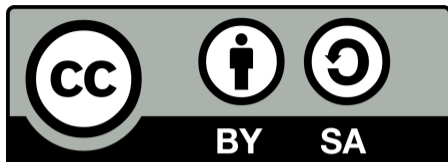
A. Géron.
Hands-On Machine Learning with Scikit-Learn, Keras and TensorFlow (3rd ed.)
O'Reilly, 2022



Ch. Giraud.
Introduction to High-Dimensional Statistics.
CRC Press, 2014



S. Bubeck.
Convex Optimization: Algorithms and Complexity.
Now Publisher, 2015



Creative Commons Attribution-ShareAlike (CC BY-SA 4.0)

- You are free to:
 - **Share:** copy and redistribute the material in any medium or format
 - **Adapt:** remix, transform, and build upon the material for any purpose, even commercially.
- Under the following terms:
 - **Attribution:** You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
 - **ShareAlike:** If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
 - **No additional restrictions:** You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.