# Statistical and Optimization Approaches in Classification
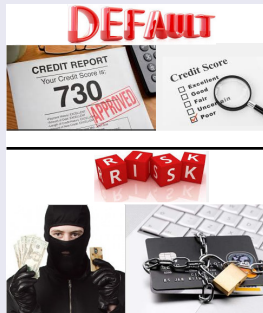
A. Fermin - MODAL'X, Université Paris Ouest

E. Le Pennec - CMAP, École polytechnique

université
**Paris Ouest**
Nanterre La Défense

*l'* X
ÉCOLE
**POLYTECHNIQUE**
UNIVERSITÉ PARIS-SACLAY

RTE - 30/06/2015

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



- Data: Client profile, Client credit history...
- Input: Client profile
- Output: Credit risk
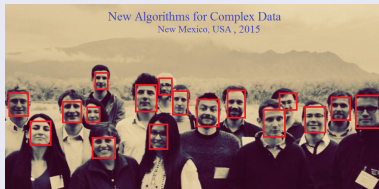
### Marketing: advertisement, recommendation...



- Data: User profile, Web site history...
- Input: User profile, Current web page
- Output: Advertisement with price, recommendation...
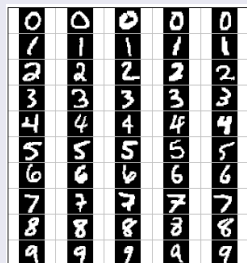
### Spam detection (Text classification)



- Data: email collection
- Input: email
- Output : Spam or No Spam

## Face Detection



- Data: Annotated database of images
- Input : Sub window in the image
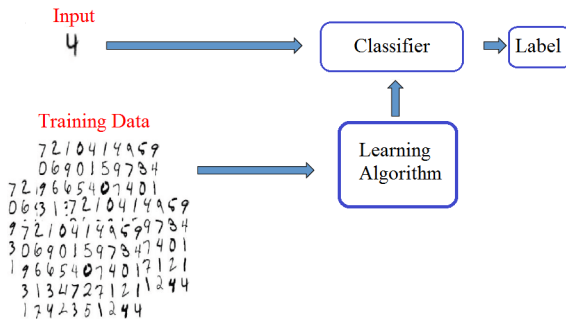- Output : Presence or no of a face...

## Number Recognition



- Data: Annotated database of images (each image is represented by a vector of $28 \times 28 = 784$ pixel intensities)
- Input: Image
- Output: Corresponding number

Classifier

Label

Training Data

Learning Algorithm

### A definition by Tom Mitchell (http://www.cs.cmu.edu/~tom/)

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.

## Big Data, Data Science and Machine Learning

- **Big Data**: buzzword to raise money (or data sets too large or too complex to be handled by the current system)
- **Data Science**: art (or science) of the generalizable extraction of knowledge from data.
- **Machine Learning**: construction and study of algorithms that can learn from and make predictions on data.

- Exciting challenges in the industrial **and** the academic worlds.

## Machine Learning

- Fundamental ingredient in data science.
- Necessity for a Data Scientist to understand the principle of the simplest methods to grasp the more sophisticated ones.

## Supervised Learning Framework

- Input measurement $\mathbf{X} = (X^{(1)}, X^{(2)}, \ldots, X^{(d)}) \in \mathcal{X}$

- Output measurement $Y \in \mathcal{Y}$.

- $(\mathbf{X}, Y) \sim \mathbf{P}$ with $\mathbf{P}$ unknown.

- Training data : $\mathcal{D}_n = \{(\mathbf{X}_1, Y_1), \ldots, (\mathbf{X}_n, Y_n)\}$ (i.i.d. $\sim \mathbf{P}$)

- Often
  - $\mathbf{X} \in \mathbb{R}^d$ and $Y \in \{-1, 1\}$ (classification)
  - or $\mathbf{X} \in \mathbb{R}^d$ and $Y \in \mathbb{R}$ (regression).

- A classifier is a function in $\mathcal{F} = \{f : \mathcal{X} \to \mathcal{Y} \text{ measurable}\}$

## Goal

- Construct a good classifier $\widehat{f}$ from the training data.

- Need to specify the meaning of good.
- Formally, classification and regression are the same problem!

## Loss function

- **Loss function** : $\ell(f(x), y)$ measure how well $f(x)$ "predicts" $y$.
- Examples:
  - Prediction loss: $\ell(Y, f(\mathbf{X})) = \mathbf{1}_{Y \neq f(\mathbf{X})}$
  - Quadratic loss: $\ell(Y, \mathbf{X}) = |Y - f(\mathbf{X})|^2$

## Risk of a generic classifier

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

$$\mathcal{R}(f) = \mathbb{E}\left[\ell(Y, f(\mathbf{X}))\right] = \mathbb{E}_X\left[\mathbb{E}_{Y|\mathbf{x}}\left[\ell(Y, f(\mathbf{X}))\right]\right]$$

- Examples:
  - Prediction loss: $\mathbb{E}\left[\ell(Y, f(\mathbf{X}))\right] = \mathbb{P}\{Y \neq f(\mathbf{X})\}$
  - Quadratic loss: $\mathbb{E}\left[\ell(Y, f(\mathbf{X}))\right] = \mathbb{E}\left[|Y - f(\mathbf{X})|^2\right]$

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

## Experience, Task and Performance measure

- Training data : $\mathcal{D} = \{(\mathbf{X}_1, Y_1), \ldots, (\mathbf{X}_n, Y_n)\}$ (i.i.d. $\sim \mathbf{P}$)
- Predictor: $f : \mathcal{X} \to \mathcal{Y}$ measurable
- Cost/Loss function : $\ell(f(\mathbf{X}), Y)$ measure how well $f(\mathbf{X})$ "predicts" $Y$
- Risk:

$$\mathcal{R}(f) = \mathbb{E}\left[\ell(Y, f(\mathbf{X}))\right] = \mathbb{E}_X\left[\mathbb{E}_{Y|\mathbf{X}}\left[\ell(Y, f(\mathbf{X}))\right]\right]$$

- Often $\ell(f(\mathbf{X}), Y) = |f(\mathbf{X}) - Y|^2$ or $\ell(f(\mathbf{X}), Y) = \mathbf{1}_{Y \neq f(\mathbf{X})}$

## Goal

- Learn a rule to construct a classifier $\widehat{f} \in \mathcal{F}$ from the training data $\mathcal{D}_n$ s.t. the risk $\mathcal{R}(\widehat{f})$ is small on average or with high probability with respect to $\mathcal{D}_n$.

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

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

## Bayes Classifier (explicit solution)

- In binary classification with $0 - 1$ loss:

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

- In regression with the quadratic loss

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

Issue: Explicit solution requires to know $\mathbb{E}\left[Y|\mathbf{X}\right]$ for all values of $\mathbf{X}$!

## Machine Learning

- Learn a rule to construct a classifier $\widehat{f} \in \mathcal{F}$ from the training data $\mathcal{D}_n$ s.t. the risk $\mathcal{R}(\widehat{f})$ is small on average or with high probability with respect to $\mathcal{D}_n$.

## 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

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

- Examples:
  - Linear regression
  - Linear discrimination with
    $$\mathcal{S} = \{\mathbf{x} \mapsto \texttt{sign}\{\beta^T \mathbf{x} + \beta_0\} / \beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}\}$$

## 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.

# Supervised Learning
## Example: More complex model

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



**Approximation error and estimation error (Bias/Variance)**

$$\mathcal{R}(\widehat{f}_{\mathcal{S}}) - \mathcal{R}(f^*) = \underbrace{\mathcal{R}(f_{\mathcal{S}}^*) - \mathcal{R}(f^*)}_{\text{Approximation error}} + \underbrace{\mathcal{R}(\widehat{f}_{\mathcal{S}}) - \mathcal{R}(f_{\mathcal{S}}^*)}_{\text{Estimation error}}$$

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

Agnostic approach

- No assumption (so far) on the law of $(\mathbf{X}, Y)$.

- 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 contains a good ideal target but the estimation error ("variance") can be large (Over-fit)

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

How to find a good function $f$ with a *small* risk
$$R(f) = \mathbb{E}\left[\ell(Y, f(X))\right] \quad ?$$
Canonical approach: $\widehat{f}_{\mathcal{S}} = \operatorname{argmin}_{f \in \mathcal{S}} \frac{1}{n} \sum_{i=1}^{n} \ell(Y_i, f(\mathbf{X}_i))$

### Problems

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

### A Statistical Point of View

**Solution:** For $\mathbf{X}$, estimate $Y|\mathbf{X}$ 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 $\ell$ by an upper bound $\ell'$ and minimize the empirical loss: SVR, SVM, Neural Network, Tree, Boosting

- Input: a data set $\mathcal{D}_n$
  Learn $Y|x$ or equivalently $p_k(\mathbf{x}) = \mathbb{P}\{Y = k | \mathbf{X} = \mathbf{x}\}$ (using the data set) and plug this estimate in the Bayes classifier

- Output: a classifier $\widehat{f} : \mathbb{R}^d \to \{-1, 1\}$

$$\hat{f}(\mathbf{x}) = \begin{cases} +1 & \text{if } \widehat{p}_{+1}(\mathbf{x}) \geq \widehat{p}_{-1}(\mathbf{x}) \\ -1 & \text{otherwise} \end{cases}$$

- Three instantiations:
  1. Generative Modeling (Bayes method)
  2. Logistic modeling (parametric method)
  3. Nearest neighbors (kernel method)

## Bayes formula

$$p_k(\mathbf{x}) = \frac{\mathbb{P}\{\mathbf{X} = \mathbf{x} | Y = k\} \, \mathbb{P}\{Y = k\}}{\mathbb{P}\{\mathbf{X} = \mathbf{x}\}}$$

**Remark**: If one knows the law of $(X, Y)$ or equivalently of $X$ given $y$ and of $Y$ then everything is easy!

- Binary Bayes classifier (the best solution)

$$f^*(\mathbf{x}) = \begin{cases} +1 & \text{if } p_{+1}(\mathbf{x}) \geq p_{-1}(\mathbf{x}) \\ -1 & \text{otherwise} \end{cases}$$

- **Heuristic**: Estimate those quantities and plug the estimations.
- By using different models for $\mathbb{P}\{\mathbf{X} | Y\}$, we get different classifiers.
- **Remark:** You can also use your favorite density estimator...

**Discriminant Analysis (Gaussian model)**

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

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

- Discriminants fonctions:

$$g_k(\mathbf{x}) = \ln(\mathbb{P}\{X|Y=k\}) + \ln(\mathbb{P}\{Y=k\})$$

$$g_k(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mu_k)^t \Sigma_k^{-1}(\mathbf{x} - \mu_k)$$

$$-\frac{d}{2}\ln(2\pi) - \frac{1}{2}\ln(|\Sigma_k|) + \ln(\mathbb{P}\{Y=k\})$$

- QDA (differents $\Sigma_k$ in each class) and LDA ($\Sigma_k = \Sigma$ for all $k$)

Beware: this model can be false but the methodology remains valid!

## Estimation

In pratice, we will need to estimate $\mu_k$, $\Sigma_k$ and $\mathbb{P}_k := \mathbb{P}\{Y = k\}$

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

- DA classifier

$$\widehat{f_G}(\mathbf{x}) = \begin{cases} +1 & \text{if } \widehat{g}_{+1} \geq \widehat{g}_{-1} \\ -1 & \text{otherwise} \end{cases}$$

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

## Naive Bayes

- Classical algorithm using a crude modeling for $\mathbb{P}\left\{X|Y\right\}$:
  - Feature independence assumption:

$$\mathbb{P}\left\{X|Y\right\} = \prod_{i=1}^{d} \mathbb{P}\left\{X^{(i)}\middle|Y\right\}$$

  - 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!

Naive Bayes with Gaussian model

Naive Bayes with kernel density estimates

- Direct modeling of $Y|x$.

### The Binary logistic model ($Y \in \{-1, 1\}$)

$$p_{+1}(\mathbf{x}) = \frac{e^{\beta^t \varphi(\mathbf{x})}}{1 + e^{\beta^t \varphi(\mathbf{x})}}$$

where $\varphi(x)$ is a transformation of the individual $\mathbf{x}$

- In this model, one verifies that
$$p_{+1}(\mathbf{x}) \geq p_{-1}(\mathbf{x}) \quad \Leftrightarrow \quad \beta^t \varphi(\mathbf{x}) \geq 0$$
- True $Y|x$ may not belong to this model $\Rightarrow$ maximum likelihood of $\beta$ only finds a good approximation!
- Binary Logistic classifier:
$$\widehat{f}_L(\mathbf{x}) = \begin{cases} +1 & \text{if } \widehat{\beta}^t \varphi(\mathbf{x}) \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

where $\widehat{\beta}$ is estimated by maximum likelihood.

- Logistic model: approximation of $\mathcal{B}(p_1(\mathbf{x}))$ by $\mathcal{B}(h(\beta^t\varphi(\mathbf{x})))$ with $h(t) = \frac{e^t}{1+e^t}$.

### Opposite of the log-likelihood formula

$$- \frac{1}{n} \sum_{i=1}^{n} \left(\mathbf{1}_{y_i=1} \log(h(\beta^t\varphi(\mathbf{x}))) + \mathbf{1}_{y_i=-1} \log(1 - h(\beta^t\varphi(\mathbf{x})))\right)$$

$$= -\frac{1}{n} \sum_{i=1}^{n} \left(\mathbf{1}_{y_i=1} \log \frac{e^{\beta^t\varphi(\mathbf{x})}}{1 + e^{\beta^t\varphi(\mathbf{x})}} + \mathbf{1}_{y_i=-1} \log \frac{1}{1 + e^{\beta^t\varphi(\mathbf{x})}}\right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \log \left(1 + e^{-y_i(\beta^t\varphi(\mathbf{x}))}\right)$$

- Convex function in $\beta$!
- **Remark:** You can also use your favorite parametric model instead of the logistic one...

Quadratic Logistic

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

### *k*-NN as local conditional density estimate

$$\widehat{p}_{+1}(\mathbf{x}) = \frac{\sum_{\mathbf{x}_i \in \mathcal{V}_\mathbf{x}} \mathbf{1}_{\{y_i = +1\}}}{|\mathcal{V}_\mathbf{x}|}$$

- KNN Classifier:
$$\widehat{f}_{KNN}(\mathbf{x}) = \begin{cases} +1 & \text{if } \widehat{p}_{+1}(\mathbf{x}) \geq \widehat{p}_{-1}(\mathbf{x}) \\ -1 & \text{otherwise} \end{cases}$$

- **Remark:** You can also use your favorite kernel estimator...

k-NN with k=5

k-NN with k=9

k-NN with k=13

k-NN with k=17

k-NN with k=25

**Error behaviour**

- Learning/training error (error made on the learning/training set) decays when the complexity of the model increases.
- Quite different behavior when the error is computed on new observations (generalization error).

- Overfit for complex models: parameters learned are too specific to the learning set!
- General situation! (Think of polynomial fit...)
- Need to use an other criterion than the training error!

Training Set          Test Set

- **Very simple idea:** use a second learning/verification set to compute a verification error.
- Sufficient to avoid over-fitting!

### Cross Validation

- Use $\frac{V-1}{V}n$ observations to train and $\frac{1}{V}n$ to verify!
- Validation for a learning set of size $(1 - \frac{1}{V}) \times n$ instead of $n$!

- Most classical variations:
    - Leave One Out,
    - $V$-fold cross validation.
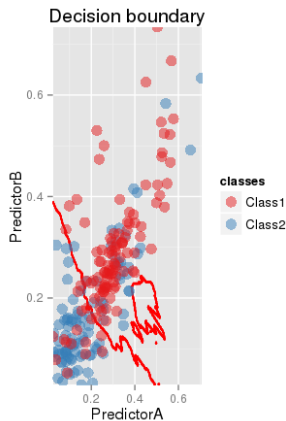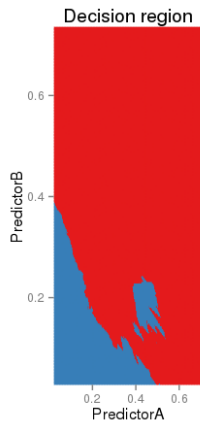- Accuracy/Speed tradeoff: $V = 5$ or $V = 10$!

# A Statistical Point of View
## Cross Validation

k-NN with k=25

How to find a good function $f$ with a *small* risk
$$R(f) = \mathbb{E}\left[\ell(Y, f(X))\right] \quad ?$$

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

### Problems

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

### A Statistical Point of View

**Solution:** For $\mathbf{X}$, estimate $Y|\mathbf{X}$ 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 $\ell$ by an upper bound $\ell'$ and minimize the empirical loss: SVR, SVM, Neural Network, Tree, Boosting

- The best solution $f^*$ is the one minimizing

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

### 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 empirical loss

$$\widehat{f} = f_{\widehat{\theta}} = \underset{f_\theta, \theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$$

- Plus convexification/regularization of the risk...
- Examples: SVM, (Deep) Neural Networks...

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

### Classical convexification

- Logistic loss: $\ell(y, f(x)) = \log(1 + e^{-yf(x)})$ (Logistic / NN)
- Hinge loss: $\ell(y, f(x)) = (1 - yf(x))_+$ (SVM)
- Exponential loss: $\ell(y, f(x)) = e^{-yf(x)}$ (Boosting...)

- Ideal solution:

$$\widehat{f} = \underset{f \in \mathcal{S}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} \ell^{0/1}(y_i, f(x_i))$$

### Logistic regression

- Use $f(x) = \langle \beta, x \rangle + b$.
- Use the logistic loss $\ell(y, f) = \log_2(1 + e^{-yf})$, i.e. the -log-likelihood.

- Different vision than the statistician but same algorithm!

Logistic

- Linear classifier: $\text{sign}(\langle \beta, x \rangle + b)$
- Separable case: $\exists (\beta, b), \forall i, y_i(\langle \beta, x \rangle + b) > 0$!

How to choose $(\beta, b)$ so that the separation is maximal?

- Strict separation: $\exists (\beta, b), \forall i, y_i(\langle \beta, x \rangle + b) \geq 1$
- Maximize the distance between $\langle \beta, x \rangle + b = 1$ and $\langle \beta, x \rangle + b = -1$.
- Equivalent to the minimization of $\|\beta\|^2$.

- What about the non separable case?
- Relax the assumption that $\forall i, y_i(\langle \beta, x \rangle + b) \geq 1$.
- Naive attempt:

$$\operatorname{argmin} \|\beta\|^2 + C\frac{1}{n}\sum_{i=1}^{n} \mathbf{1}_{y_i(\langle \beta, x \rangle + b) \leq 1}$$

- Non convex minimization.

**SVM: better convex relaxation!**

$$\operatorname{argmin} \|\beta\|^2 + C\frac{1}{n}\sum_{i=1}^{n} \max(1 - y_i(\langle \beta, x \rangle + b), 0)$$

- Convex relaxation:

$$\text{argmin} \|\beta\|^2 + C\frac{1}{n}\sum_{i=1}^{n}\max(1 - y_i(\langle\beta, x\rangle + b), 0)$$

$$= \text{argmin} \frac{1}{n}\sum_{i=1}^{n}\max(1 - y_i(\langle\beta, x\rangle + b), 0) + \frac{1}{C}\|\beta\|^2$$

- **Prop:** $\ell^{0/1}(y_i, \text{sign}(\langle\beta, x\rangle + b)) \leq \max(1 - y_i(\langle\beta, x\rangle + b), 0)$

---

Penalized convex relaxation (Tikhonov!)

$$\frac{1}{n}\sum_{i=1}^{n}\ell^{0/1}(y_i, \text{sign}(\langle\beta, x\rangle + b))$$

$$\leq \frac{1}{n}\sum_{i=1}^{n}\max(1 - y_i(\langle\beta, x\rangle + b), 0) + \frac{1}{C}\|\beta\|^2$$

$$\Phi : \mathbb{R}^2 \to \mathbb{R}^3$$
$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2} x_1 x_2, x_2^2)$$

- Non linear separation: just replace $x$ by a non linear $\Phi(x)$...

### Kernel trick

- Computing $k(x, y) = \langle \Phi(x), \Phi(y) \rangle$ may be easier than computing $\Phi(x)$, $\Phi(y)$ and then the scalar product!

- $\Phi$ can be specified through its definite positive kernel $k$.

- Examples: Polynomial kernel $k(x, y) = (1 + \langle x, y \rangle)^d$, Gaussian kernel $k(x, y) = e^{-\|x-y\|^2/2}$,...

- RKHS setting!

- Can be used in (logistic) regression and more...

# An Optimization Point of View
## SVM



Support Vector Machine with polynomial kernel

Support Vector Machine with Gaussian kernel

Activation Neuron Configuration

I = Input
O = Output
B = Bias

Activation Fonction

## Artificial neuron

- Structure:
  - Mix inputs with a weighted sum,
  - Apply a (non linear) activation function to this sum,
  - Eventually 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.

### Neural network structure

- Cascade of artificial neurons organized in layers
- Thresholding decision only at the output layer

- Most classical case use logistic neurons and the -log-likelihood as the criterion to minimize.
- Classical (stochastic) gradient descent algorithm (Back propagation)
- Non convex and thus may be trapped in local minima.

Neural Network

## Deep Neural Network structure

- Deep cascade of layers!

- No conceptual novelty but initialization becomes a crucial issue.

- Bunch of solutions proposed on a greedy initialization of the layers starting from the deepest one.

- Very impressive results!

H2O NN

**Family of Machine Learning algorithm combining:**

- a (deep) multilayered structure,
- a clever (often unsupervised) initalization,
- a more classical final fine tuning optimization.

- Examples: Deep Neural Network, Deep (Restricted) Boltzman Machine, Stacked Encoder...
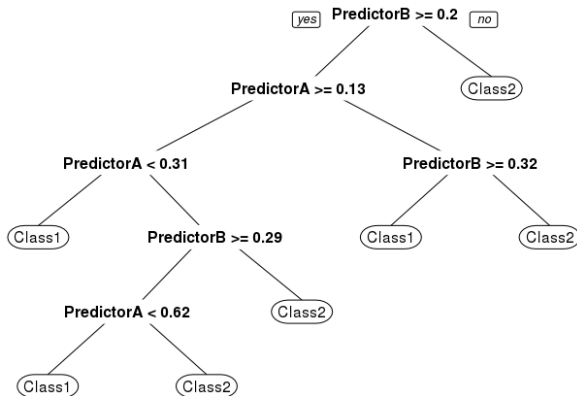- Appears to be very efficient but lack of theoretical fundation!

### Tree principle

- Construction of a recursive partition through a tree structured set of questions (splits around a given value of a variable)
- For a given partition, statistical approach **and** optimization approach yields the same classifier!
- A simple majority vote in each leaf

- Quality of the prediction depends on the tree (the partition).
- Issue: Minim. of the (penalized) empirical error 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...

## Various definition of *homogeneous*

- CART: empirical loss based criterion
$$C(R, \overline{R}) = \sum_{x_i \in R} \ell(y_i, y(R)) + \sum_{x_i \in \overline{R}} \ell(y_i, y(\overline{R}))$$

- CART: Gini index (classification)
$$C(R, \overline{R}) = \sum_{x_i \in R} p(R)(1 - p(R)) + \sum_{x_i \in \overline{R}} p(\overline{R})(1 - p(\overline{R}))$$

- C4.5: entropy based criterion (Information Theory)
$$C(R, \overline{R}) = \sum_{x_i \in R} H(R) + \sum_{x_i \in \overline{R}} H(\overline{R})$$

- CART with Gini is probably the most used technique...
- Other criterion based on $\chi^2$ 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 one minimizing the criterion

- Variations: split at all categories of a categorical variables (ID3), split at a fixed position (median/mean)
- 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!

- Model select. 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.

- Limits over-fitting...

- 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

# Model Selection
## Outline

- Ideal solution:

$$f^*(x) = \arg\max \mathbb{P}\left\{Y|x\right\}$$

### Logistic

- Model $Y|X$ with a logistic model.
- Estimate its parameters with a Maximum Likelihood approach.
- Plug the estimate in the Bayes classifier.

- Model hyperparameters:
  - Features
  - Parametric model...

- Ideal solution:

$$f^*(x) = \arg\max \mathbb{P}\left\{Y|x\right\}$$

**Generative Modeling**

- Estimate $X|Y$ with a density estimator as well as $\mathbb{P}\left\{Y\right\}$
- Deduce using the Bayes formula an estimate $Y|X$.
- Plug the estimate in the Bayes classifier.

- Model hyperparameters:
    - Features
    - Generative model

- Ideal solution:

$$f^*(x) = \arg\max \mathbb{P}\{Y|x\}$$

### Kernel methods

- Estimate $Y|X$ with a kernel conditional density estimator.
- Plug the estimate in the Bayes classifier.

- Model hyperparameters:
  - Features
  - Bandwidth and kernel

- Ideal solution:

$$f^* = \underset{f \in \mathcal{S}}{\operatorname{argmin}} \, \mathbb{E}\left[\ell^{0/1}(Y, f(X))\right]$$

### Logistic

- Replace $\ell^{0/1}$ by the logistic loss.
- Add a penalty $\lambda \|f\|_p$
- Compute the minimizer.

- Model hyperparameters:
    - Features
    - Penalty and regularization parameter.

- Ideal solution:

$$f^* = \underset{f \in \mathcal{S}}{\operatorname{argmin}} \, \mathbb{E}\left[\ell^{0/1}(Y, f(X))\right]$$

### SVM

- Replace the expectation by its empirical counterpart.
- Replace $\ell^{0/1}(y, f) = \mathbf{1}_{y=f}$ by $\ell'(y, f) = (1 - yf)_+$.
- Add a penalty $\lambda\|f\|_{\mathcal{S}}^2$.
- Compute the minimizer.

- Model hyperparameters:
  - Features
  - $\mathcal{S}$ RKHS structure: features mapping and metric
  - Regularization parameters $\lambda$

- Ideal solution:

$$f^* = \underset{f \in \mathcal{S}}{\operatorname{argmin}} \, \mathbb{E}\left[\ell^{0/1}(Y, f(X))\right]$$

### NN

- Neuron: $x \mapsto \sigma(\langle \beta, x \rangle + b)$
- Neural Network: Convolution system of neurons.
- Replace $\ell^{0/1}(y, f)$ by a smooth/convex loss.
- Minimize the empirical loss using the backprop algorithm (gradient descent)

- Model hyperparameters:
  - Features
  - Net architecture, activation function
  - Initialization strategy
  - Optimization strategy (and regularization strategy)

- Ideal solution:

$$f^*(x) = \arg\max \mathbb{P}\{Y|x\} \quad \text{and} \quad f^* = \underset{f \in \mathcal{S}}{\operatorname{argmin}} \, \mathbb{E}\left[\ell^{0/1}(Y, f(X))\right]$$

### Single tree

- Greedy Partition construction.
- Local conditional density estimation / loss minimization.
- Suboptimal tree optimization through a relaxed criterion

### Bagging/Random Forest

- Averaging of several predictors (statistical point of view)

### Boosting

- Best interpretation as a minimization of the exponential loss $\ell(y, f) = e^{-yf}$ (optimization point of view)

## Models

- How to design models? (Model/feature design)
- How to chose among several models? (Model/feature selection)

- Key to obtain good performance!

## Approximation error and estimation error (Bias/Variance)

$$\mathcal{R}(\widehat{f}_{\mathcal{S}}) - \mathcal{R}(f^*) = \underbrace{\mathcal{R}(f_{\mathcal{S}}^*) - \mathcal{R}(f^*)}_{\text{Approximation error}} + \underbrace{\mathcal{R}(\widehat{f}_{\mathcal{S}}) - \mathcal{R}(f_{\mathcal{S}}^*)}_{\text{Estimation error}}$$

- Approximation error can be large for not suitable model $\mathcal{S}$!
- Estimation error can be large if the model is complex!

- Need to find the good balance automatically!

- Empirical error biased toward complex models!

## Selection criterion

- **Cross validation:** Very efficient (and almost always used in practice!) but slightly biased as it target uses only a fraction of the data.

- **Penalization approach:** use empirical loss criterion but penalize it by a term increasing with the complexity of $\mathcal{S}$

$$R_n(\widehat{f_{\mathcal{S}}}) \rightarrow R_n(\widehat{f_{\mathcal{S}}}) + \text{pen}(\mathcal{S})$$

and choose the model with the smallest penalized risk.

- How to combine several predictors (models)?
- Two strategies: mixture or sequential

## Mixture

- Model averaging
- Data dependent model averaging (learn mixture weights)

## Stagewise

- Modify learning procedure according to current results.
- Boosting, Cascade...

- Data Science and Big Data: Much more than a hype!

- **Big data** is an all-encompassing term for any collection of data sets so large and complex that it becomes difficult to process using traditional data processing applications.
- **Data science** is the study of the generalizable extraction of knowledge from data, yet the key word is science.
- **Statistics** is the study of the collection, analysis, interpretation, presentation and organization of data.
- **Machine Learning** explores the construction and the study of algorithms that can learn from and make predictions on data.
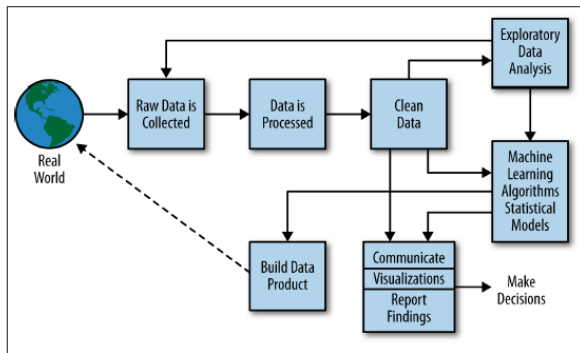
*Figure 2-2. The data science process*

**Doing Data Science: Straight talk from the frontline**

- Rachel Schutt, Cathy O'Neil - O'Reilly
- Art of data driven decision / evaluation.

## Data everywhere

- Huge volume,
- Huge variety...

## Affordable computation units

- Cloud computing
- Graphical Processor Units (GPU)...

- Growing academic and industrial interest!

# Data Science and Big Data
## Big Data is (quite) Easy

## Example of *off the shelves* solution



```scala
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"Training: $numTraining, test: $numTest.")
  examples.unpersist(blocking = false)

  val updater = params.regType match {
    case L1 => new L1Updater()
    case L2 => new SquaredL2Updater()
  }

  val algorithm = new LogisticRegressionWithSGD()
    algorithm.optimizer
      .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



```
export AWS_ACCESS_KEY_ID=<your-access-keyid>
export AWS_SECRET_ACCESS_KEY=<your-access-key-secret>
cellule/spark/ec2/sparl-ec2 -i cellule.pem -k cellule -s <number of machines> launch <cluster-name>
ssh -i cellule.pem root@<your-cluster-master-dns>
spark-ec2/copy-dir ephemeral-hdfs/conf
ephemeral-hdfs/bin/hadoop distcp s3n://celluledecalcul/dataset/raw/train.csv /data/train.csv
scp -i cellule.pem cellule/challenge/target/scala-2.10/target/scala-2.10/challenges_2.10-0.0.jar

cellule/spark/bin/spark-submit \
        --class fr.cc.challenge.Preprocess \
        challenges_2.10-0.0.jar \
        /data/train.csv \
        /data/train2.csv

cellule/spark/bin/spark-submit \
        --class fr.cc.sparktest.LogisticRegression \
        challenges_2.10-0.0.jar \
        /data/train2.csv
```

$\Rightarrow$ Logistic regression for arbitrary large dataset!

More Ideas Based on Your Browsing History

You looked at

You might also consider

Thriving in the Knowledge Age: New... Paperback by John H. Falk
$29.95

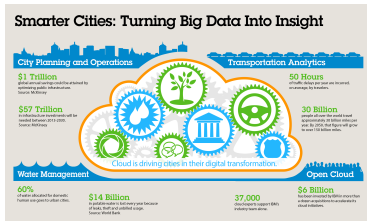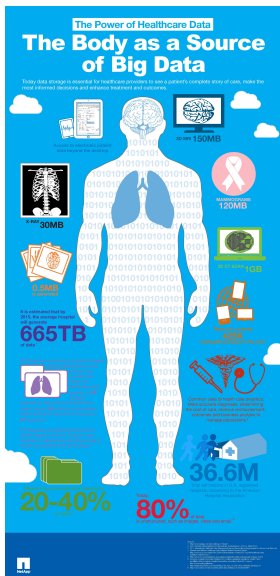Museum Administration: An Introduction Paperback by Hugh H. Genoways
$31.95 $28.75

Exhibit Labels: An Interpretive Approach Paperback by Beverly Serrell
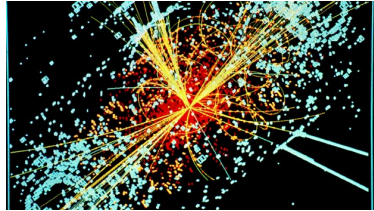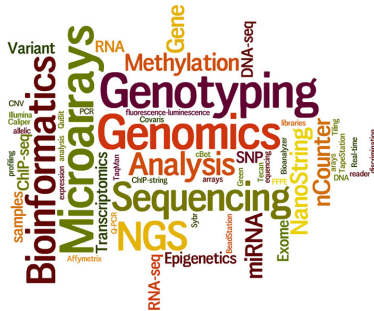$34.95 $27.85

› Find similar items

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

# Data Science and Big Data
## A Complex Ecosystem!



BIG DATA LANDSCAPE, VERSION 3.0

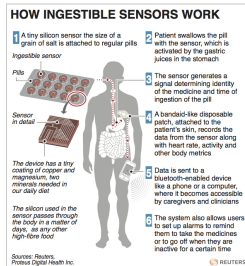© Matt Turck (@mattturck), Sutian Dong (@sutiandong) & FirstMark Capital (@firstmarkcap)

- Applied math **AND** Computer science
- Huge importance of domain specific knowledge: physics, signal processing, biology, health, marketing...

### Some joint math/computer science challenges

- Data acquisition
- Unstructured data and their representation
- Huge dataset and computation
- High dimensional data and model selection
- Learning with less supervision
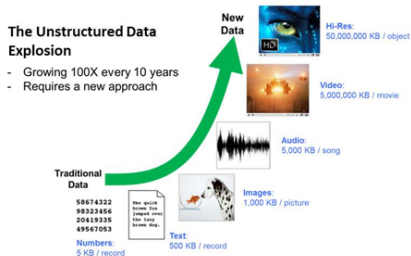- Visualization
- Software(s)...

## Some challenges

- How to measure new things?
- How to choose what to measure?
- How to deal with distributed sensors?
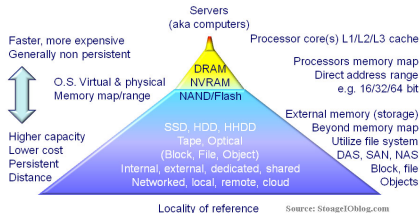- How to look for new sources of informations?
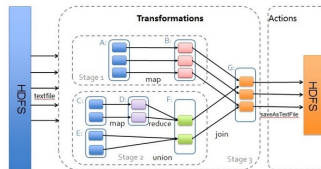
### Some challenges

- How to store efficiently the data?
- How to describe (model) them to be able to process them?
- How to combine data of different nature?
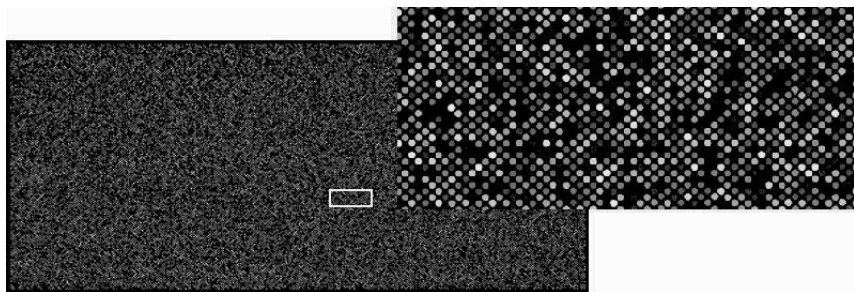- How to learn dynamics?

## Some challenges

- How to take into account the locality of the data?
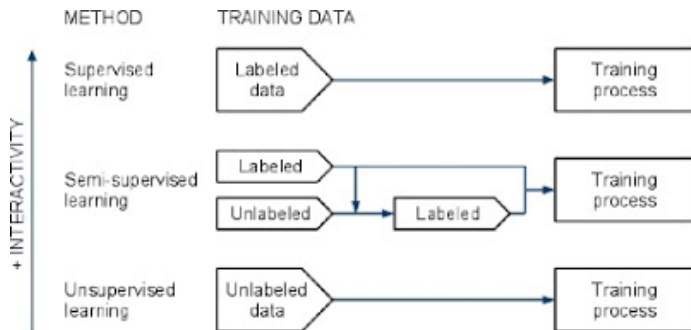- How to construct distributed architectures?
- How to design adapted algorithms?

### Some challenges

- How to describe (model) the data?
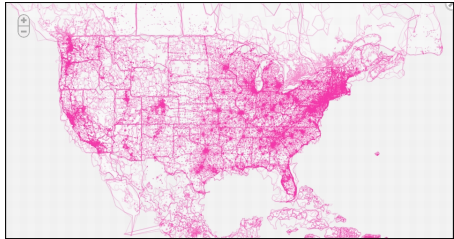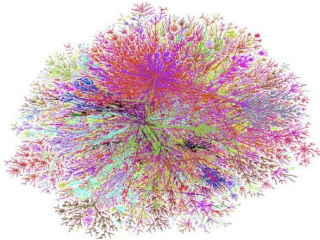- How to reduce the data dimensionality?
- How to select/mix models?

### Some challenges

- How to learn with the less possible interactions?
- How to learn simultaneously several related tasks?

## Some challenges

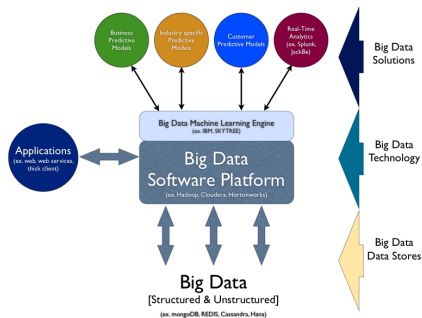- How to look at the data?
- How to present results?
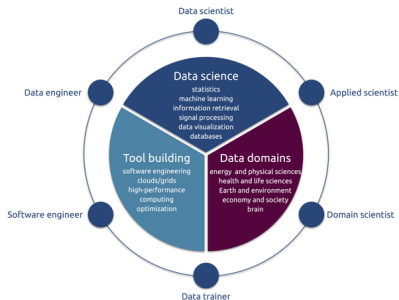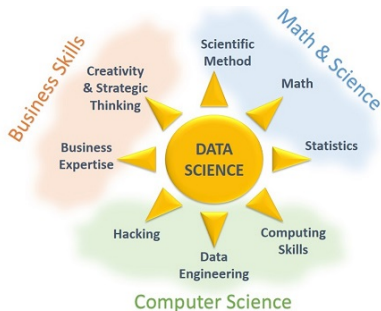- How to help taking better informed decision?

## Some challenges

- How to construct a consistent ecosystem?
- How to construct interoperable systems?

### Challenges

- No one masters all the skills!
- Importance of teams.
- Training...

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

G. James, D. Witten, T. Hastie and R. Tibshirani (2013)
An Introduction to Statistical Learning with Applications in R
*Springer Series in Statistics.*

B. Schölkopf, A. Smola (2002)
Learning with kernels.
*The MIT Press*

R. Schutt, and C. O'Neil (2014)
Doing Data Science: Straight talk from the frontline
*O'Reilly*