Data Science or Big Data

Erwan Le Pennec Associate Professor Applied Math Dpt



X Forum - 19/10/2016

Data is the new oil?

Big Data Everywhere!





The 5 Vs of Big Data

Big Data Everywhere!





Lots of Words!

Big Data Everywhere!





Doing Data Science

Big Data Everywhere!





Figure 2-2. The data science process

- Doing Data Science: Straight talk from the frontline.
 - Rachel Schutt, Cathy O'Neil
 - O'Reilly

Data Science



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

Data Science

Data Science



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...

Big Data is (quite) Easy

Data Science



Example of off the shelves solution





<pre>ief run(params: Params) { val.conf = new SparkKonf() i .setAppName(s"BinaryClassification with \$params") val.sc = new SparkContext(conf)</pre>	
Logger.getRootLogger.setLevel(Level.WARN)	
<pre>val examples = MLUtils.loadLibSVMFile(sc, params.input).cache()</pre>	
<pre>vol splits = examples_randomSplit(Array(0.0, 0.2)) vol training = splits(0).cache() vol test = splits(1).cache() vol set = splits(1).cach</pre>	
val updater = params.regType match { case L1 ⇒ new L1Updater() case L2 ⇒ new SquaredL2Updater() }	
<pre>via laporthm = nmx isstitl@pressionHibSO() algorthm.cpitiar .ethwafterstions(parmas.nufterstions) .eth(spit(parmas.stp)) .eth(spit(parmas.stp)) .eth(spit(parma).stp)) .eth(spit(parma).stp)) .eth(spit(parma).stp)) .eth(spit(parma).stp)) .eth(spit(parma).stp) .eth(spit(parma).stp)) .eth(spit(parma).stp)) </pre>	
<pre>val prediction = model.predict(test.map(features)) val predictionAndLabel = prediction.zip(test.map(label))</pre>	
<pre>val metrics = new BinaryClassificationMetrics(predictionAndLabel) val myMetrics = new MyBinaryClassificationMetrics(predictionAndLabel)</pre>	
<pre>println(s"Empirical CrossEntropy = \${myMetrics.crossEntropy()}.") println(s"Test areaUnderPR = \${metrics.areaUnderPR()}.") println(s"Test areaUnderROC = \${metrics.areaUnderROC()}.")</pre>	
sc.stop()	

Big Data is (quite) Easy

Data Science





```
export AWS_ACCESS_KEY_ID=<your-access-keyid>
export AWS_SECRET_ACCESS_KEY=
export AWS_SECRET_ACCESS_KEY=
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!

Data Science is (quite) Complex!

Data Science





Data Science is (quite) Complex!

Data Science





New Interdisciplinary Challenges

Data Science Challenges



- 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)
- Domain specific issue!

Data Acquisition





- 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?

Unstructured Data

Data Science 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?

Visualization

Data Science Challenges





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

Software(s)





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

Huge Dataset (Big Data)

Data Science Challenges





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

Domain Specific Knowledge

Data Science Challenges





- How to optimize a given process?
- How to find new processes / new outcomes?
- How to completly change a chain of value?
- How to measure the performance?

Recommendation System

Data Products



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 \$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 ...

Smart city

Data Products







of water allocated for domestic human use goes to urban cities. \$14 Billion in potable water is lost every year because of leaks, theft and unbilled usage. Source: World Bank 37,000 cloud experts support IBM's industry team alone. has been invested by IBM in more than a dozen acquisitions to accelerate its cloud initiatives.

IBM Intelligent Operations software is designed with cities, for cities, to provide the tools to monitor, visualize and analyze vital city services such as water and wastewater systems, transportation, infrastructure planning, permit management and emergency response.



Health

Data Products





Physics

Data Products





Big Data?





Hardware Constraints

- All the computations are done in a core using data stored somewhere nearby.
- Constrainst:
 - Data access / storage (Locality of Reference).
 - Multiple core architecture (Parallelization).
 - Cluster (Distribution)

Locality of Reference

Big Data





Memory Issue

- Data should be as **close** as possible from the core.
- Ideal case: dataset in the memory of a single computer.
- Useless if data used only once... (bottleneck = disk)
- **Split and Apply:** split the data in piece and work independently on each piece.
- Memory growth faster than data growth (Death of big data?)
- Memory required may be larger than dataset (interactions...)







Speed Issue

- Parallelization: Modern computer have several cores.
- HPC / DS setting: CPU bound tasks / IO bound tasks.
- Data science: Often embarrassingly parallel setting (no interaction between tasks).
- Not always acceleration due to IO limitation!

Distribution

Big Data





True Big Data Setting

- Computation in a **cluster**:
 - Distribution of the data(DS),
 - or/and distribution of the computation (HPC)
- Hadoop/Spark realm.
- Locally **parallel in memory** computation are faster... if data used more than once.
- Real challenge when not embarrassingly parallel (interaction...)

Data Scientists!

Data Scientist





Data Scientist

- Mix of Math/Statistics/ML, Computer Science, Business (and Communication) skills.
- Hard to be an expert of everything!
- Importance of balanced teams.

More than one Type of Data Scientists?

Data Scientist





Data Science Teaching at X



• Focus on the Data Scientist profile...

3 years program for X student

- 2A: Fundamentals (Math, Stat., Learn., CS) + 1 year project
- 3A: Dedicated track (Appl. Math/CS) + Projects + Internship
- 4A: Data Sciences Master

Data Sciences track of the **Master** Mathematics and Application of **Paris Saclay**

- Operated by Polytechnique in collaboration with Telecom, ENSAE, Paris Sud and ENS Cachan.
- Data Scientist training: statistical learning, machine learning, optimization, Big Data technologies...

Data Science at X



X/HEC Msc Big Data for Business

• 2 year program with both the technical skills and the strategic mindset for a business career with a Big Data expertise.

Data Sciences Starter Program

- 18 days continuous training program.
- RNCP certification.

Data Science Initiative at X

- More than 12 permanent researchers, 20 Phd students/Post doc, 10 engineers!
- Teaching and Research (projects, collaborations, chairs...)
- Examples:
 - Collaboration with CNAM.
 - Data Scientist chair (X, Keyrus, Orange, Thales)
 - Data Science for Insurance Sector chaire (Axa)
 - . . .

Don't Believe the Hype?

Conclusion





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



Data Science is here to stay...

- How to enhance decisions based on data!
- Business/Communication as important as ML/Stat/CS.
- If you are interested, learn and try!:
 - Initial training: Masters, Lectures...,
 - Online training: Books, Mooc...,
 - Continuous training: DSSP, Telecom...
 - **Practice:** R/Python, Hadoop/Spark, Kaggle/Challenge Data, Hackathon/RAMP...
- Lot of **opportunities** for a variety of **profiles**:
 - Data Scientist (strong in Stat/ML and Algorithm)
 - Domain Scientist (strong in Business/Sciences)
 - Data Engineer (strong in Development/Hardware)