

# Data Science Starter Program

## Introduction to Data Science

E. Le Pennec, A. Fermin



Spring 2015

# Introduction to Data Science

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

# Data Science in the media

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

Data Science in the media

Le Monde

[INTERNATIONAL](#) [POLITIQUE](#) [SOCIÉTÉ](#) [ÉCONOMIE](#) [CULTURE](#) [IDÉES](#) [SPORT](#) [SCIENCES](#) [TECHNO](#) [STYLE](#) [VOUS](#) [ÉDITION ABONNÉS](#)

# M Idées

**IDÉES** Les débats Think tanks Points de vue Editoriaux Opinions du Monde Analyses Idées chroniques Chats Blogs Forums

## Les données, puissance du futur

LE MONDE | 07.01.2013 à 15h10 • Mis à jour le 07.01.2013 à 18h03

Par Stéphane Grumbach, Stéphane Frénot

Abonnez-vous à partir de 1 € [Réagir](#) [Classer](#) [Imprimer](#) [Envoyer](#) Partager [Facebook](#) [Twitter](#) [LinkedIn](#)

[Recommander](#) [Envoyer](#) 467 personnes le recommandent.



**Les plus partagés**

- Une équipe de scientifiques filme un calamar géant par 900 mètres de fond dans le Pacifique 2212
- Infirmiers et aides-soignants refusent d'être des "pigeons" 1664
- Messi remporte son 4e Ballon d'or consécutif 987
- Mariage homosexuel : Wauquiez veut "forcer" le débat sur un référendum 629
- La première Eglise athée ouvre à Londres 603

**Nous suivre**

Retrouvez le meilleur de notre communauté

[Facebook](#) [Twitter](#) [R&D](#) [Mobile](#) [RSS](#)

# Data Science in the media

## NY Times

The New York Times Technology | Personal Tech | Business Day

[Log In](#) | [Register Now](#)



Search Bits

Go

OCTOBER 24, 2012, 9:00 AM | [4 Comments](#)

## Big Data in More Hands

By QUENTIN HARDY

[FACEBOOK](#)

[TWITTER](#)

[GOOGLE+](#)

[SAVE](#)

[E-MAIL](#)

[SHARE](#)

[PRINT](#)

Business people, Big Data is coming for you.

Software that captures lots of data and uses it to make predictions has mostly been the province of engineers skilled in arcane databases and statisticians capable of developing complex algorithms. As the business gets bigger, however, software makers are domesticating their products in the hope they will prove attractive to a broader population.

[Cloudera](#), which offers a popular version of the open source database called Hadoop, released software on Wednesday that makes it possible to run queries from a more mainstream SQL programming language interface. SQL, thanks to its adoption by Oracle, Microsoft and others, is known to millions of business analysts.

"This enables us to talk to a whole other class of customer," said Mike Olson, the chief executive of Cloudera. "The knock against Hadoop was that it is too complex."

There is a reason for that. Hadoop is one of several so-called unstructured databases that were created at Yahoo and Google, after those two companies found they had previously unimaginable amounts of data about activities like people's Web-surfing habits. Put into databases designed to handle this unstructured behavior, then analyzed, this information was

PREVIOUS POST

[◀ Google Shifts Pitch for Its New Chromebooks](#)

NEXT POST

[▶ In Contest for Rescue Robots, Darpa Offers \\$2 Million Prize](#)

### AROUND THE WEB »

THE NEXT WEB

[Google says Maps redirect on Windows Phone was a product decision, and will be removed](#)



BLOOMBERG  
[HTC Posts Lowest Net Income in Eight Years After Revenue Drops](#)



**SCUTTLEBOT** *News from the Web, annotated by our staff*

**Google's Schmidt arrive in North Korea**

REUTERS | From Mountain View to...errr, Pyongyang? - *Somini Sengupta*

**AP provides sponsored tweets during electronics show**

AP.ORG | The Associated Press is renting out its Twitter feed, with 1.5 million followers, to advertisers during C.E.S. - *Joshua Brustein*

**A history of griefing**

EDGE-ONLINE.COM | Meet the cult of gamers who want to ruin your day - just for kicks. - *Jenna Wortham*

**A Million First Dates**

THE ATLANTIC | Is online romance threatening monogamy? - *Jenna Wortham*

[SEE MORE ▶](#)

# Data Science in the media

World Bank

The screenshot shows the World Bank's homepage with a specific news article highlighted.

**THE WORLD BANK**  
Working for a World Free of Poverty

English | Español | Français | 中文 | Русский | GO

Search

ABOUT DATA RESEARCH LEARNING NEWS PROJECTS & OPERATIONS PUBLICATIONS COUNTRIES TOPICS

## World Bank Live

### What Happens When Big Data Meets Official Statistics? - Live Webcast

What happens when official statistics meets...

**#bigstats**  
December 19th 2.30pm  
World Bank HQ  
MC13 -121  
[bigstats.eventbrite.com](http://bigstats.eventbrite.com)

SHARE

**ABOUT**  
World Bank Live is a space to discuss key development topics in real time. Chat live with experts, watch livestreams and participate in events, ask tough questions.

Subscribe to alerts on upcoming events

E-mail: \*

# Data Science in the media

Criteo



Solutions CPOP Plateforme Editeurs Success Stories Actualités Carrières Contact Login annonceur | FR ▾

Accueil > Actualités & Événements

## CriteoLabs : soirée d'inauguration

### Criteo inaugure à Paris l'un des premiers centres de R&D en publicité prédictive d'Europe

- › Fleur Pellerin, Ministre déléguée chargée des PME, de l'Innovation et de l'Economie Numérique, apporte son soutien à cette entreprise innovante du secteur numérique, véritable « success story » à la française.
- › Criteo inaugure CriteoLabs, son nouveau centre de R&D de 10.000 m<sup>2</sup> au cœur de Paris.
- › Avec à terme 300 ingénieurs, ce site est déjà l'un des premiers centres européens de R&D en algorithmes appliqués à la publicité en ligne. Pour accompagner sa forte croissance, Criteo recrute cette année 250 nouveaux collaborateurs.



Jean-Baptiste Rudelle, CEO et Pascal Gauthier COO



Arrivée de Fleur Pellerin



Visite de l'Observatoire de la Publicité à Paris

Criteo inaugure à Paris l'un des plus gros pôles européens de R&D dédiés à la publicité prédictive, CriteoLabs. Sur 10.000 m<sup>2</sup>, ce nouveau centre a vocation à accueillir 300 ingénieurs et à permettre ainsi à Criteo de garantir son avancée technologique sur ses 30 marchés d'exportation, des Etats-Unis, à l'Europe, en passant par l'Asie. Cette année, l'entreprise compte ainsi recruter 250 nouveaux collaborateurs, dont une centaine d'ingénieurs.

Ce nouveau siège, que Criteo a choisi délibérément de situer à Paris, vient ponctuer un développement continu, qui a permis à l'entreprise d'atteindre des résultats remarquables, 3 ans seulement après son lancement commercial :

- › 600 salariés présents dans 15 bureaux dans le monde
- › 2 000 annonceurs, parmi les plus importants e-commerçants mondiaux tels que Dell, Macy's, John Lewis, Marks & Spencers, Zalando, La Redoute, Les 3 Suisses, etc.
- › 4 000 éditeurs
- › Plus de 200 millions de dollars de CA en 2011

# Data Science in the media

We are in the press as well...

# L'USINE digitale

Quand le numérique réinvente l'industrie

TOUTE L'INFO | L'USINE NOUVELLE | INSCRIVEZ-VOUS À LA NEWSLETTER | DIGITAL AVENUE | Rechercher dans L'Usine Digitale | S'abonner | Suivre

INTERNET | LOGICIELS & APPLICATIONS | HARDWARE | CLOUD ET DATA | INDUSTRIES | ECONOMIE NUMÉRIQUE | ANNUAIRE DE START-UP

USINE DIGITALE > L'USINE CAMPUS

## Polytechnique forme les professionnels au big data

Par Cécile Maillard - Publié le 05 juin 2014, à 15h49

► [L'Usine Campus](#), [Numérique](#), [Formation](#), [France](#), [L'actu des campus](#),



Polytechnique forme les professionnels au big data © Wikimedia Commons

L'Ecole polytechnique ouvre à la rentrée une formation courte pour les doctorants et salariés qui ont besoin de comprendre quelles belles opportunités leur ouvrent les big data. Les écoles d'ingénieurs et de management sont de plus en plus nombreuses à proposer des formations d'analyse des données massives.

"Nous avons une pression des entreprises, qui nous questionnent de plus en plus sur ce que propose l'Ecole polytechnique sur les big data", explique Frank Pacard, directeur de l'enseignement et de la recherche à l'X. Pour répondre à leur demande, Polytechnique ouvre, en octobre 2014, un parcours baptisé "Data

Science Starter Program". "C'est un programme d'évangélisation, qui expliquera à des doctorants, post-doctorants ou salariés, tout ce qu'ils peuvent tirer de l'exploitation des données, dans leurs métiers et entreprises."

### A LIRE SUR LE MÊME SUJET

Les quatre conseils d'IBM et SAP pour éviter les pièges du big data

### RECEVOIR NOTRE NEWSLETTER :

E-Mail

OK

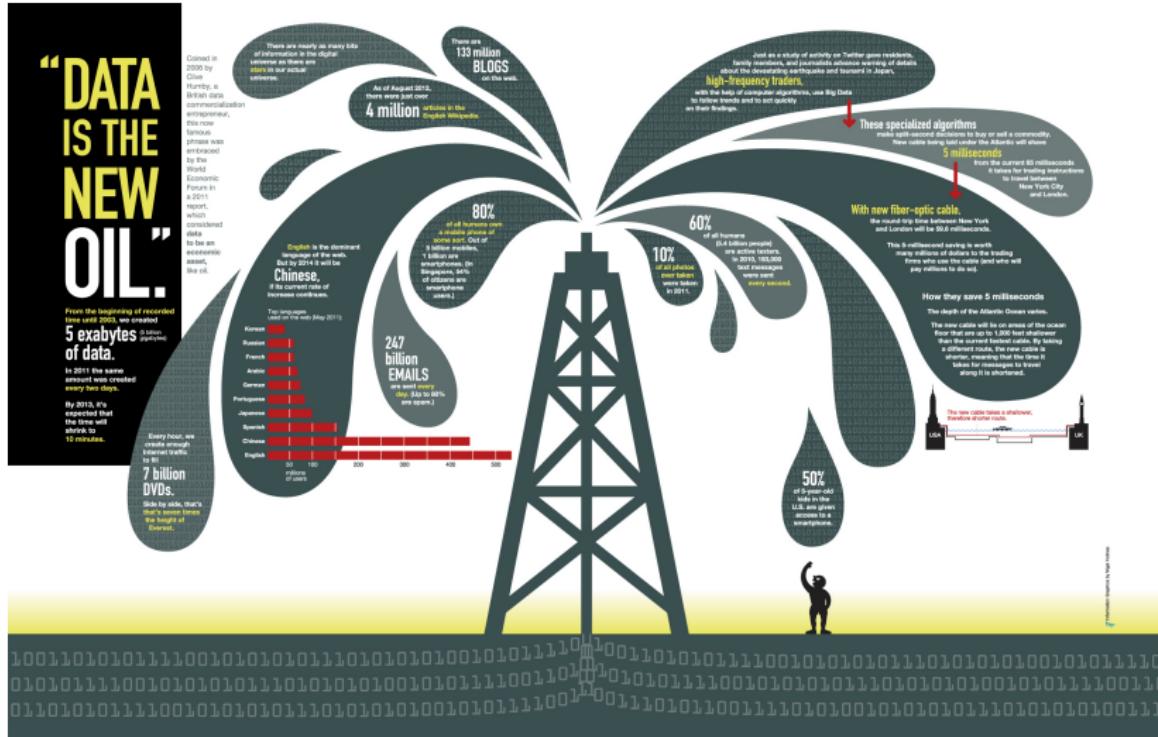
### LES PLUS LUS DE LA RUBRIQUE «NUMÉRIQUE»

- La SNCF roule au digital
- L'accéléromètre, ce mouchard caché dans nos smartphones
- Le marché européen pourrait faire gagner 18 millions d'abonnés à Netflix d'ici à 2018
- La dataviz, facteur de compétitivité pour les entreprises
- Nous allons révéler la face numérique de l'industrie

### LES AUTRES ACTUALITÉS DE LA RUBRIQUE «NUMÉRIQUE»

# Data Science in the media

Data is the new oil?



# From Data to Product

## Outline

- 1 Data Science in the media
- 2 From Data to Product**
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

# From Data to Product

## Web search

Google moteur de recherche

Web Actualités Images Vidéos Maps Plus Outils de recherche

Environ 10 100 000 résultats (0,24 secondes)

**Moteur de recherche - Mozbot France - La recherche facile ...**  
[www.mozbot.fr/](http://www.mozbot.fr/) ▾  
Moteur de recherche Mozbot en partenariat avec Brioude-Internet, Abondance et Google : résultats, synonymes, expressions connexes, statistiques mots clés, ...

**Actualités correspondant à moteur de recherche**

 Le moteur de recherche DuckDuckGo bloqué en Chine  
Le Monde - il y a 3 heures  
Selon le site spécialisé TechInAsia, le moteur de recherche serait bloqué depuis le 4 septembre dans le pays. DuckDuckGo, qui se présente ...

Canoë L'Allemagne souhaite que Google dévoile les algorithmes ...  
Clubic.com - il y a 5 jours

Plus d'actualités pour "moteur de recherche"

**Moteur de recherche — Wikipédia**  
[fr.wikipedia.org/wiki/Moteur\\_de\\_recherche](http://fr.wikipedia.org/wiki/Moteur_de_recherche) ▾  
Un moteur de recherche est une application web permettant de retrouver des ressources (pages web, articles de forums Usenet, images, vidéo, fichiers, etc.) ...

**Moteur de Recherche SEEK.fr™**  
[www.seek.fr/](http://www.seek.fr/) ▾  
Moteur de recherche alternatif français respectant la vie privée via un métamoteur utilisant les principaux moteurs de recherche ainsi qu'un annuaire ...  
Metamoteur Web SEEK.fr - A Propos de Seek - Horoscope - Seek annuaire

# From Data to Product

## Recommendation system

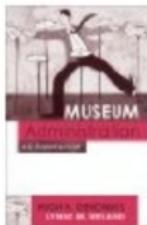
### More Ideas Based on Your Browsing History

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...**

# From Data to Product

## Advertisement

**Outlet**  
» Descubrelos

Libros universitarios y de estudios superiores a precios bajos  
» Descubrelos



Innovatoren und Kleinunternehmer nutzen ihre Möglichkeiten bei Amazon  
» Ihre Geschichten



Jetzt neu:  
Schnell & einfach Ersatzteile finden  
» Hier klicken



365 Tage im Jahr Licht bei 0€ Stromkosten  
» Hier klicken



**Neuheiten von Makita**  
» Hier klicken



fire PHONE + 12 MONTHS OF PRIME  
NOW ONLY \$0.99  
with a two-year contract [Shop now](#)



Fall Outlet Event  
» Shop now



FALL COATS  
» See more



New from iRobot:  
Roomba 870 Vacuum Cleaning Robot  
» Learn more



Save Big on Outdoor Fire Pits from Strathwood  
» Shop now



Rentrée des Conservatoires  
-10% sur une sélection d'instruments\*  
\*Voir conditions [Cliquez ici](#)



Vos courses en livraisons gratuites et régulières  
» Economisez en vous Abonnant [Cliquez ici](#)



PROMOTIONS CHAUSSURES -30% -40% -50%...  
» J'en profite



PROMOTIONS SACS À MAIN  
» J'en profite



# From Data to Product

## Intrusion detection



# From Data to Product

## Crime prevention



ABOUT HOW PREDPOL WORKS PROVEN RESULTS TECHNOLOGY PRESS CONTACT US BLOG

# PREDICTIVE POLICING®

The Predictive Policing Company.

PredPol's cloud-based software enables law enforcement agencies to better prevent crime in their communities by generating predictions on the places and times that future crimes are most likely to occur.

# From Data to Product Marketing



chiefmartec.com Marketing Technology Landscape

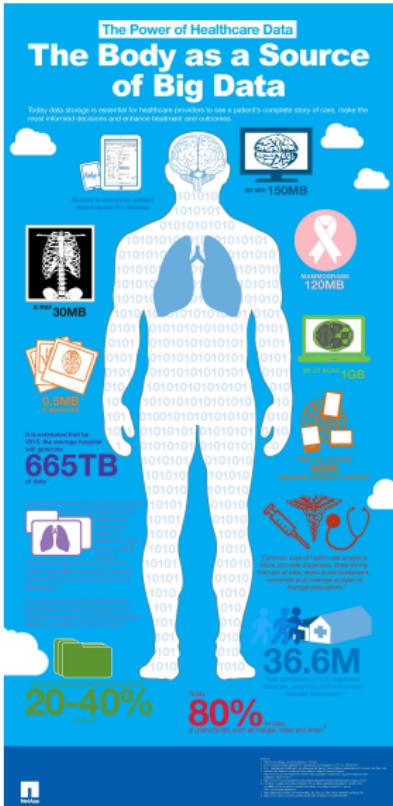
January 2014

The image is a dense grid of company logos, each representing a different brand or service within the marketing and technology sectors. The grid is organized into several sections, each with a distinct color and title. The sections include:

- MARKETING EXPERIENCES**: Email Marketing, Creative & Design, Call Centers, Marketing Apps.
- MARKETING OPERATIONS**: Marketing Data, Marketing Analytics.
- MARKETING ANALYTICS**: Marketing Resource Mgmt, Channel/Local Mktg, Dashboards, Web & Mobile Analytics.
- MARKETING AUTOMATION**: Testing & Optimization, Content Marketing, Personalization, Loyalty & Gamification, Marketing Automation / Integrated Marketing.
- MARKETING TECHNOLOGY**: Display Advertising, Video Ads & Marketing, Events & Webinars, Testing & Optimization, Sales Enablement, Agile & Project Mgmt.
- DATA MANAGEMENT**: Data Management Platforms/Customer Data Platforms, Tag Management, Lead Management, Cloud Connectors, APIs.
- MIDDLEWARE**: Data Integration, Business Process Management, Application Integration.
- BACKBONE PLATFROMS**: CRM, Marketing Automation / Integrated Marketing, Web Site / WCM / WEM.
- INFRASTRUCTURE**: Databases, Big Data, Cloud, Mobile App Dev, Web Dev, Internet.
- ENVIRONMENT**: Marketing Environment.

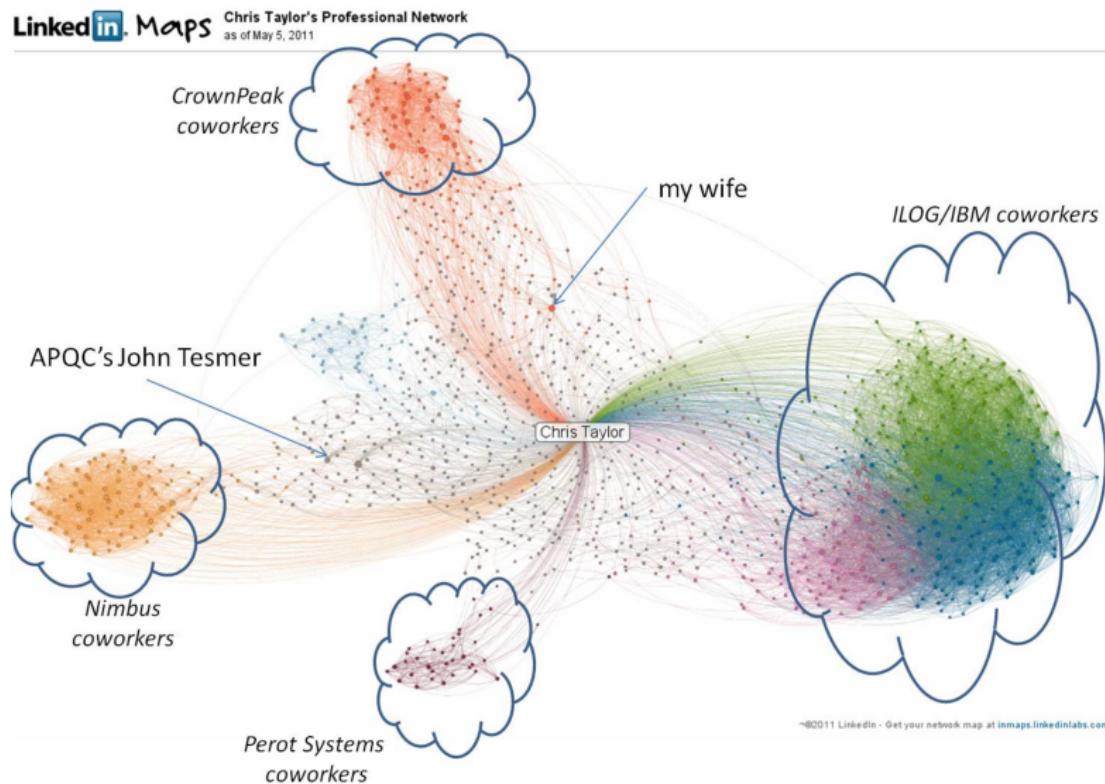
The logos are arranged in a grid format, with each logo being a small square or rectangle containing the company's name and a unique icon. The colors of the logos vary widely, creating a vibrant and diverse visual representation of the industry.

## From Data to Product Health



# From Data to Product

## LinkedIn



# From Data to Product

## Smart city

### Smarter Cities: Turning Big Data Into Insight



#### City Planning and Operations

**\$1 Trillion**

global annual savings could be attained by optimizing public infrastructure.  
Source: McKinsey

#### Transportation Analytics

**50 Hours**

of traffic delays per year are incurred, on average, by travelers.

**\$57 Trillion**

in infrastructure investments will be needed between 2013-2030.  
Source: McKinsey

**30 Billion**

people all over the world travel approximately 30 billion miles per year. By 2050, that figure will grow to over 150 billion miles.

#### Water Management

**60%**

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

**\$6 Billion**

has been invested by IBM in more than a dozen acquisitions to accelerate its cloud initiatives.



Cloud is driving cities in their digital transformation.

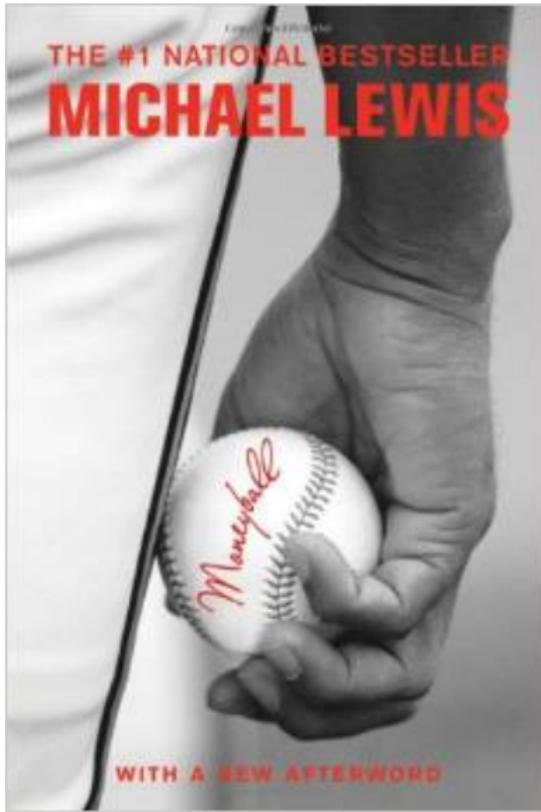
#### Open Cloud

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.



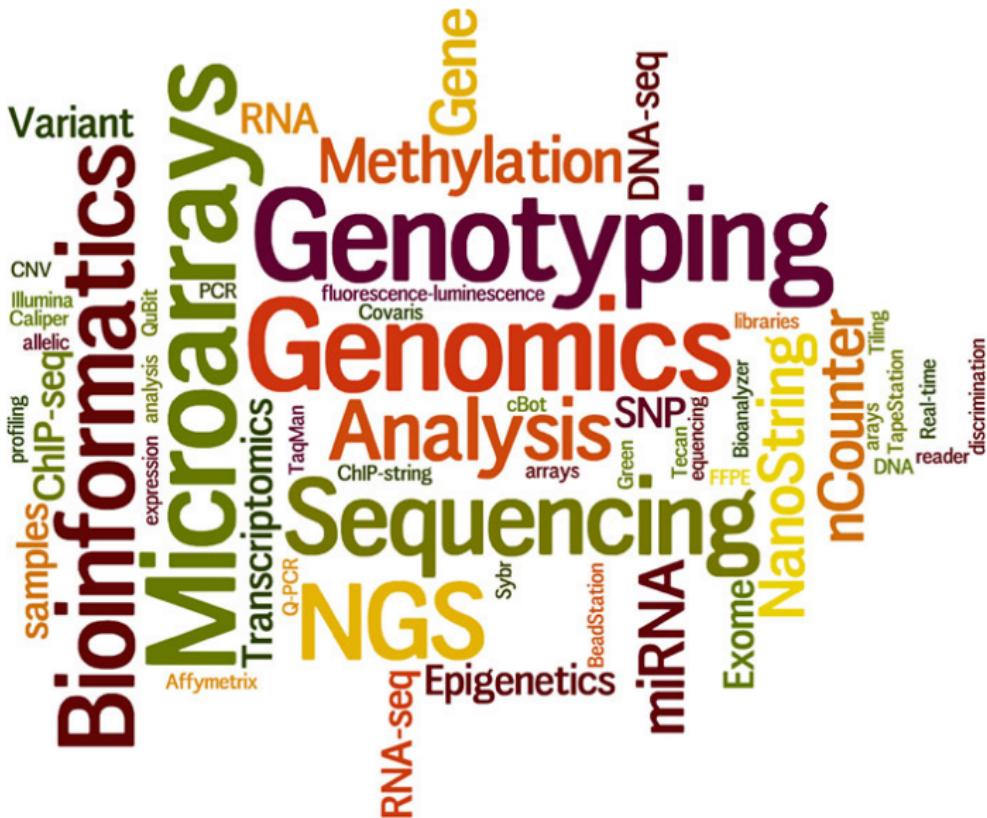
# From Data to Product

## Sports

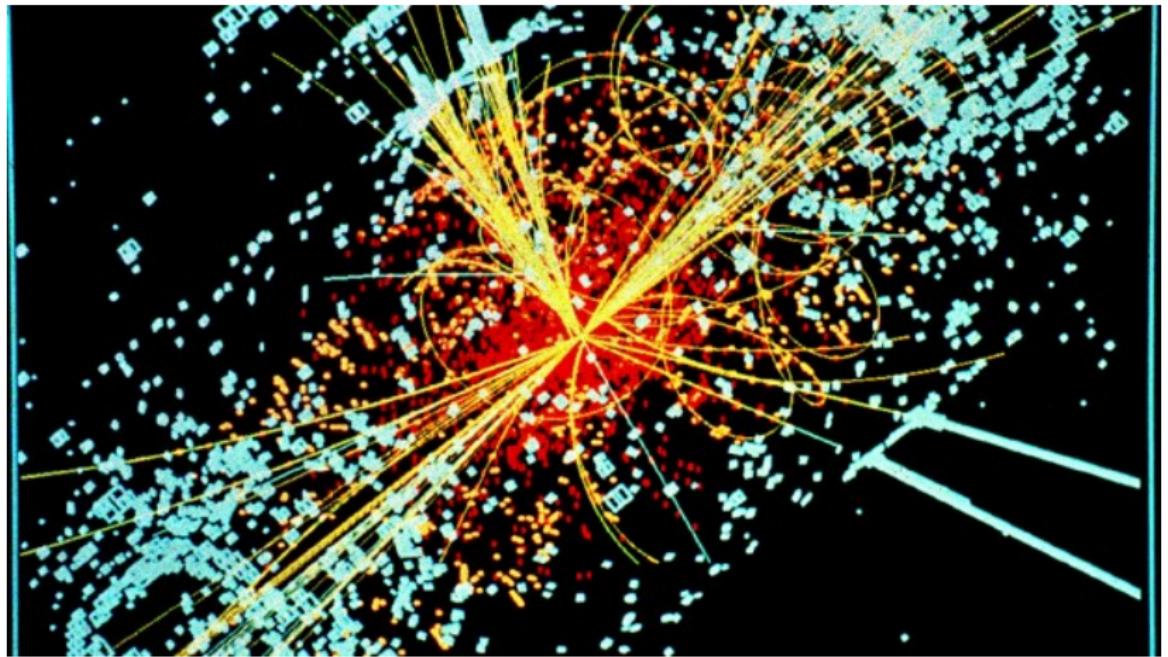


# From Data to Product

## Genomics



# From Data to Product Physics



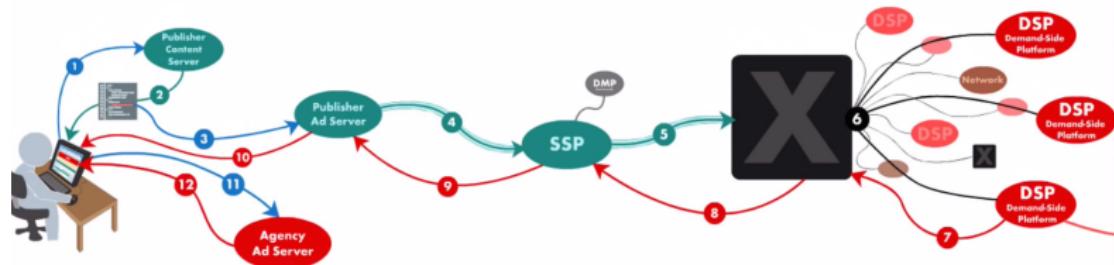
# An example: Real Time Bidding

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding**
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

## An example: Real Time Bidding

### An example: Real Time Bidding



- A **customer** visits a webpage with his browser: a complex process of content selection and delivery begins.

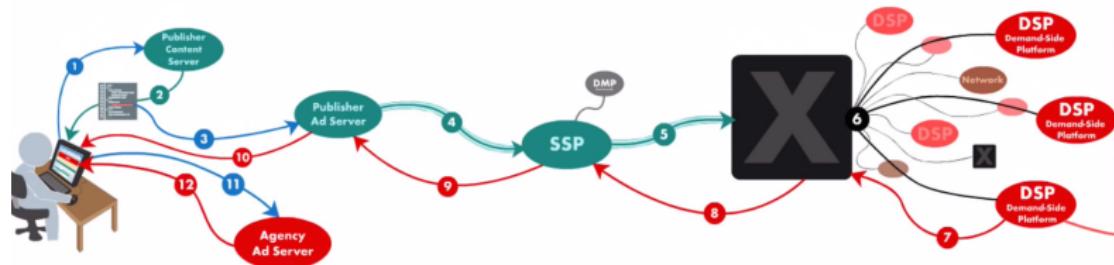
An **advertiser** might want to display an ad to this customer on the webpage he is going to

The webpage belongs to a **publisher**. The publisher delivers contents: news, music, information, sports, etc. This content draws an **audience**

The publisher sells ad space to advertisers who want to reach that audience

## An example: Real Time Bidding

### An example: Real Time Bidding



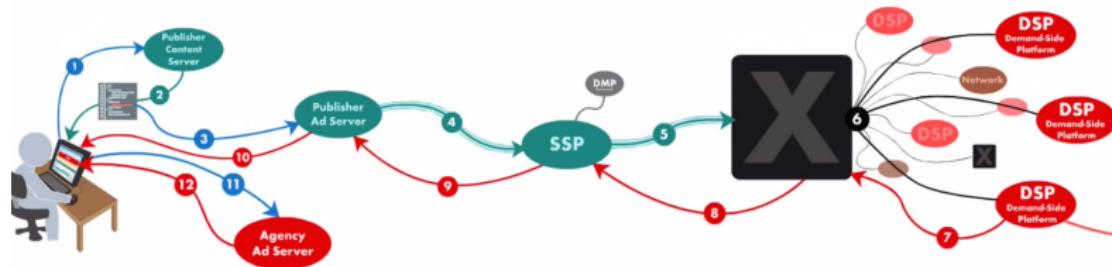
- 1. 2. The customer visits a publisher's webpage: the browser opens a connection to the **publisher's content server**. It returns the content for the page (html code).

The html code describing this content is retrieved by the browser, and it starts to render and interpret it.

But... there is a line in this html code that says “follow this URL to retrieve ad content”

## An example: Real Time Bidding

### An example: Real Time Bidding

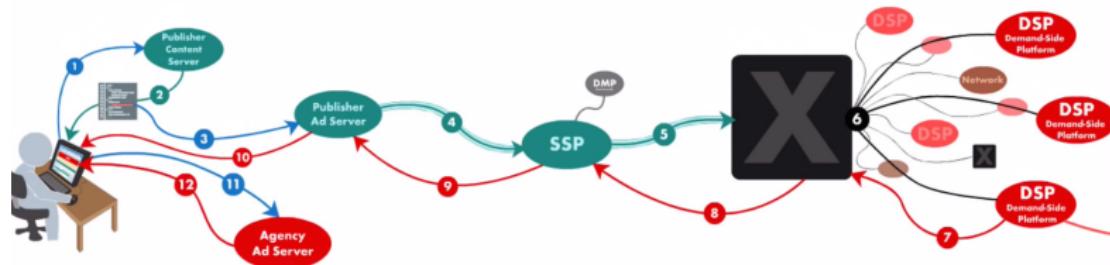


- 3. The publisher has an **ad server**: it answers the request by considering possibilities: can I put an ad for my premium buyers ? Do I have data about this consumer viewing my content? (it could help me to decide to which buyer I could give this display opportunity). Only logical rules apply (no machine-learning here).

The ad display opportunity is not premium, and this space or type of customer is not already reserved by a buyer. **Publisher's ad server** puts this opportunity of ad display in the open ad market.

## An example: Real Time Bidding

### An example: Real Time Bidding



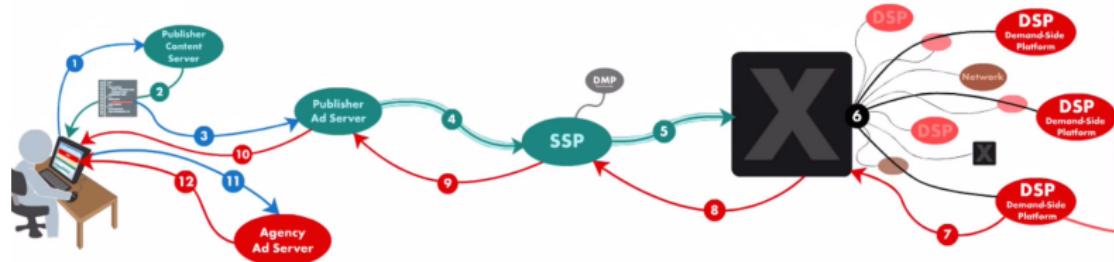
- 4. The **publisher ad server** connects to an **SSP** (Supply-Side Platform). This platform monetizes its programmable display inventory.

The SSP asks: have I already seen this consumer before ? Do I have additional data on him? The SSP requests extra information to a **DMP** (Data-Management Platform) about the user: profiling, audience segments, etc. Here machine learning is applied.

- 5. Using this information, the SSP sends the ad request to an **ad-exchange**

## An example: Real Time Bidding

### An example: Real Time Bidding



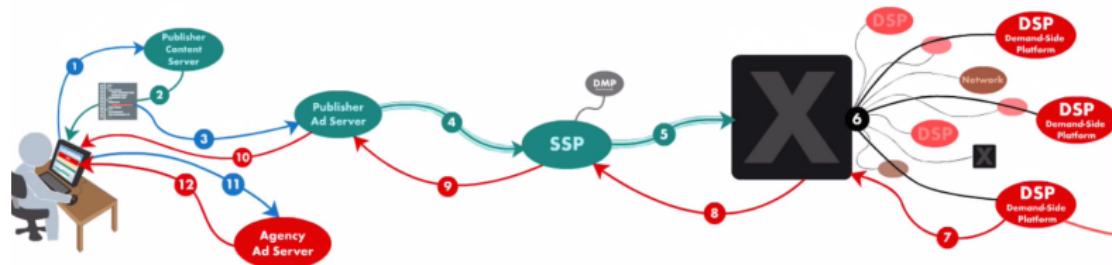
- 6. Meanwhile, the ad-exchange is connected and exchanges with many potential buying systems: DSP (Demand-Side Platform), ad-networks, even other ad-exchange networks.

Ad-network and DSP can have **pre-cached bid**: I'm paying 1\$ for 1000 displays of 25years-old males in France, I buy 100 displays as soon as the price is below some threshold (like a broker).

If no pre-cached bids, the ad-exchange says: no direct buyer for this display. Let's us an auction rule! The **RTB** (Real Time Bidding) begins.

## An example: Real Time Bidding

### An example: Real Time Bidding

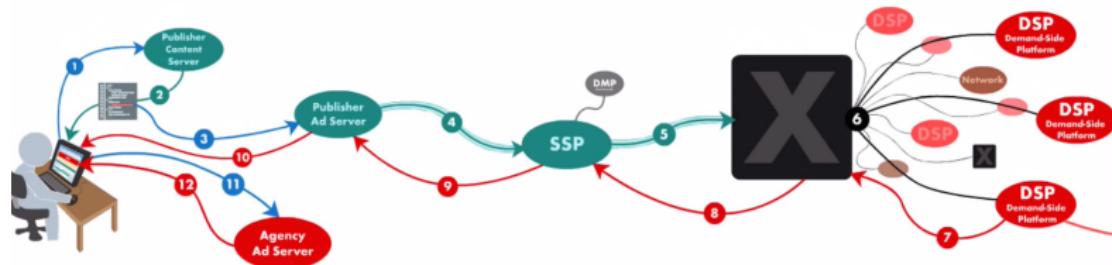


- 6. RTB: buyers have **10ms** (!) to give a price to the ad-exchange. Buyers assess in real-time how willing they are to display an ad to this customer.

Machine learning is used here, but only the **prediction** step, e.g. to assess the probability that the customer will click on some ads. The model must contain few parameters to answer quickly: the use of feature selection in the training step is crucial here.

## An example: Real Time Bidding

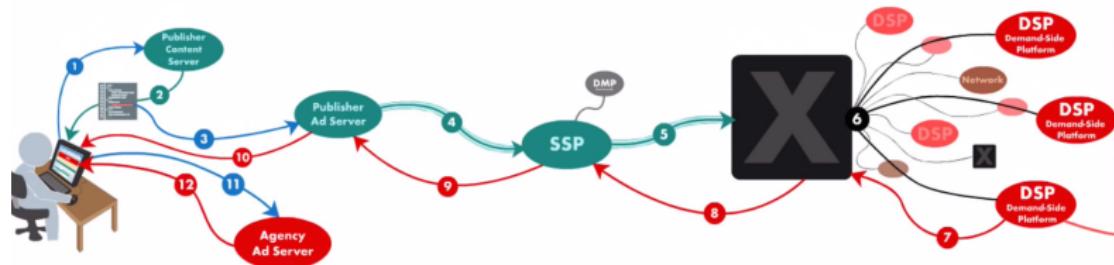
### An example: Real Time Bidding



- 7. The ad-exchange selects the highest bidder. The winning DSP gives instruction to the ad-exchange to retrieve the ad creative.
- 8. The ad-exchange passes these instructions to the SSP
- 9. The SSP send the request to the publisher ad server
- 10. The publisher ad server responds to the still existing http connection of the browser,
- 11. 12. and tells to the browser to go to the **agency's ad server** to download the ad.

## An example: Real Time Bidding

### An example: Real Time Bidding



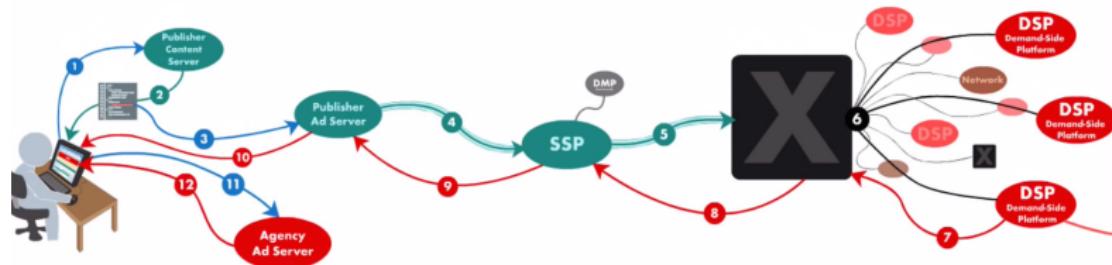
- Now the ad can be displayed in the browser.

Full process takes < 100ms !

- Where is data science:
  - DMP side, to cluster audience into marketing segments and to profile customers: clustering and classification
  - Buyer's side (DSP, ad-network) to compute the price proposed for RTB. Need to estimate the probability of a click on ads: regression and classification

## An example: Real Time Bidding

### An example: Real Time Bidding



- Some numbers for a large web-advertisement company:
  - 10 million prediction of click probability per second
  - answers within 10ms
  - stores 20Terabytes of data daily

# Data Science ecosystem

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

# Data Science ecosystem

## A new Context

### Data everywhere

- Huge volume,
- Huge variety...

### Affordable computation units

- Cloud computing
- Graphical Processor Units (GPU)...
- Growing academic and industrial interest!

# Data Science ecosystem

## Big Data is (quite) Easy

### Example of *off the shelves* solution



```
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()
```

# Data Science ecosystem

## Big Data is (quite) Easy

### 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
```

⇒ Logistic regression for arbitrary large dataset!

# Data Science ecosystem

## Big Data is (quite) Easy

### Example of *off the shelves* solution



```
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()
```

# Data Science ecosystem

## Big Data is (quite) Easy

### 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

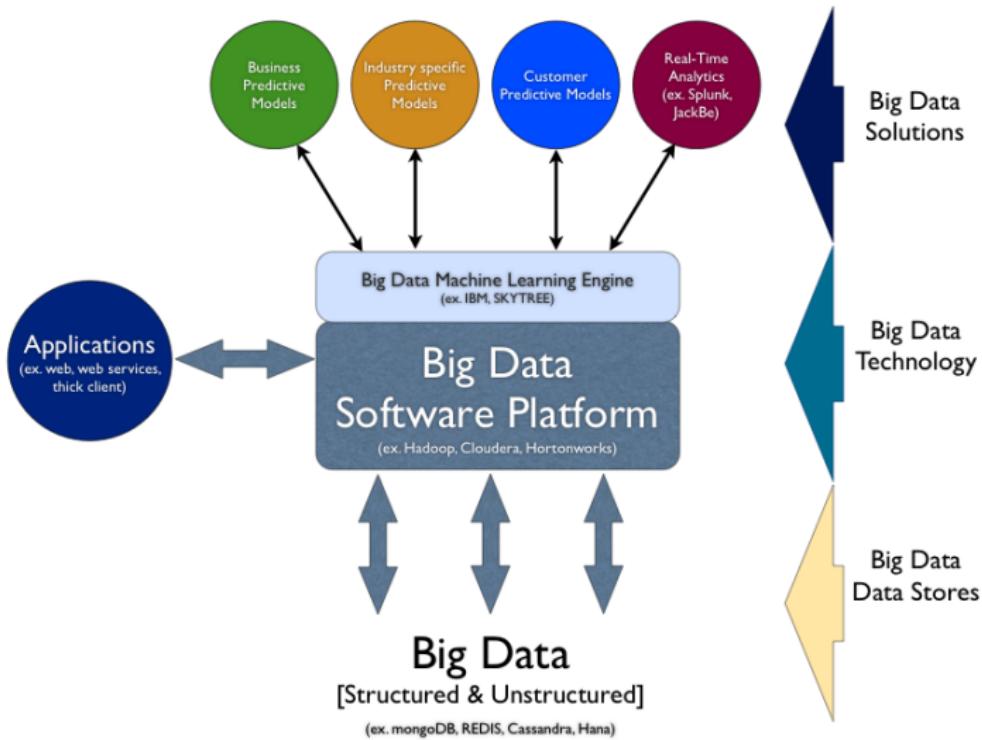
cellule/spark/bin/spark-submit \
--class fr.cc.challenge.Preprocess \
challenges_2.10-0.0.jar \
/data/train.csv \
/data/train2.csv

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

⇒ Logistic regression for arbitrary large dataset!

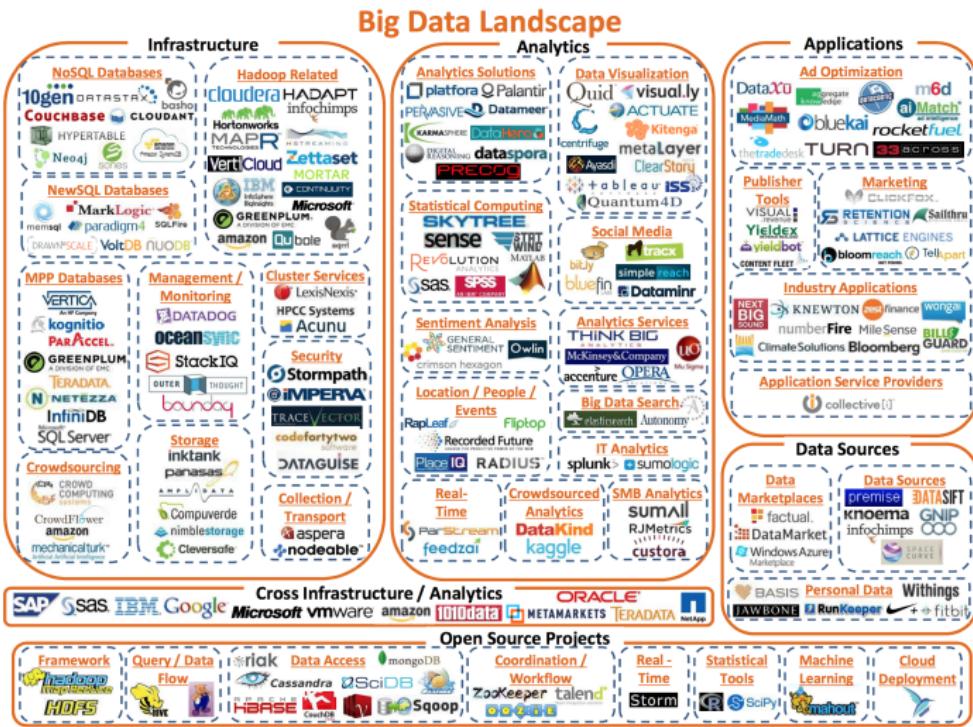
# Data Science ecosystem

A Complex Ecosystem!



## Data Science ecosystem

## A Complex Ecosystem!



Matt Turck (@mattturck) and Shivon Zilis (@shivonz)

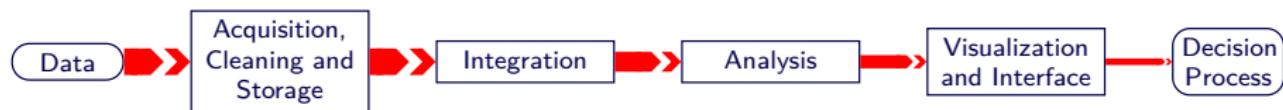
# Data cycle

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

## Data cycle

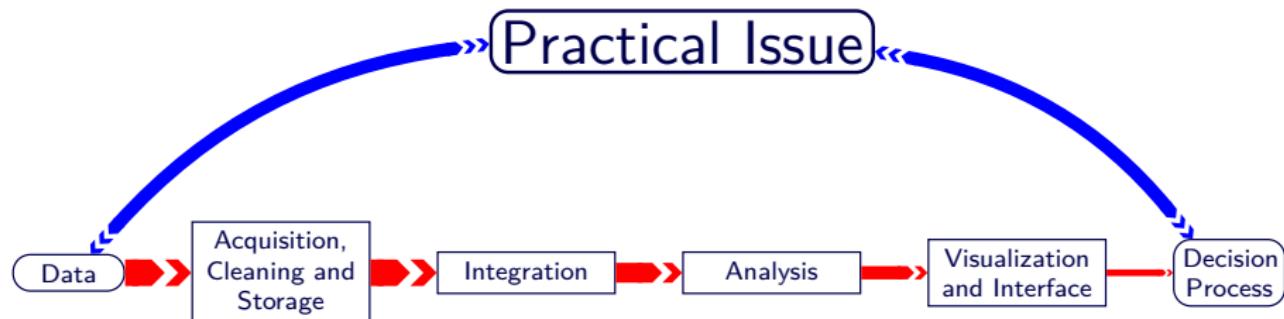
## Data Cycle



- Data/Information flow vision

Data cycle

Data Cycle



- Data/Information flow vision
- Goal oriented

- Raw material:

- Structured and unstructured data (**Variety**)
- Data quality issue (**Veracity**)
- Quantity (**Volume** and **Velocity**)

- Various sources:

- Open data,
- Proprietary data

## Data cycle

### Acquisition, Cleaning, Storage and Integration

- Get the data from the sources.
- Storage issue and availability for processing.
- Cleaning and formating
- Integration: Data preparation for analysis
- Time consuming!

- Extract information from the data
- Statistics/Machine learning
- Big Data: hardware is the limit (time/volume)

- Reporting part: **V**isualization, text...
- Also used for data exploration
- Very important aspect!

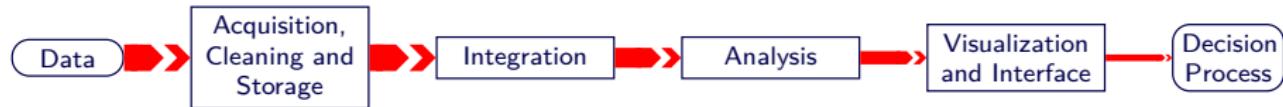
## Data cycle

### Decision and goal oriented analysis

- Better decisions: **Value**
- Need to answer a problem/question!
- Need to formalize the problem: no answer without a question!
- Feedback everywhere...

# Data cycle

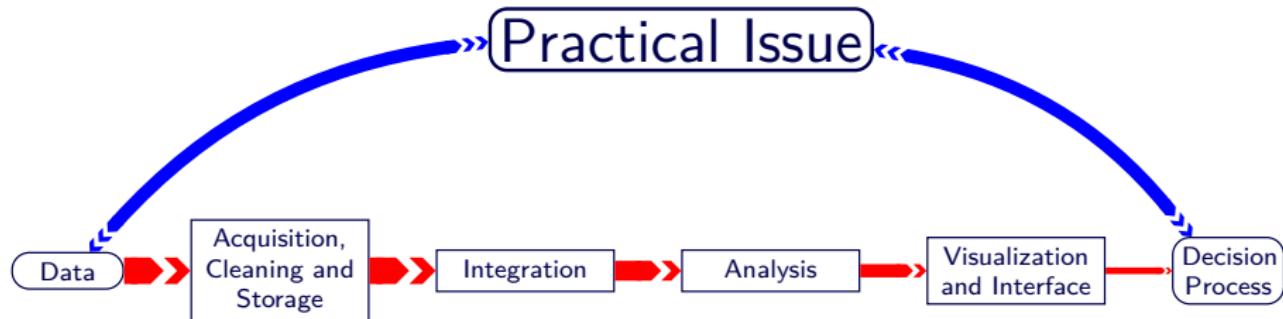
## Real Data Cycle



- Data/Information flow vision

## Data cycle

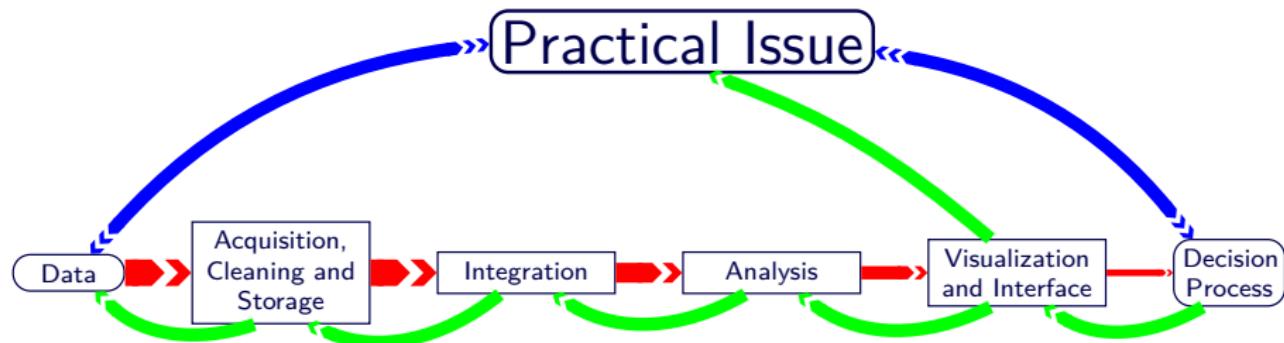
### Real Data Cycle



- Data/Information flow vision
- Goal oriented

## Data cycle

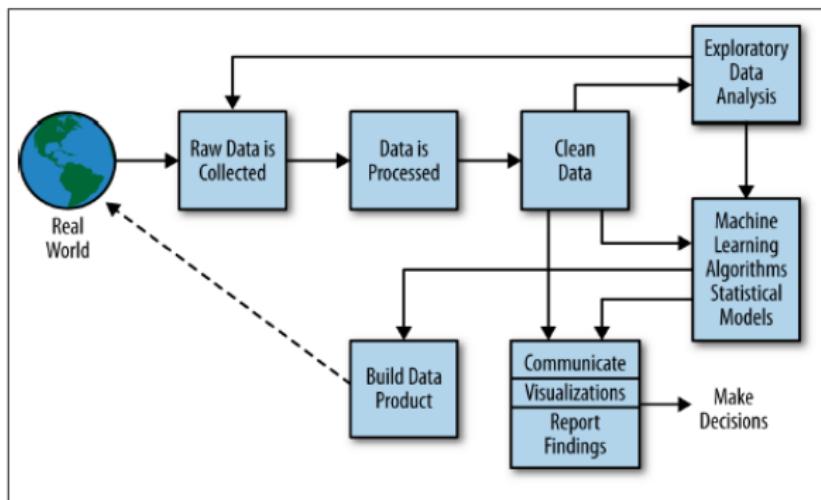
### Real Data Cycle



- Data/Information flow vision
- Goal oriented
- Iterative and interactive process

# Data cycle

## Doing Data Science



*Figure 2-2. The data science process*

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

# Data Science project

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project**
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

# Data Science project

A 7 step program



## 1. Identify the problem

- Type of problems and metric used to measure success
- Identify key people within your organization and outside
- Get specifications, requirements, priorities, budgets
- How accurate the solution needs to be?
- Do we need all the data?
- Outsourcing?

# Data Science project

A 7 step program



## 2. Identify available data sources

- Extract and check sample data / Perform Exploratory Data Analysis
- Assess quality of data, and value available in data
- Identify data glitches, find work-around
- Data quality improvement?
- Verify with field expert that you understand the data
- Infrastructure?

# Data Science project

A 7 step program



## 2. Identify available data sources

- Extract and check sample data / Perform Exploratory Data Analysis
- Assess quality of data, and value available in data
- Identify data glitches, find work-around
- Data quality improvement?
- Verify with field expert that you understand the data
- Infrastructure?

## 3. Identify if additional data sources are needed

- What? How much? How to?
- Real time?
- Do we need experimental design?

# Data Science project

A 7 step program



## 4. Data preparation and analyses

- Data preparation and cleaning
- Explore methodologies
- Select variables and models
- Detect / remove outliers
- Validate chosen methodology
- Measure accuracy, provide confidence intervals
- Provide visualization and ask for feedback

# Data Science project

A 7 step program



## 4. Data preparation and analyses

- Data preparation and cleaning
- Explore methodologies
- Select variables and models
- Detect / remove outliers
- Validate chosen methodology
- Measure accuracy, provide confidence intervals
- Provide visualization and ask for feedback

## 5. Implementation, development

- FSSRR: Fast, simple, scalable, robust, re-usable
- Debugging
- Need to create an API to communicate with other apps?

# Data Science project

A 7 step program



## 6. Communicate results

- Integration and visualization
- Discuss potential improvements (with cost estimates)
- Provide training
- Code and methodology documentation

# Data Science project

A 7 step program



## 6. Communicate results

- Integration and visualization
- Discuss potential improvements (with cost estimates)
- Provide training
- Code and methodology documentation

## 7. Maintenance

- Test the model or implementation; stress tests
- Regular updates
- Outsourcing?

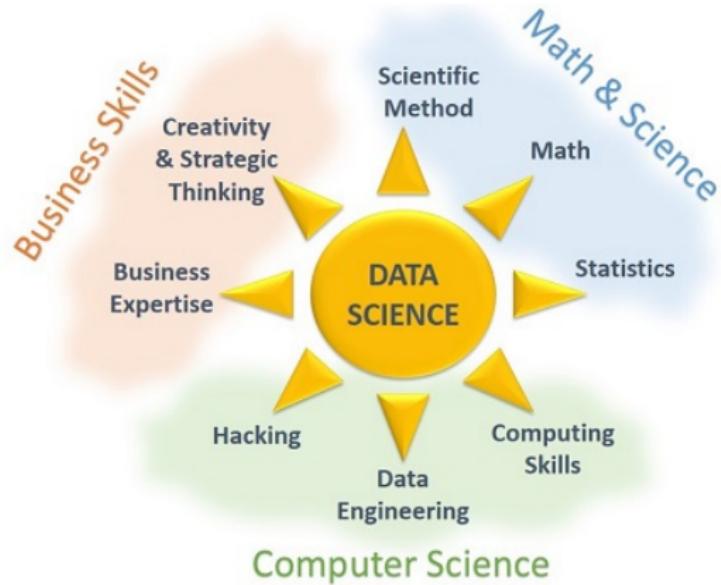
# Data scientists

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

- Business:
  - Business analysis, market knowledge, product usage,  
...
- Data Management:
  - Data collection, storage, cleaning, filtering,  
integration,...
- Statistic and Machine Learning:
  - Data modeling, inference, prediction, pattern  
recognition,...
- Programming:
  - Software development, Large-scale or parallel data  
processing,...
- Interface and Data Visualization:
  - HCI design, visualization, story-telling,...

# Data scientists Profiles

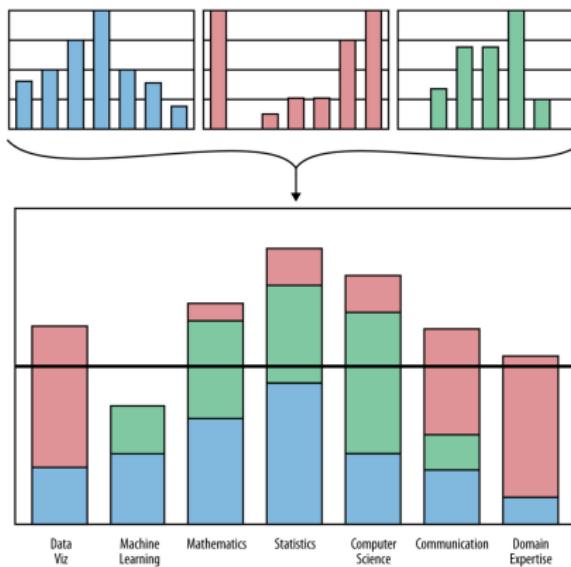


- No one masters all the skills!

# Data scientists

## Data science team

No one person can be the perfect data scientist, so we need teams.



- Gather people having different skills

# Data scientists

## Main types of data scientists

There are the ones...

- Strong in **statistics**: develop new statistical theories for big data: statistical modeling, experimental design, sampling, clustering, data reduction, confidence intervals, testing, modeling, predictive modeling, etc.
- Strong in **mathematics**: operations research, analytic business (inventory management and forecasting, pricing optimization, supply chain, quality control, yield optimization)
- Strong in **data engineering**, Hadoop, database/memory/file systems optimization and architecture, API's, Analytics as a Service, optimization of data flows

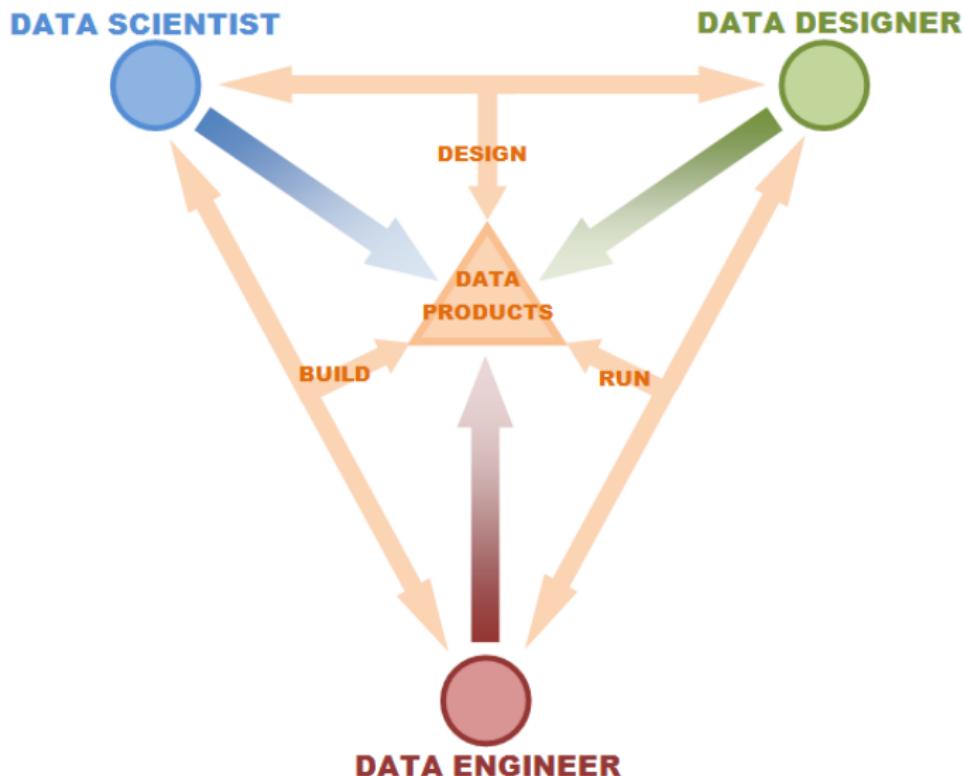
# Data scientists

## Main types of data scientists

- Strong in **computer science** (algorithms, computational complexity, optimization)
- Strong in **business**, ROI optimization, decision sciences (dashboards design, metric mix selection and metric definitions, ROI optimization, high-level database design)
- Strong in **production code** development, software engineering

Data scientists

More than data scientists?



# Big Data, Data Science, Statistics

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics**
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

# Big Data, Data Science, Statistics

## Wikipedia

### Big data

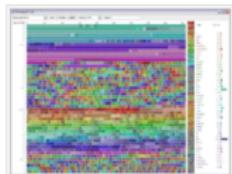
From Wikipedia, the free encyclopedia

This article is about large collections of data. For the band, see [Big Data \(band\)](#).

**Big data**<sup>[1][2]</sup> is the term for a collection of [data sets](#) so large and complex that it becomes difficult to process using on-hand database management tools or traditional data processing applications. The challenges include capture, curation, storage,<sup>[3]</sup> search, sharing, transfer, analysis<sup>[4]</sup> and visualization. The trend to larger data sets is due to the additional information derivable from analysis of a single large set of related data, as compared to separate smaller sets with the same total amount of data, allowing correlations to be found to "spot business trends, determine quality of research, prevent diseases, link legal citations, combat crime, and determine real-time roadway traffic conditions."<sup>[5][6][7]</sup>

As of 2012, limits on the size of data sets that are feasible to process in a reasonable amount of time were on the order of [exabytes](#) of data.<sup>[8]</sup> Scientists regularly encounter limitations due to large data sets in many areas, including [meteorology](#), [genomics](#),<sup>[9]</sup> [connectomics](#), complex physics simulations,<sup>[10]</sup> and biological and environmental research.<sup>[11]</sup> The limitations also affect [Internet search](#), [finance](#) and [business informatics](#). Data sets grow in size in part because they are increasingly being gathered by ubiquitous information-sensing mobile devices, aerial sensory technologies ([remote sensing](#)), software logs, cameras, microphones, [radio-frequency identification](#) readers, and [wireless sensor networks](#).<sup>[12][13]</sup> The world's technological per-capita capacity to store information has roughly doubled every 40 months since the 1980s;<sup>[14]</sup> as of 2012, every day 2.5 [exabytes](#) ( $2.5 \times 10^{18}$ ) of data were created.<sup>[15]</sup> The challenge for large enterprises is determining who should own big data initiatives that straddle the entire organization.<sup>[16]</sup>

Big data is difficult to work with using most relational database management systems and desktop statistics and visualization packages, requiring instead "massively parallel software running on tens, hundreds, or even thousands of servers".<sup>[17]</sup> What is considered "big data" varies depending on the capabilities of the organization managing the set, and on the capabilities of the applications that are traditionally used to process and analyze the data set in its domain. "For some organizations, facing hundreds of gigabytes of data for the first time may trigger a need to reconsider data management options. For others, it may take tens or hundreds of terabytes before data size becomes a significant consideration."<sup>[18]</sup>



A visualization created by IBM of Wikipedia edits. At multiple [terabytes](#) in size, the text and images of Wikipedia are a classic example of big data.

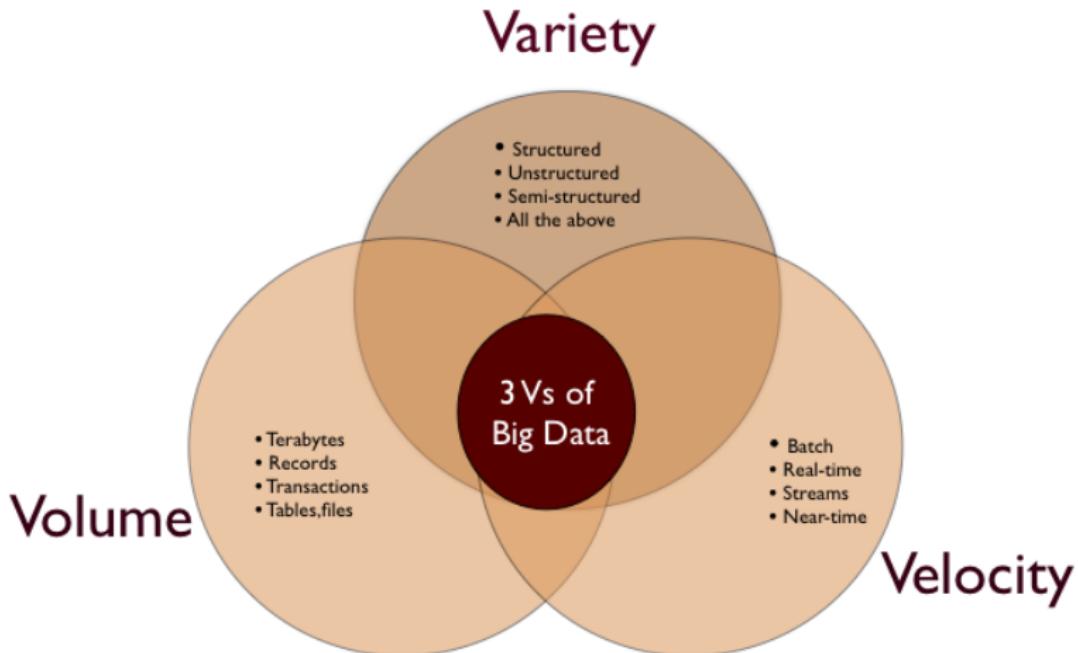
- **Big data** is an all-encompassing term for any collection of data sets so large and complex that it becomes difficult to process using traditional data processing applications.
- **Data science** is the study of the generalizable extraction of knowledge from data, yet the key word is science.
- **Statistics** is the study of the collection, analysis, interpretation, presentation and organization of data.

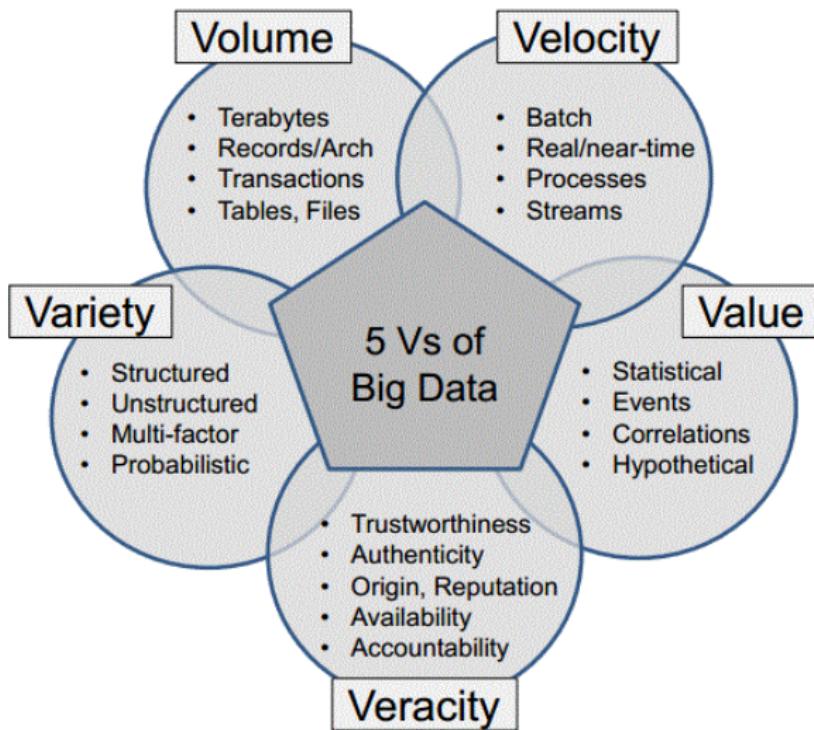
# Big Data, Data Science, Statistics

## Data science evolution

### Main Paradigmatic Changes in Big Data Analytics Environment

	Statistical Data Analysis <1985 <i>(Pure Statistical Inference)</i>	Business Intelligence 1985-2008 <i>(Constrained Data Mining)</i>	Big Analytics >2008 -up to now <i>(Unconstrained Data Mining)</i>
Data types	Homogeneous Structured Data (proprietary)	Homogeneous Structured & Homogeneous Unstructured Data, separately	Mix of Heterogeneous Unstructured & Structured Data (proprietary + open data)
Data storing	Line & column dimensions fixed Flat Files, Hierarchical DBs, & first Relational DBs	Column dimensions fixed SQL DBs: MySQL, DB2, ORACLE &OLAP Cubes	No dimensions fixed NoSQL DBs: Column oriented DBs, object oriented DBs etc.
Volume Cost/volume	<b>Exponential cost decrease</b>		<b>Exponential volume increase</b>
Basic Analytical Principles	Hypotheses driven mode: Power use of sampling Techniques	Mix Hypotheses driven & Data driven: Dimensions Reduction & Populations Segmentations	Full Data driven mode: Power use of learning techniques, mainly unsupervised
Main Algorithmic approaches	Regression Analysis, Factorial Analysis, Statistical Inference thru sampling, Linear general Models, Decision Trees.Etc.	Clustering (K- means, K Neighbours), Classification & Support Vector Machines Multi layers Neural Nets, Scoring Techniques, Sequential Patterns, etc.	Deep adaptive learning techniques, Auto encoded neural Nets Huge Graph Modularization, & Visual Analytics, Full unsupervised linear Clustering, etc.
New types of Business deliverables	Score Cards, Decisional Models based on sampling	Populations Profiling: CRM, Churn & Attrition Analysis, Loyalty & Propensity Programs,Cross selling	Near real time analysis for: "individually" adapted online marketing & sales, machine learning for various purposes, automated maintenance programs





# Big Data, Data Science, Statistics

Data science or statistics?

A vocabulary problem:

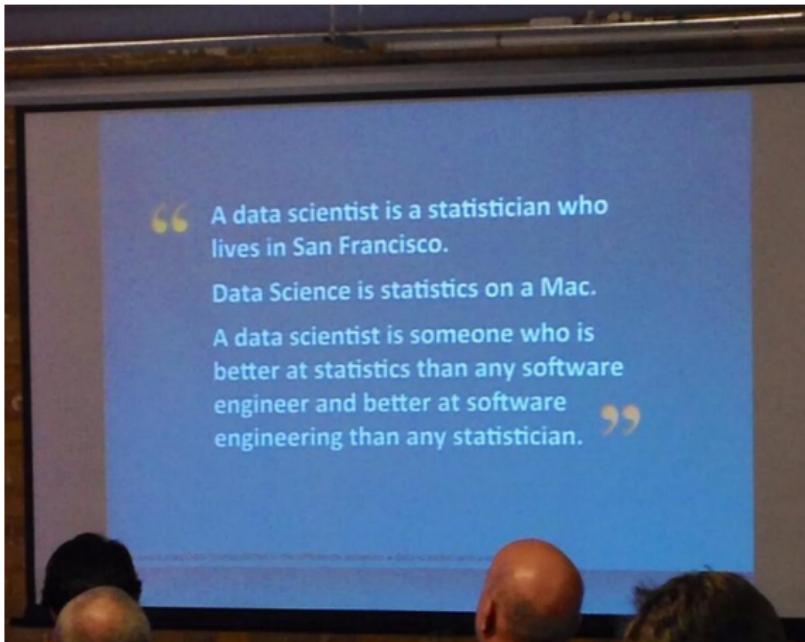
data scientist      or      statistician?

statistics      or      data science?

# Big Data, Data Science, Statistics

Data science or statistics?

A possible answer:



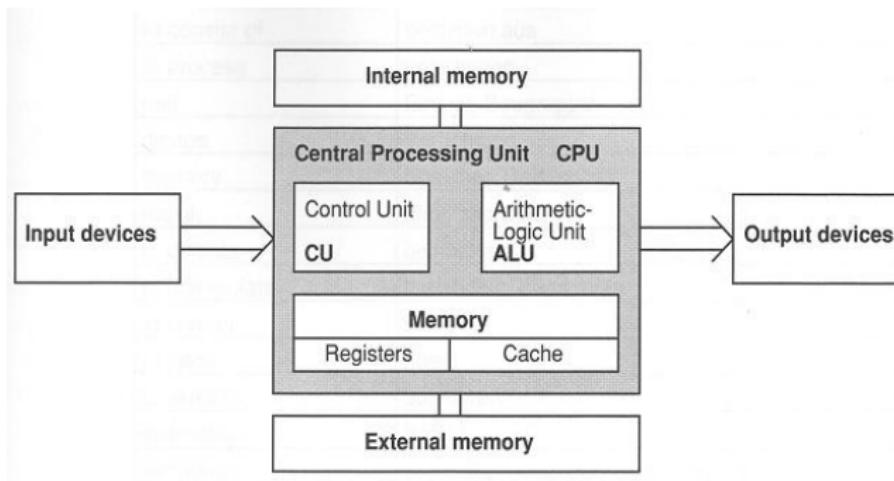
# Computing and Distributed Computing

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing**
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

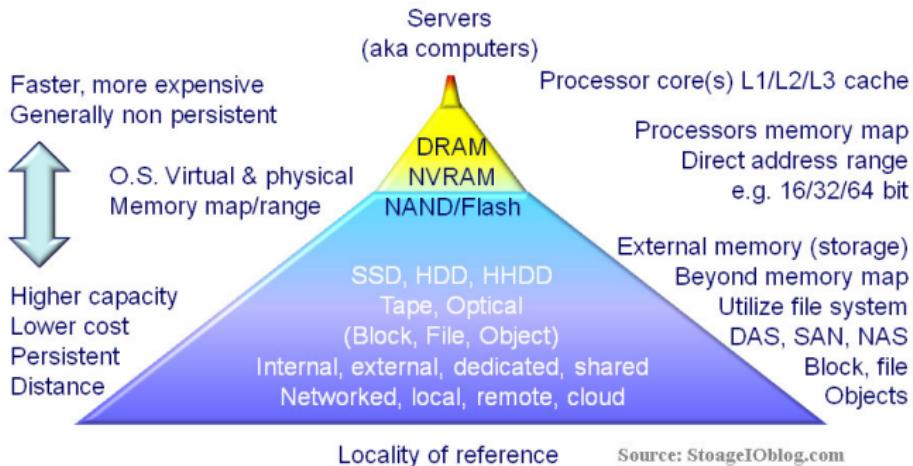
# Computing and Distributed Computing

## Computer Architecture



- Everything should go through the CPU...

# Computing and Distributed Computing Memories



CPU register	64 b × 16
Level 1 cache access	8-128 kb
Level 2 cache access	32-1024 kb
Level 3 cache access	1-8 MB
Main memory access	2-16 GB
Solid-state disk I/O	250 GB-1 TB (4TB)
Rotational disk I/O	500 GB-4 TB

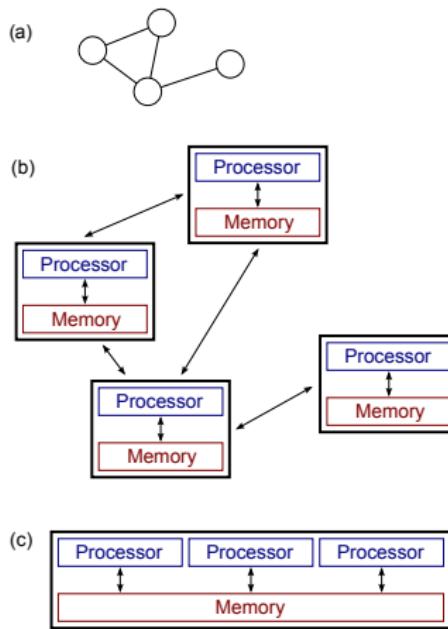
# Computing and Distributed Computing

## Memories

1 CPU cycle	0.3 ns	1 s
Level 1 cache access	0.9 ns	3 s
Level 2 cache access	2.8 ns	9 s
Level 3 cache access	12.9 ns	43 s
Main memory access	120 ns	6 min
Solid-state disk I/O	50-150 $\mu$ s	2-6 days
Rotational disk I/O	1-10 ms	1-12 months
Internet: SF to NYC	40 ms	4 years
Internet: SF to UK	81 ms	8 years
Internet: SF to Australia	183 ms	19 years
OS virtualization reboot	4 s	423 years
SCSI command time-out	30 s	3000 years
Hardware virtualization reboot	40 s	4000 years
Physical system reboot	5 m	32 millenia

# Computing and Distributed Computing

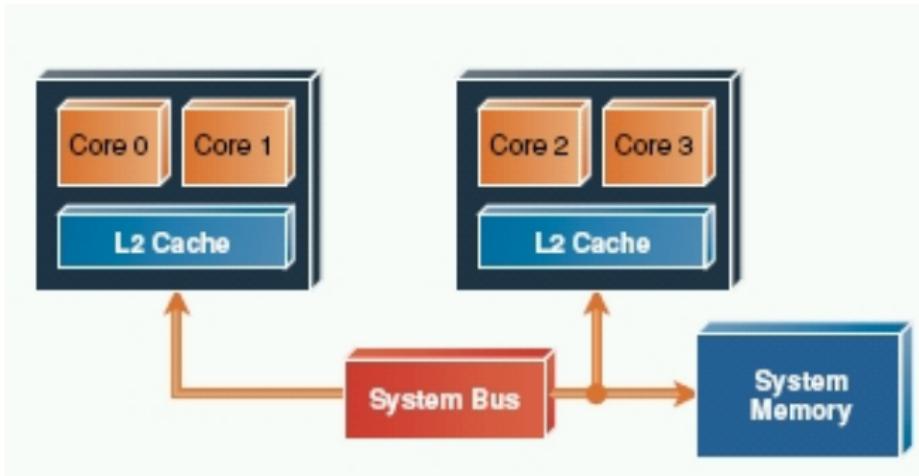
## Distributed/Parallel Computing



- Distributed (a/b)
- Parallel (c)

# Computing and Distributed Computing

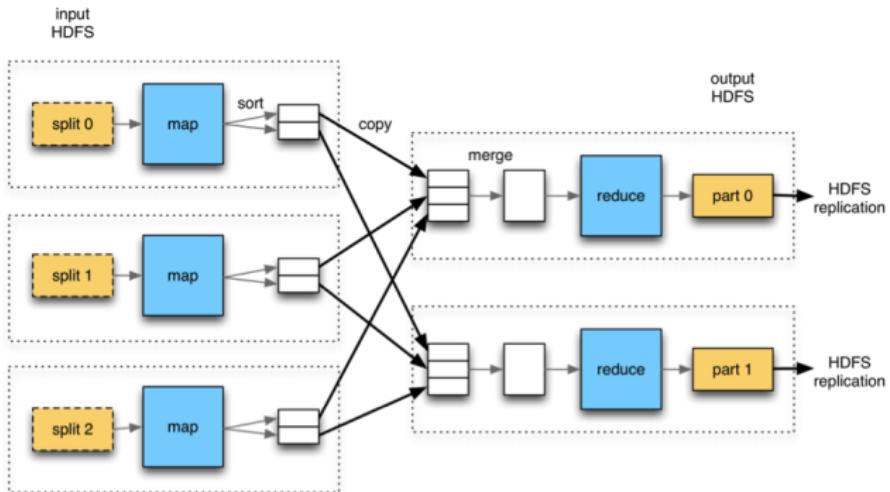
## MultiCore



- Several processors/cores with the same shared ram.
- No too expensive transfer between core.
- Strategies:
  - Independent batch
  - Parallelization technique limiting information transfer...
- System memory limitation!

# Computing and Distributed Computing

## Hadoop and Map/Reduce

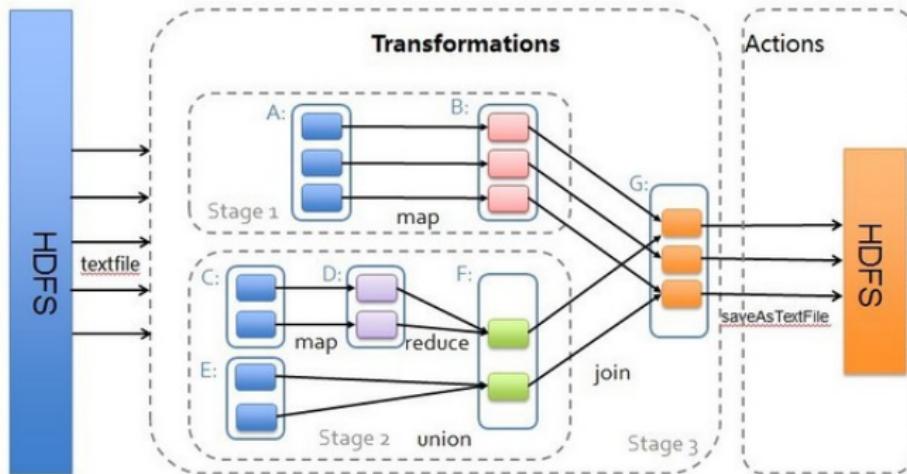


- Data transfer through disk and networked file system!
- Hadoop: Node failure handling and ecosystem.

# Computing and Distributed Computing

## Spark

### Spark: Transformations & Actions



- Strategy: keep everything as much as possible in memory...

# Computing and Distributed Computing

## GP-GPU



- Combine different processor types...
- CPU < DSP < FPGA < ASICS

# Data Science Challenges

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges**
- 11 Vocabulary of Data Science
- 12 Bibliography

## Data Science Challenges

### New Interdisciplinary Challenges

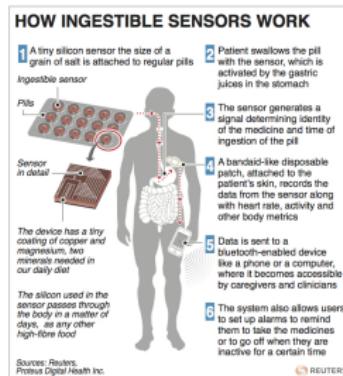
- Applied math **AND** Computer science
- Huge importance of domain specific knowledge: physics, signal processing, biology, health, marketing...

#### 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

# Data Science Challenges

## Data acquisition



## Some challenges

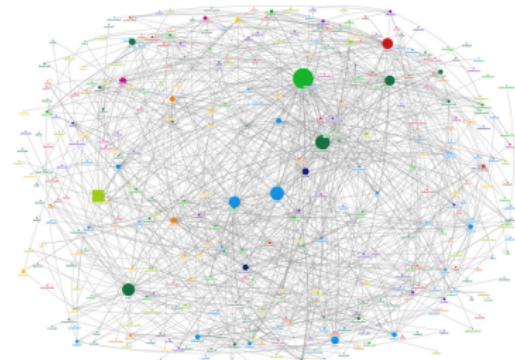
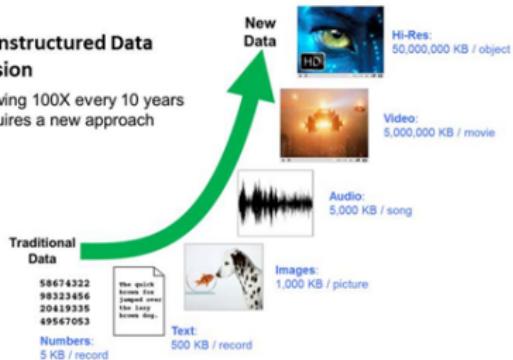
- How to measure new things?
- How to choose what to measure?
- How to deal with distributed sensors?
- How to look for new sources of informations?

# Data Science Challenges

## Unstructured Data

### The Unstructured Data Explosion

- Growing 100X every 10 years
- Requires a new approach

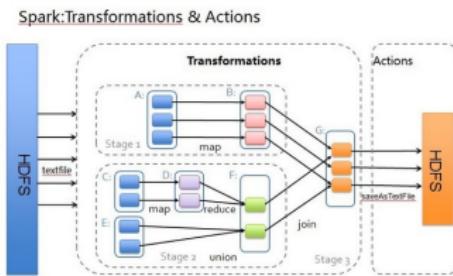
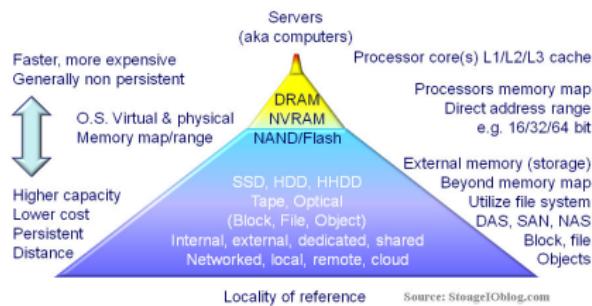


## Some challenges

- 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?

# Data Science Challenges

## Huge Dataset



## Some challenges

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

# Data Science Challenges

## High Dimensional Data

Main Paradigmatic Changes in Big Data Analytics Environment			
	Statistical Data Analysis <1985 (Pure Statistical Inference)	Business Intelligence 1985-2008 (Constrained Data Mining)	Big Analytics >2008 up to now (Unconstrained Data Mining)
Data types	Homogeneous Structured Data (proprietary)	Homogeneous Structured & Homogeneous Unstructured Data, separately	Mix of Heterogeneous Unstructured & Structured Data (proprietary + open data)
Data storing	Line & column dimensions fixed Flat Files, Hierarchical DBs, & first Relational DBs	Column dimensions fixed SQL DBs: MySQL, DB2, ORACLE & OLAP Cubes	No dimensions fixed NoSQL DBs; Column oriented DBs, object oriented DBs etc.
Volume Cost/volume	<b>Exponential cost decrease</b>		
Basic Analytical Principles	Hypotheses driven mode: Power use of sampling Techniques	Mix Hypotheses driven & Data driven: Dimensions Reduction & Populations Segmentation	Full Data driven mode: Power use of learning techniques, mainly unsupervised
Main Algorithmic approaches	Regression Analysis, Factorial Analysis, Statistical Inference thru sampling, Linear general Models Decision Trees, Etc.	Clustering (K-means, K Neighbours), Classification & Support Vector Machines Multi layers Neural Nets, Scoring Techniques, Sequential Patterns, Populations Profiling: CRM, Churn & Attrition Analysis, Loyalty & Propensity Programs, Cross selling	Deep adaptive learning techniques, Auto encoded neural Nets Huge Graph Modularization, & Visual Analytics, Full unsupervised linear Clustering, etc.
New types of Business deliverables	Score Cards, Decisional Models based on sampling		Near real time analysis for “individuality” adapted online marketing & sales, machine learning for various purposes, automated maintenance programs

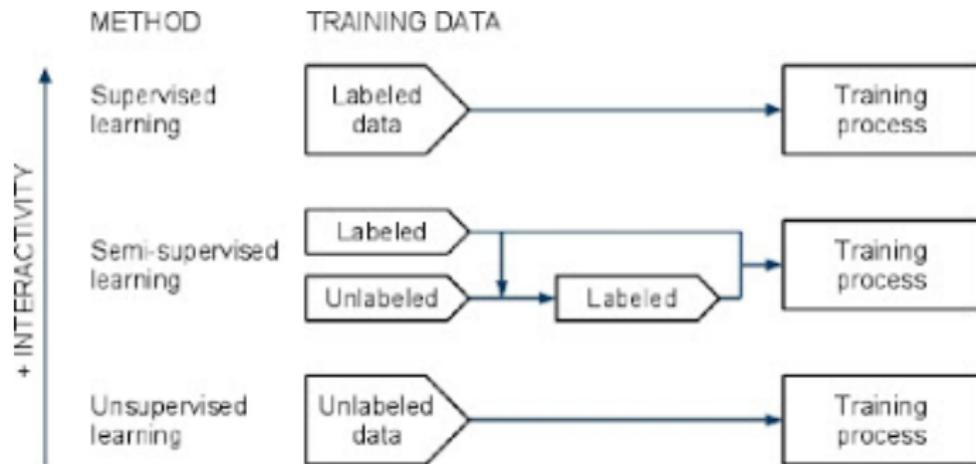
THALES

## Some challenges

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

# Data Science Challenges

## Learning and Supervision

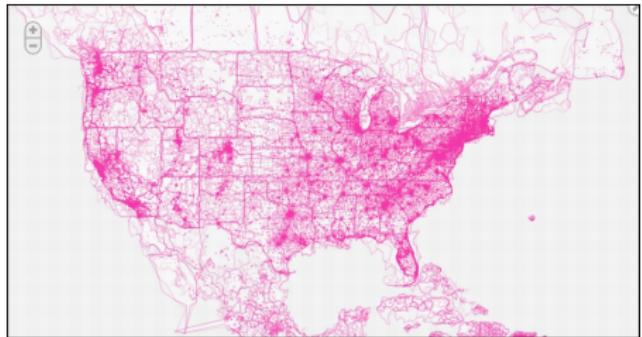
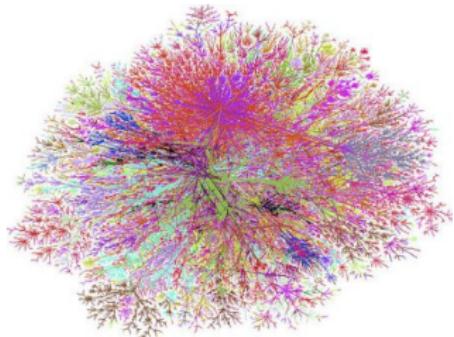


## Some challenges

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

# Data Science Challenges

## Visualization



### Some challenges

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

# Vocabulary of Data Science

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

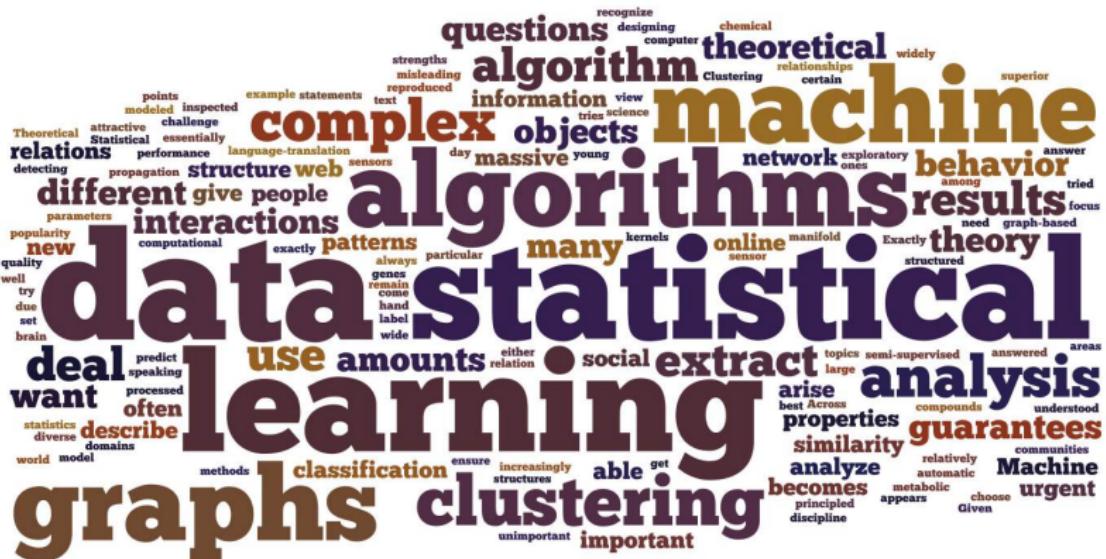
# Vocabulary of Data Science

Lots of words



# Vocabulary of Data Science

## Lots of words



# Vocabulary of Data Science

## Main Fields

**Data mining.** Extract patterns from data by combining methods from statistics, machine learning and data processing technologies.

*Example:* market basket analysis to model the purchase behavior of customers.

**Machine learning.** Design and develop algorithms allowing computers to learn from data, in order to take intelligent decisions automatically. *Example:* Natural language processing.

**Statistics.** Collection, organization, and interpretation of data. Mathematical methods to construct quantitative assessments of errors and risks when taking decisions, estimating parameters and doing predictions. *Example:* quantitative assessments of relationships between variables, computing confidence intervals for model parameters, hypothesis testing.

# Vocabulary of Data Science

## Main Fields

**Natural language processing (NLP).** Specialization of machine learning and linguistics that builds algorithms to analyze human (natural) language. *Example:* sentiment analysis on social networks.

**Network analysis.** Characterize relationships among nodes in a graph or a network, understand the communities, the influence of nodes on the others, understand how information travels in the network. *Example:* identify key opinion leaders in a social network, identify the information flows in a large company

**Predictive modeling.** Use of a mathematical model to predict an outcome, e.g. regression, classification, etc. *Example.* Predict the probability that a customer will churn.

# Vocabulary of Data Science

## Main Fields

**Supervised learning.** Machine learning techniques that infer a function or a relationship from a set of training data. *Examples:* classification, regression

**Unsupervised learning.** Machine learning techniques that finds structure in unlabeled data. *Example:* clustering is a part of unsupervised learning

**Visualization.** Techniques used for representation of data by creating images, diagrams, animations, in order to communicate, understand, explore and improve understanding of data.

# Vocabulary of Data Science

## Machine Learning

**Labels.** Characteristics / categories of interest in points of data.  
This is the information one wants to predict in supervised learning.

**Features.** A set of information about a point of data (a customer, a company, a country, etc.)

**Clustering.** Algorithms used in unsupervised learning to assign a group to each data point. Groups are called clusters. *Example:* customer segmentation in a e-commerce platform

**Classification.** Algorithms used in supervised learning to predict the labels of data points. It relies on the training of a model or algorithm using labeled data. *Keywords:* Logistic Regression, SVM, CART, Boosting, etc.

**Feature selection.** Algorithms in machine learning that select features that best explain an outcome. *Example.* In biology, find the genetic informations that best explains a patient's response to a drug.

# Vocabulary of Data Science

## Mathematical Concepts

**Generalization.** Ability of a predictive algorithm to generalize: give good predictive results on a sample different than to one used to train the algorithm.

**Parameters.** A set of coefficients (vectors, matrices), that specifies a model. *Example:* the mean and standard deviation of a Gaussian distribution

**Statistical model.** A mathematical formulation that attempts to explain how data is generated. *Example:* data is generated by a multivariate Gaussian distribution

**Likelihood.** The probability that data is generated by a model for some parameters choice

# Vocabulary of Data Science

## Mathematical Concepts

**Optimization.** Study and design of numerical algorithms used to (but not only) minimize or maximize functions. *Example:* optimization of a likelihood is called the **training** step in machine learning.

**Goodness-of-fit.** A quantity that assesses the closeness of a (trained) model to data. *Example.* The least-squares error for linear regression.

**Over-fitting.** Something that must be avoided to have a good predictive performance on new data.

**Cross-validation.** Splitting of data into several subsets. A model is trained on a subset and tested on another, to check its goodness-of-fit both on data used for training, but on new data as well.

**Structured data.** Data structured in fixed fields. *Example:* relational databases or excel spreadsheets.

**Semi-structured data.** Data not structured in fixed fields but contain markers to separate data elements. *Example:* XML or HTML-tagged text.

**Unstructured data.** Data not structured in fixed fields. *Example:* books, articles, body of e-mail messages, audio, image and video data, etc.

**Metadata.** Data that describes the content and context of data: creation, purpose, time and date, author, etc.

**Data fusion and data integration.** Set of techniques that integrate and analyze data from multiple sources, instead of analyzing single sources of data. *Example:* combine analysis of social network data with NLP and real-time sales data, to assess the effect of a marketing campaign on customer sentiment and purchases

**Cloud computing.** A computing paradigm where computing resources are configured as a distributed system, which provides a service through a network.

**Distributed system.** Several computers, communicating through a network, used to solve a common storing or computational problem. Aim is higher performance at a lower cost, higher reliability and scalability.

**Relational database.** A database consisting of collections of tables (relations), namely data are stored in rows and columns. SQL is the most widely used language for managing relational databases.

**Non-relational database.** A database that does not store data in tables (rows and columns).

## Vocabulary of Data Science

### Databases and Data Processing

**SQL.** Acronym for Structured Query Language. It is a computer language designed for managing data in relational databases.

*Example.* insert, query, update, delete data, manage database structures, and control access to data in the database.

**NoSQL.** A group of database management systems. Data is not stored in tables like in Relational database. It does not rely on the mathematical relationship between tables. It gives a way of storing and retrieving unstructured data quickly.

**Hbase.** A distributed, non-relational database. It is managed as a project of the Apache Software foundation and a part of Hadoop.

## Vocabulary of Data Science

### Databases and Data Processing

**Hadoop.** A framework that supports large scale data processing by allowing the decomposition of large tasks into smaller tasks, that are executed in parallel, on independently slices of the data and then finally merged to answer to the task.

**MapReduce.** A software framework introduced by Google for processing huge datasets on a distributed system. Implemented in Hadoop. It supports large scale data processing by decomposing large tasks into smaller tasks, executed in parallel, on independent parts of data and finally merged to answer to the task.

**Stream processing.** Technologies designed to process large real-time streams of event data. *Example:* high-frequency algorithmic trading, analysis of Twitter data streams

# Bibliography

## Outline

- 1 Data Science in the media
- 2 From Data to Product
- 3 An example: Real Time Bidding
- 4 Data Science ecosystem
- 5 Data cycle
- 6 Data Science project
- 7 Data scientists
- 8 Big Data, Data Science, Statistics
- 9 Computing and Distributed Computing
- 10 Data Science Challenges
- 11 Vocabulary of Data Science
- 12 Bibliography

# Bibliography

## Bibliography

-  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*