

# House of Finance Days 2016

## Machine Learning and Finance

---

Introduction to machine learning, application to  
Hawkes processes

S. Gaïffas



## Stéphane Gaïffas

- Professeur Chargé de Cours
- CMAP, Ecole polytechnique
- <http://www.cmap.polytechnique.fr/~gaiffas/>
- [stephane.gaiffas@polytechnique.edu](mailto:stephane.gaiffas@polytechnique.edu)
- Enseignement + Recherche
- *mots-clés*: Machine learning, Statistique, Big-data

- 1 Teasers
  - Data Science in the media
  - Examples
  - Big Data is (quite) Easy ?
- 2 Supervised learning
  - Introduction
  - An example: RTB
  - Loss functions, linearity
  - Logistic regression
  - Penalization
- 3 Optimization
  - Introduction
  - Gradient descent
  - Stochastic gradient descent
  - Variance reduction
  - Conclusion
- 4 Hawkes processes
  - Introduction
  - Model
  - Large dimension
  - Maximum Likelihood Estimation
  - Mean-Field approximation

- 1 Teasers
  - Data Science in the media
  - Examples
  - Big Data is (quite) Easy ?
- 2 Supervised learning
  - Introduction
  - An example: RTB
  - Loss functions, linearity
  - Logistic regression
  - Penalization
- 3 Optimization
  - Introduction
  - Gradient descent
  - Stochastic gradient descent
  - Variance reduction
  - Conclusion
- 4 Hawkes processes
  - Introduction
  - Model
  - Large dimension
  - Maximum Likelihood Estimation
  - Mean-Field approximation

# Le Big Data, nouvel eldorado des entreprises

Par Direct Matin, publié le 26 Septembre 2014 à 08:32



Les mégadonnées représentent un marché de plusieurs milliards d'euros[© infocux technologies]

**Considéré comme le "nouveau pétrole du XXIe siècle", le Big data attise toutes les convoitises.**

EN COMPLÉMENT



# M Idées

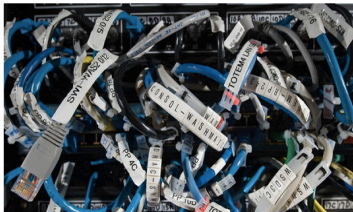
## Les données, puissance du futur

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 f t v in

Recommander Envoyer 467 personnes le recommandent.



### Les plus partagés

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

### Nous suivre

Retrouvez le meilleur de notre communauté



# 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 grieving**

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 »

The New York Times  
**Sunday Review** | The Opinion Pages

Search All NYTimes.com

NEWS ANALYSIS

## The Age of Big Data

By STEVE LOHR

Published: February 11, 2012 | 82 Comments

GOOD with numbers? Fascinated by data? The sound you hear is opportunity knocking.



Mo Zhou was snapped up by I.B.M. last summer, as a freshly minted Yale M.B.A., to join the technology company's fast-growing ranks of data consultants. They help businesses make sense of an explosion of data — Web traffic and social network comments, as well as software and sensors that monitor shipments, suppliers and customers — to guide decisions, trim costs and lift sales. "I've always had a love of numbers," says Ms. Zhou, whose job as a data analyst suits her skills.

Multimedia

To exploit the data flood, America will need many more like her. A report last year by the [McKinsey Global Institute](#), the

RECOMMEND

TWITTER

LINKEDIN

COMMENTS (82)

SIGN IN TO E-MAIL

PRINT

REPRINTS

SHARE

Log in to see what your friends are sharing on nytimes.com. [Privacy Policy](#) | [What's This?](#)

Log In With Facebook

### What's Popular Now

Despite New Health Law, Some See Sharp Rise in Premiums



The Big Fail



MOST E-MAILED

RECOMMENDED FOR YOU

1. OFF THE DRIBBLE  
 Studemire Commemorates Brother's Death



2. CRITIC'S NOTEBOOK  
 The Rainbow That Follows 'Jersey Shore'

3. TAKING NOTE  
 Opinion Report: Tax Reform

4. THE LEARNING NETWORK  
 Fill-In | Trendy Spot Urges Tourists to Ride In and Spend, 'Gangnam Style'





THE WORLD BANK  
Working for a World Free of Poverty

English | Español | Français | العربية | Русский | 中文 | ▶

Search  GO



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...

?  
**BIG DATA**

**#bigstats**

December 19th 2.30pm  
World Bank HQ  
MC13 -121

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: \*



## Focus on the issues

[Deloitte Research](#)[Deloitte University Press](#)[Books](#)[Email Alerts](#)[Podcasts](#)[Tools](#)[Video library](#)[Browse by industry](#)[Browse by service](#)

## Billions and billions: Big data becomes a big deal



### The podcast

#### Deloitte global podcasts

Big data becomes a big deal

To use our embedded media player, please install the latest version of **Adobe Flash Player**. You can also **download the podcast file**.

Big data projects had a total industry revenue of only \$100 million in 2009. However 2012 will see 90 per cent of Fortune 500 companies kick off a big data initiative, which will trigger industry revenue of between \$1 billion and 1.5 billion. Big data is still in its infancy, mostly used for meteorology and physics simulations, but interest is gaining pace as data warehouses start to overflow and the need for "real-time" analysis puts strain on traditional analytics tools. Internet companies have led the way with exploring big data but fast follower sectors are likely to include the public sector, financial services, retail, entertainment, and media. This could trigger a talent shortage with up to 190,000 skilled professionals needed to cope with demand in the US alone over the next five years. Meanwhile companies launching initiatives need to take a disciplined and targeted approach to big data.

#### Podcast highlights:

- What does "big data" mean?
- Where will the industry growth come from?
- What does the trend mean for traditional data companies?
- What does the accessibility of "big data" mean for the way companies are currently doing business?

#### Related links

[Read the Prediction](#)[More Technology Predictions](#)

#### Stay connected

[Contact us](#)[Submit RFP](#)[Global blog](#)[Global podcasts](#)[Social media](#)[RSS feed](#)

Accueil &gt; Actualités &amp; Évènements

## 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 m2 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

Criteo inaugure à Paris l'un des plus gros pôles européens de R&D dédiés à la publicité prédictive, CriteoLabs. Sur 10.000 m2, 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



Vieilles photos de CriteoLabs lors de l'inauguration de Fleur Pellerin

# Data is the new oil?

## "DATA IS THE NEW OIL"

From the beginning of recorded time until 2011, we created **5 exabytes** of data.

In 2011 the same amount was created every two days.

By 2015, it's expected that the time will shrink to 15 minutes.

Every hour, we create enough Internet traffic to fill **7 billion DVDs**.

Side by side, that's like a superhighway the height of Everest.

Collected in 2009 by Clive Maxfield, a British data commercialization entrepreneur, this now famous phrase was embraced by the World Economic Forum in a 2011 report, which considered data to be an economic asset, like oil.

There are nearly as many bits of information in the digital universe as there are stars in our entire universe.

As of August 2012, there were just over **4 million** videos in the English language.

There are **133 million BLOGS** on this web.

**80%** of all business.com domains point to white sites. Out of \$3 billion in advertising, \$1 billion are spent on those 40% of white and nonwhite sites.

English is the dominant language of the web. But by 2014 it will be **Chinese**, if its current rate of increase continues.

Top language used on the web May 2013



**247 billion EMAILS** are sent every day (24 in 2010, and nearly 1 billion in 2011).



Just as a study of activity on Twitter gave residents, family members, and journalists advance warning of details about the devastating earthquake and tsunami in Japan,

**high-frequency traders**, with the help of computer algorithms, use Big Data to follow trends and to act quickly on their findings.

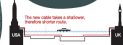
These specialized algorithms make split-second decisions to buy or sell a commodity. These cables being laid under the Atlantic will shave **5 milliseconds** from the current 60 milliseconds it takes for trading information to travel between New York City and London.

With new fiber-optic cable, the round-trip time between New York and London will be 50 milliseconds.

This 5-millisecond saving is worth many millions of dollars to the trading firms who use the cables and will pay millions to do so.

How they save 5 milliseconds

The depth of the Atlantic Ocean varies. The new cable will lie on the sea floor that are up to 1,000 feet shallower than the current fastest cable. By taking a different route, the new cable is shorter, meaning that the time it takes for messages to travel along it is shortened.



**60%** of all tweets (6.4 billion per day) are instant messages. In 2011, 100,000 text messages were sent every second.

**10%** of all photos ever taken were taken in 2011.

**50%** of all e-mail sent in the U.S. are spam.





**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 TechnAsia, le **moteur de recherche** serait bloqué depuis le 4 septembre dans le pays. DuckDuckGo, qui se présente ...

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

## More Ideas Based on Your Browsing History

You looked at



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

> [Find similar items](#)

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**

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

# Display advertising

## Outlet

> Descubrela

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

Libros universitarios  
y de estudios  
superiores  
a precios bajos



> Descubrelas

**Neuheiten**  
von Makita



> Hier klicken

**fire** + 12 MONTHS  
OF PRIME  
PHONE



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



Économisez  
en vous Abonnant > [Cliquez ici](#)

**PROMOTIONS  
CHAUSSURES**

-30% -40% -50%...



> [J'en profite](#)

**PROMOTIONS  
SACS À MAIN**



> [J'en profite](#)



A background map of a city street grid with several red squares overlaid, indicating predicted crime hotspots. A red-bordered inset on the left shows a zoomed-in view of a residential block with these squares.

# PREDICTIVE POLICING®

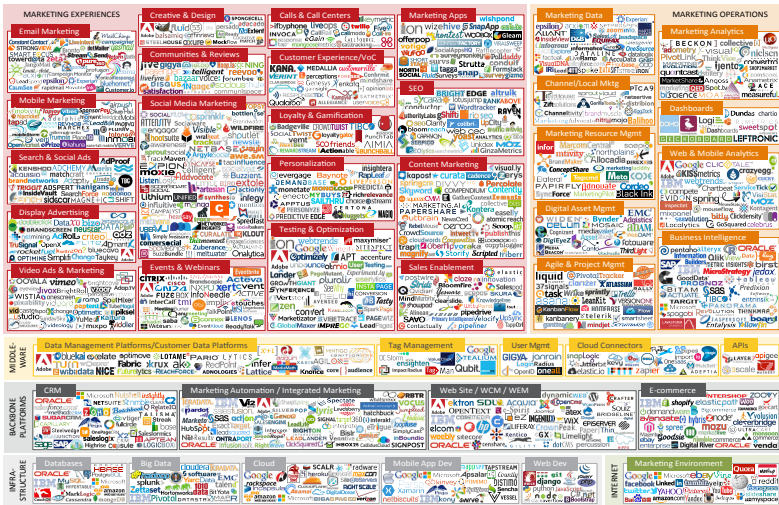
The Predictive Policing Company.

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.



## chiefmartec.com Marketing Technology Landscape

January 2014

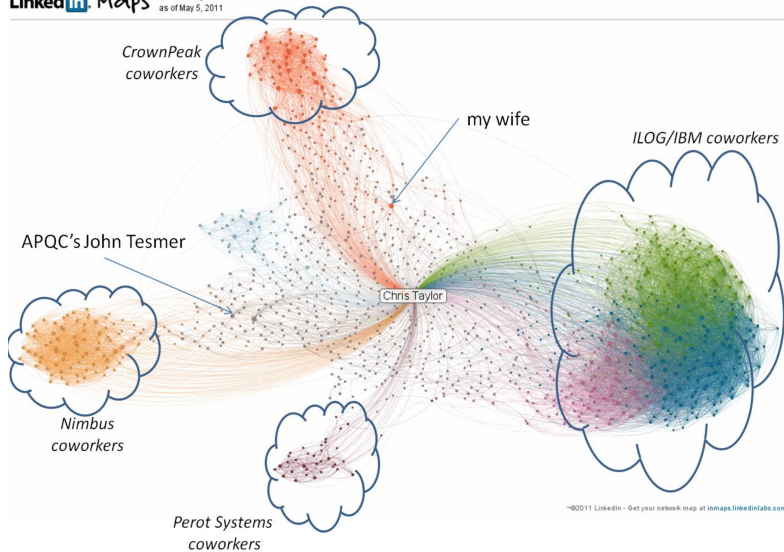




# Social networks analytics

LinkedIn Maps

Chris Taylor's Professional Network  
as of May 5, 2011



## 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

**\$57 Trillion**

in infrastructure investments will be needed between 2013-2030.

Source: McKinsey

### Transportation Analytics

**50 Hours**

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

**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.

Cloud is driving cities in their digital transformation.

### 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

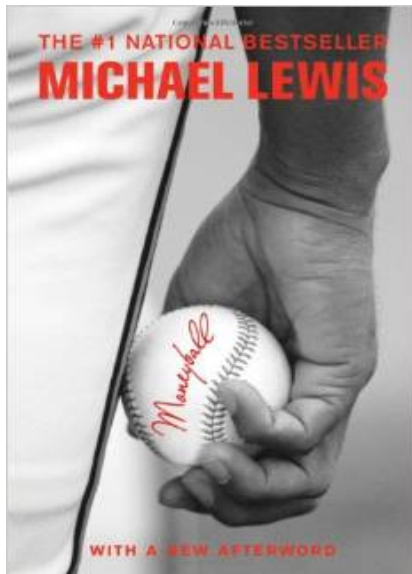
### Open Cloud

**\$6 Billion**

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

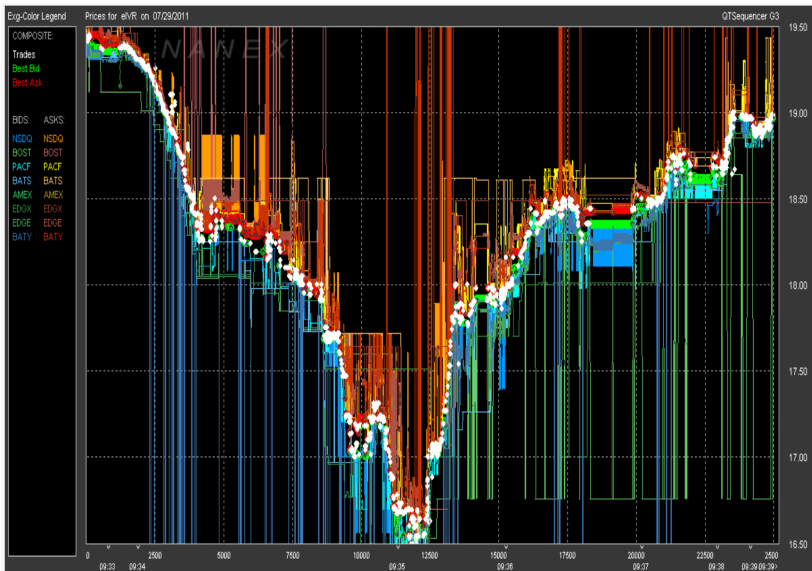
IBM Intelligent Operations software is designed with cities, for cities, to provide the tools to monitor, visualize and analyze vital city services such as water and wastewater systems, transportation, infrastructure planning, permit management and emergency response.







# High frequency trading



## Big data

- Capacity to store information has doubled every 40 months since the 1980s
- In 2012, 2.5 exabytes ( $2.5 \times 10^{18}$ ) created per **day**
- Big companies such as Google, Amazon, Facebook, Apple (GAFA) but also banking, marketing, pharmaceuticals, insurance, telecoms, personalized medicine, bioinformatics, etc.

## 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()
    .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()
}
```

# 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!

# A Complex Ecosystem

## Big Data Landscape



Matt Turk (@mattturk) and Shivon Zilis (@shivonz)

# Data science or statistics?



**Jeremy Jarvis**

@jeremyjarvis

 Follow

"A data scientist is a statistician who lives in San Fransisco"

[#monkigras](#) [pic.twitter.com/HypLL3Cnye](http://pic.twitter.com/HypLL3Cnye)

12:13 PM - 30 Jan 2014

  1,475  841



**Big Data Borat**

@BigDataBorat

 Follow

Data Science is statistics on a Mac.

3:32 PM - 27 Aug 2013

  611  273



**Josh Wills**

@josh\_wills

 Follow

Data Scientist (n.): Person who is better at statistics than any software engineer and better at software engineering than any statistician.

6:55 PM - 3 May 2012

  1,360  821

- 1 Teasers
  - Data Science in the media
  - Examples
  - Big Data is (quite) Easy ?
- 2 Supervised learning
  - Introduction
  - An example: RTB
  - Loss functions, linearity
  - Logistic regression
  - Penalization
- 3 Optimization
  - Introduction
  - Gradient descent
  - Stochastic gradient descent
  - Variance reduction
  - Conclusion
- 4 Hawkes processes
  - Introduction
  - Model
  - Large dimension
  - Maximum Likelihood Estimation
  - Mean-Field approximation

## Setting

- Data  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$  for  $i = 1, \dots, n$
- $x_i$  is the *features* of  $i$
- $y_i$  is the *label* of  $i$
- $y_i \in \mathbb{R}$  (regression)  $y_i \in \{-1, 1\}$  (binary classification)
- Usually, assume  $(x_i, y_i)$  are i.i.d

## Aim

- Based on  $(x_i, y_i)$ , learn a function that predicts  $y$  based on a new  $x$  (generalization property)

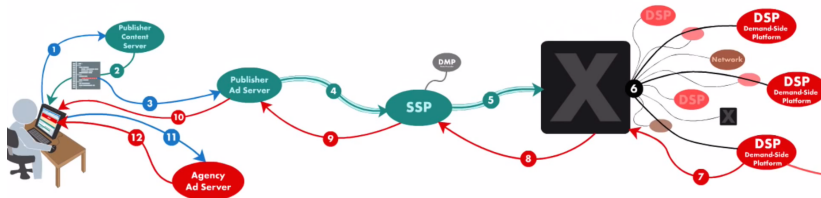
## Scaling

- *High-dimension*:  $d$  is large, say  $d \geq 10^4$
- *Big data*:  $n$  is large, say  $n \geq 10^6$

## Scenarios

- $d$  is large,  $n$  is small: computational biology
- $d$  is small,  $n$  is large: marketing
- $d$  is large,  $n$  is large: web-advertisement, ad display

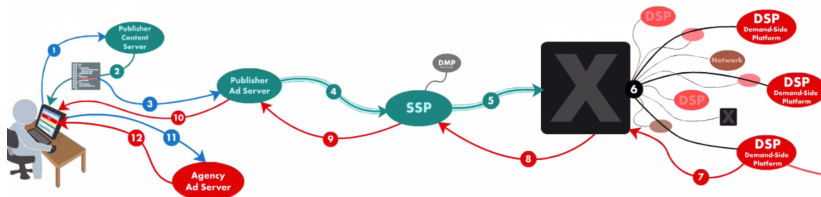
# An example: Real Time Bidding



- 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 on the webpage where the user is going. The webpage belongs to a **publisher**.
- The publisher sells ad space to advertisers who want to reach customers

In some cases, an auction starts: **RTB** (Real Time Bidding)

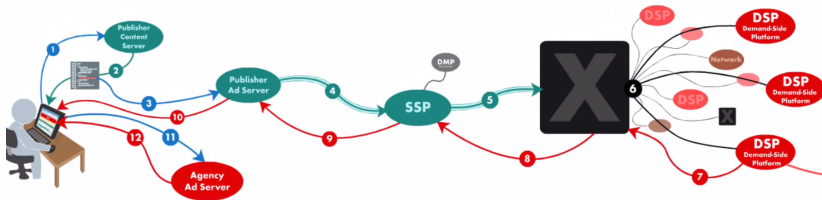
# An example: Real Time Bidding



- Advertisers have **10ms (!)** to give a price: they need to assess quickly how willing they are to display the ad to this customer
- Machine learning is used here to predict the probability of click on the ad. Time constraint: few model parameters to answer quickly
- feature selection / dimension reduction is crucial here

Full process takes  $< 100\text{ms}$

# An example: Real Time Bidding



## Some figures

- 10 million prediction of click probability per second
- answers within 10ms
- stores 20Terabytes of data daily

## What to do ?

Minimize with respect to  $f : \mathbb{R}^d \rightarrow \mathbb{R}$

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

where

- $\ell$  is a **loss** function:  $\ell(y_i, f(x_i))$  small means  $y_i$  is close to  $f(x_i)$
- $R_n(f)$  is called **goodness-of-fit** or **empirical risk**

Computation of  $f$  is called **training** or **estimation**

A problem

- $n$  and  $d$  are large: training is too time-consuming for a complex function  $f$

A simplification

- Choose a **linear** function  $f$  :

$$f(x) = x^\top \theta = \sum_{j=1}^d x_j \theta_j,$$

for a parameter  $\theta \in \mathbb{R}^d$  to be **trained**

Remark

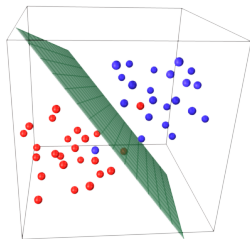
- linear with respect to  $x_j$ , but **you** choose the features  $x_j$
- usually not linear w.r.t the raw features: **feature engineering**

## Logistic regression

- The most widespread approach
- Assumes that the **log-odd ratio** is linear:

$$\log \left( \frac{\mathbb{P}[Y = 1|X = x]}{\mathbb{P}[Y = 0|X = x]} \right) = \langle x, \theta \rangle$$

- Leads to a **linear** separation between the 1s and –1s



Goodness of fit =  $-\log$ -likelihood, equal in this case to

$$R_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i x_i^\top \theta}).$$

So, let's just find out (more on that later...)

$$\hat{\theta} \in \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} R_n(\theta).$$

Now, given a new  $x$ , predict  $y$  using

$$\hat{y} = \operatorname{sign}(x^\top \hat{\theta})$$

End of story ?

- No !
- Overfitting problem

Minimizing only

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i^\top \theta)$$

is generally a bad idea. Minimize instead

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i^\top \theta) + \lambda \operatorname{pen}(\theta) \right\}$$

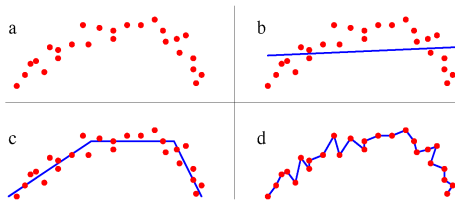
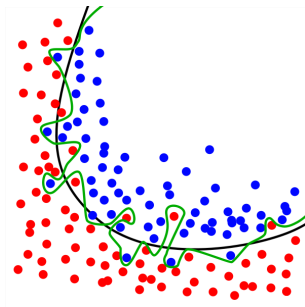
where

- $\operatorname{pen}$  is a **penalization** function, it forbids  $\theta$  to be “too complex”
- $\lambda > 0$  is a **tuning** or **smoothing** parameter, that **balances** goodness-of-fit and penalization

Why using penalization?

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i^\top \theta) + \lambda \operatorname{pen}(\theta) \right\}$$

Penalization, for a well-chosen  $\lambda > 0$ , allows to avoid **overfitting**



Penalization most widely used is

$$\text{pen}(\theta) = \frac{1}{2} \|\theta\|_2^2 = \frac{1}{2} \sum_{j=1}^d \theta_j^2.$$

Penalizes the energy of  $\theta$ , measured by squared  $\ell_2$ -norm

## Sparsity inducing penalization.

- It would be nice to find a model where  $\hat{\theta}_j = 0$  for many coordinates  $j$
- few features are useful for prediction, model is simpler, faster prediction (e.g. RTB example)
- We say that  $\hat{\theta}$  is **sparse**
- How to do it ?

Tempting to use

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i^\top \theta) + \lambda \|\theta\|_0 \right\},$$

where

$$\|\theta\|_0 = \#\{j : \theta_j \neq 0\}.$$

But, to do it exactly, you need to try **all** possible subsets of non-zero coordinates of  $\theta$ :  $2^d$  possibilities. Impossible!

A solution:  $\ell_1$ -penalization

$$\text{pen}(\theta) = \|\theta\|_1 = \sum_{j=1}^d |\theta_j|.$$

- Convex relaxation principle. Also called Lasso
- In a noiseless setting, in a certain regime,  $\ell_1$ -minimization gives the “same solution” as  $\|\cdot\|_0$

Why do  $\ell_1$ -penalization leads to sparsity?

Solution of

$$\underset{a \in \mathbb{R}}{\text{argmin}} \left\{ \frac{1}{2}(a - b)^2 + \lambda|a| \right\},$$

for  $\lambda > 0$  and  $b \in \mathbb{R}$  is given by

$$a_* = \text{sign}(b)(|b| - \lambda)_+$$

where  $a_+ = \max(0, a)$ .

A minimizer

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(y_i, \mathbf{x}_i^\top \theta) + \lambda \|\theta\|_1 \right\}$$

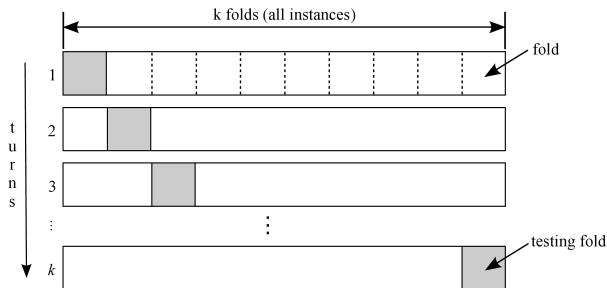
is typically sparse ( $\hat{\theta}_j = 0$  for many  $j$ ).

- for  $\lambda$  large (larger than some constant)  $\hat{\theta}_j = 0$  for all  $j$
- for  $\lambda = 0$  then there is no penalization
- Between the two, the “sparsity” depends on the value of  $\lambda$ :  
once again, it is a regularization or penalization parameter

How to choose it?

## V-Fold cross-validation

- Most standard cross-validation technique
- Take  $V = 5$  or  $V = 10$ . Pick a random partition  $I_1, \dots, I_V$  of  $\{1, \dots, n\}$ , where  $|I_v| \approx \frac{n}{V}$  for any  $v = 1, \dots, V$



For each  $v = 1, \dots, V$

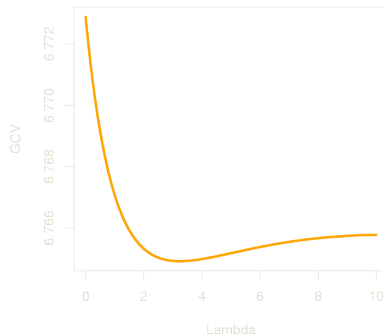
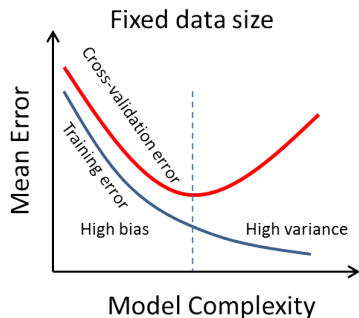
- Put  $D_{v,\text{train}} = \cup_{v' \neq v} I_{v'}$  and  $D_{v,\text{test}} = I_v$
- Find

$$\hat{\theta}_{v,\lambda} \in \underset{\theta}{\operatorname{argmin}} \left\{ \frac{1}{|D_{v,\text{train}}|} \sum_{i \in D_{v,\text{train}}} \ell(y_i, x_i^\top \theta) + \lambda \operatorname{pen}(\theta) \right\}$$

Take

$$\hat{\lambda} \in \underset{\lambda}{\operatorname{argmin}} \sum_{v=1}^V \sum_{i \in D_{v,\text{test}}} \ell(y_i, x_i^\top \hat{\theta}_{v,\lambda})$$

# Cross-validation



- Training error:

$$\lambda \mapsto \sum_{v=1}^V \sum_{i \in D_{v,\text{train}}} \ell(y_i, x_i^\top \hat{\theta}_{v,\lambda})$$

- Testing, validation or cross-validation error:

$$\lambda \mapsto \sum_{v=1}^V \sum_{i \in D_{v,\text{test}}} \ell(y_i, x_i^\top \hat{\theta}_{v,\lambda})$$

- 1 Teasers
  - Data Science in the media
  - Examples
  - Big Data is (quite) Easy ?
- 2 Supervised learning
  - Introduction
  - An example: RTB
  - Loss functions, linearity
  - Logistic regression
  - Penalization
- 3 Optimization
  - Introduction
  - Gradient descent
  - Stochastic gradient descent
  - Variance reduction
  - Conclusion
- 4 Hawkes processes
  - Introduction
  - Model
  - Large dimension
  - Maximum Likelihood Estimation
  - Mean-Field approximation

- Optimization is the main problem in practice
- Need efficient optimization for the training step
- Recent advances on this topic lately, aimed at machine learning applications

Let's put for short

$$f(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i^\top \theta) \quad \text{and} \quad g(\theta) = \lambda \|\theta\|_1$$

How to minimize  $f + g$  ?

A key point: the **descent lemma**

- If  $f$  convex and  $\nabla f$  is  $L$ -Lipschitz, then for any  $\theta, \theta^k \in \mathbb{R}^d$

$$f(\theta) + g(\theta) \leq f(\theta^k) + \nabla f(\theta^k)^\top (\theta - \theta^k) + \frac{L}{2} \|\theta - \theta^k\|_2^2 + g(\theta)$$

and remark that

$$\begin{aligned}
 & \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left\{ f(\theta^k) + \nabla f(\theta^k)^\top (\theta - \theta^k) + \frac{L}{2} \|\theta - \theta^k\|_2^2 + g(\theta) \right\} \\
 &= \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left\{ \frac{L}{2} \left\| \theta - \left( \theta^k - \frac{1}{L} \nabla f(\theta^k) \right) \right\|_2^2 + g(\theta) \right\} \\
 &= \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{2} \left\| \theta - \left( \theta^k - \frac{1}{L} \nabla f(\theta^k) \right) \right\|_2^2 + \frac{1}{L} g(\theta) \right\} \\
 &= T_{g/L} \left( \theta^k - \frac{1}{L} \nabla f(\theta^k) \right)
 \end{aligned}$$

where  $T_{g/L}(\theta) = \operatorname{sign}(\theta) \odot (|\theta| - \lambda/L)_+$

## Proximal gradient descent algorithm

- **Input:** starting point  $\theta^0$ , Lipschitz constant  $L > 0$  for  $\nabla f$
- For  $k = 1, 2, \dots$  until *converged* do
  - $\theta^k = T_{g/L}\left(\theta^{k-1} - \frac{1}{L}\nabla f(\theta^{k-1})\right)$
- **Return** last  $\theta^k$

Simplest algorithm, a plethora of others:

- Coordinate descent
- L-BFGS-B
- Primal-Dual
- ...

What if  $n$  (and  $d$ ) is large?

- Each iteration of a full gradient method has complexity  $O(nd)$
- A large dataset makes a modern computer look old: go back to “old” algorithms (Robbins and Monro 1951)

Idea: we want to minimize an average of losses...

- If I choose uniformly at random  $I \in \{1, \dots, n\}$ , then

$$\mathbb{E}[\nabla f_I(\theta)] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\theta) = \nabla f(\theta)$$

- $\nabla f_I(\theta)$  is an *unbiased* but very noisy estimate of the full gradient  $\nabla f(\theta)$

Robbins and Monro 1951

## Stochastic Gradient Descent (SGD)

- **Input:** starting point  $\theta^0$ , sequence of learning rates  $\{\eta_t\}_{t \geq 0}$
- For  $t = 1, 2, \dots$  until *convergence* do
  - Pick at random (uniformly)  $i$  in  $\{1, \dots, n\}$
  - Put

$$\theta^t = \theta^{t-1} - \eta_t \nabla f_i(\theta^{t-1})$$

- **Return** last  $\theta^t$
- Each iteration has complexity  $O(d)$  instead of  $O(nd)$  for full gradient methods
- The step size must be decreasing
- Fast in the early iterations
- Very slow convergence to a precise minimizer

Recent results improve this:

- Bottou and LeCun (2005)
- Shalev-Shwartz et al (2007, 2009)
- Nesterov et al. (2008, 2009)
- Bach et al. (2011, 2012, 2014, 2015)
- T. Zhang et al. (2014, 2015)

- Put  $X = \nabla f_l(\theta)$  with  $l$  uniformly chosen at random in  $\{1, \dots, n\}$
- Variance of  $X$  as an approximation of  $\mathbb{E}X$  is large
- Reduce it by finding  $C$  s.t.  $\mathbb{E}C$  is easy to compute and such that  $C$  is highly correlated with  $X$
- Put  $Z = X - C + \mathbb{E}C$ , so that  $\mathbb{E}Z = \mathbb{E}X$  and

$$\text{Var } Z = \text{Var } X + \text{Var } C - 2 \text{cov}(X, C)$$

## Stochastic Variance Reduced Gradient (SVRG)

**Input:** starting point  $\theta^0$ , learning rate  $\eta > 0$

Put  $\tilde{\theta} \leftarrow \theta^0$

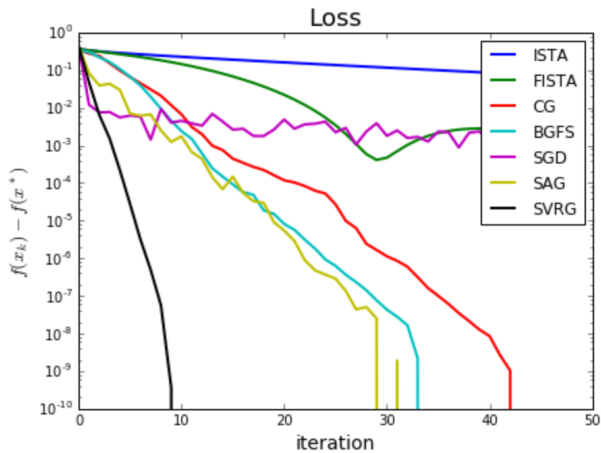
For  $k = 1, 2, \dots$  until *convergence* do

- Put  $\theta_1^k \leftarrow \tilde{\theta}$
- Compute  $\nabla f(\tilde{\theta})$
- For  $t = 1, \dots, n$ 
  - Pick uniformly at random  $i$  in  $\{1, \dots, n\}$
  - Apply the step

$$\theta_{t+1}^k \leftarrow \theta_t^k - \eta(\nabla f_i(\theta_t^k) - \nabla f_i(\tilde{\theta}) + \nabla f(\tilde{\theta}))$$

- Set  $\tilde{\theta} \leftarrow \theta_{n+1}^k$

# Numerical comparison: batch VS stochastic



What we've seen in this glimpse of supervised learning:

- Model assumptions, simplifying assumptions
- Overfitting, penalization, sparsity, cross-validation
- Optimization matters!

Beyond supervised learning?

Let us describe an example: Hawkes processes

- 1 Teasers
  - Data Science in the media
  - Examples
  - Big Data is (quite) Easy ?

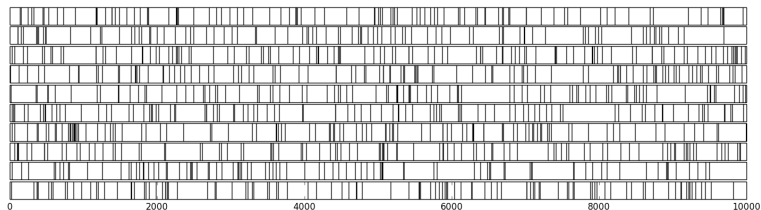
- 2 Supervised learning
  - Introduction
  - An example: RTB
  - Loss functions, linearity
  - Logistic regression
  - Penalization

- 3 Optimization
  - Introduction

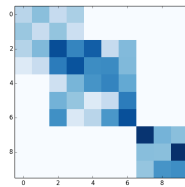
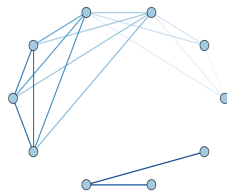
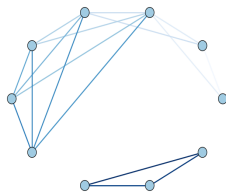
- Gradient descent
- Stochastic gradient descent
- Variance reduction
- Conclusion

- 4 Hawkes processes
  - Introduction
  - Model
  - Large dimension
  - Maximum Likelihood Estimation
  - Mean-Field approximation

## Observation



## What we want to infer



One Hawkes component can model (on a high frequency, data can be as precise as 10 microseconds)

- an arrival of a particular type of order (with a particular size, any agent)
- a move of the price of a particular size and direction and of a particular asset
- an order of a particular agent

NB:

- possible to combine all that and do some high-dimensional models
- possible to add exogenous events with some impacts (e.g., news)

- A two-components model of the high-frequency (mid)price dynamics (one component for upward jumps and one for downward jumps). Basket of assets=high dimension
- Multiple agents models: an event is a particular type of order of a particular size by a particular agent, including components for the moves of the price. Analysis of the interactions of all the agents and their impact on the price

## Model: Multivariate Hawkes Process (MHP)

- A  $d$ -dimensional counting process  $N = [N_1, \dots, N_d]^\top$
- $d$  is “large”
- Observed on  $[0, T]$
- $N_j$  has intensity  $\lambda_j$ , namely

$$\mathbb{P}(N_j \text{ has a jump in } [t, t + dt] \mid \mathcal{F}_t) = \lambda_j(t)dt$$

for  $j = 1, \dots, d$  where  $\mathcal{F}_t$  some filtration

## Model: Multivariate Hawkes Process (MHP)

- MHP assumes the following autoregressive structure:

$$\lambda_j(t) = \mu_j(t) + \int_{(0,t)} \sum_{k=1}^d \varphi_{j,k}(t-s) dN_k(s),$$

- $\mu_j(t) \geq 0$  baseline intensity of the  $j$ -th coordinate
- $\varphi_j : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  self-exciting component
- Write this in matrix form

$$\lambda(t) = \boldsymbol{\mu} + \int_{(0,t)} \boldsymbol{\varphi}(t-s) dN(s),$$

with  $\boldsymbol{\mu} = [\mu_1, \dots, \mu_d]^\top$  and  $\boldsymbol{\varphi}(t) = [\varphi_{j,k}(t)]_{1 \leq j, k \leq d}$ .

- Notation:

$$\int_{(0,t)} \varphi(t-s) dN_k(s) = \sum_{i: 0 < T_{i,k} < t} \varphi(t - T_{i,k})$$

Introduced by Hawkes in 1971

- **Earthquakes and geophysics** : Kagan and Knopoff (1981), Zhuang, Harte, Werner, Hainzl and Zhou (2012)
- **Genomics** : Reynaud-Bouret and Schbath (2010)
- **High-frequency Finance** : Bacry Delattre Hoffmann and Muzy (2013)
- **Terrorist activity** : Porter and White (2012)
- **Neurobiology** : Hansen, Reynaud-Bouret and Rivoirard (2012)
- **Social networks** : Carne and Sornette (2008), Simma and Jordan (2010), Zhou Song and Zha (2013)
- And even **FPGA-based implementation** : Guo and Luk (2013)

### **Parametric estimation** (Maximum likelihood)

- First work : Ogata 78
- Simma and Jordan (2010), Zhou Song and Zha (2013)  
→ Expected Maximization (EM) algorithms, with priors

### **Non parametric estimation**

- Marsan Lengliné (2008), generalized by Lewis, Mohler (2010)  
→ EM for penalized likelihood function  
→ Monovariate Hawkes processes, Small amount of data, No theoretical results
- Reynaud-Bouret and Schbath (2010)  
→ Developed for small amount of data (Sparse penalization)
- Bacry and Muzy (2014)  
→ Larger amount of data

Dimension  $d$  is large:

- Need a simple parametric model on  $\mu$  and  $\varphi$
- For inference: we want a **tractable** and **scalable** optimization problem

A recent work [Bacry, G., Mastromatteo, Muzy 2016]:

- Focus on optimization (the training task) for parametric Hawkes models
- Exploits a recent mean-field property of this model
- Improves state-of-the-art solvers in this case

# A simple parametrization of the MHP

Simple parametrization:

- Constant baselines  $\mu_j(\cdot) \equiv \mu_j$
- Take

$$\varphi_{j,k}(t) = a_{j,k} \alpha e^{-\alpha t}$$

- $a_{j,k}$  = level of interaction between agents  $j$  and  $k$
- $\alpha$  = lifetime of instantaneous excitation (assumed known here...)

The matrix

$$\mathbf{A} = [a_{j,k}]_{1 \leq j, k \leq d}$$

is understood has a **weighted adjacency matrix** of mutual excitement of agents  $\{1, \dots, d\}$

NB: we can consider actually

$$\varphi_{j,k}(t) = \sum_{l=1}^L a_{j,k,l} \alpha_l e^{-\alpha_l t}$$

We end up with intensities

$$\lambda_{j,\theta}(t) = \mu_j + \int_{(0,t)} \sum_{k=1}^d a_{j,k} \alpha e^{-\alpha(t-s)} dN_k(s)$$

for  $j \in \{1, \dots, d\}$  where

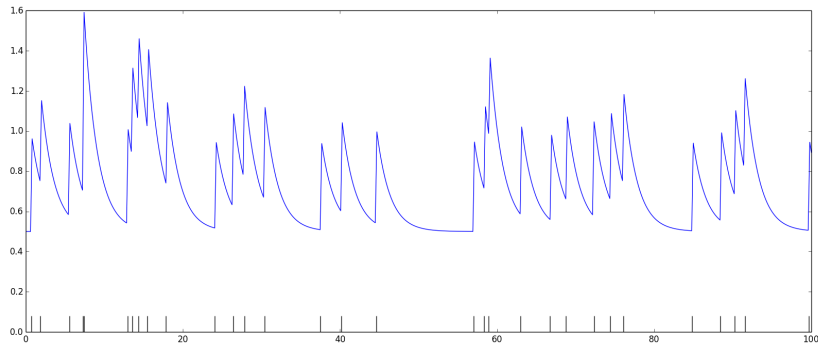
$$\theta = [\mu, \mathbf{A}]$$

with

- baselines  $\mu = [\mu_1, \dots, \mu_d]^\top \in \mathbb{R}_+^d$
- interactions  $\mathbf{A} = [a_{j,k}]_{1 \leq j, k \leq d} \in \mathbb{R}_+^{d \times d}$

# A simple parametrization of the MHP

For  $d = 1$ , intensity  $\lambda_\theta$  looks like this:

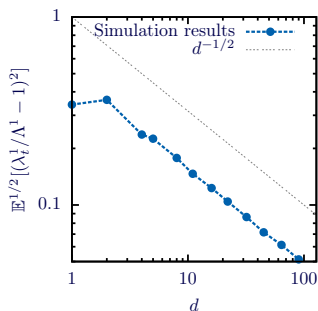
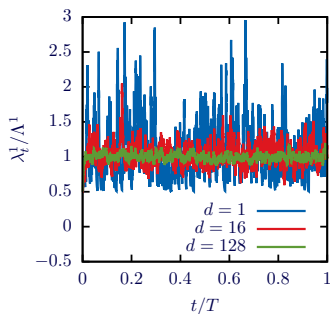


$$-\ell_T(\theta) = \sum_{j=1}^d \left\{ \int_0^T (\lambda_{j,\theta}(t) - 1) dt - \int_0^T \log \lambda_{j,\theta}(t) dN_j(t) \right\}$$

with

$$\lambda_{j,\theta}(t) = \mu_j + \sum_{k=1}^d a_{j,k} \int_{(0,t)} \alpha \exp(-\alpha(t-s)) dN_k(s)$$

- For inference, we exploit the fact that  $d$  is large
- Using a Mean-Field approximation! [Delattre, Fournier and Hoffmann 2015]



Idea: when  $d$  is large, we have dilute interactions, ie

$$\lambda_{j,\theta}(t) \approx \frac{N_j([0, T])}{T} = \bar{\Lambda}_j$$

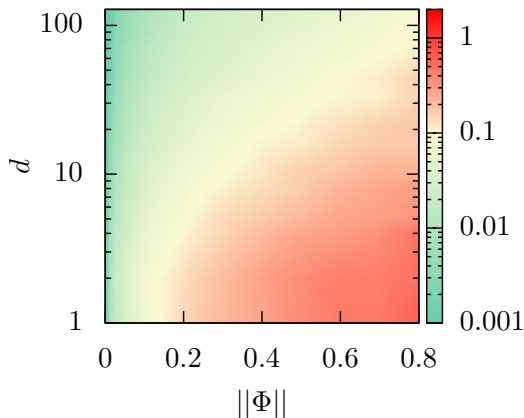
So, we use the following approximation for inference

$$\log \lambda_{j,\theta}(t) \approx \log \bar{\Lambda}_j + \frac{\lambda_{j,\theta}(t) - \bar{\Lambda}_j}{\bar{\Lambda}_j} - \frac{(\lambda_{j,\theta}(t) - \bar{\Lambda}_j)^2}{2(\bar{\Lambda}_j)^2}$$

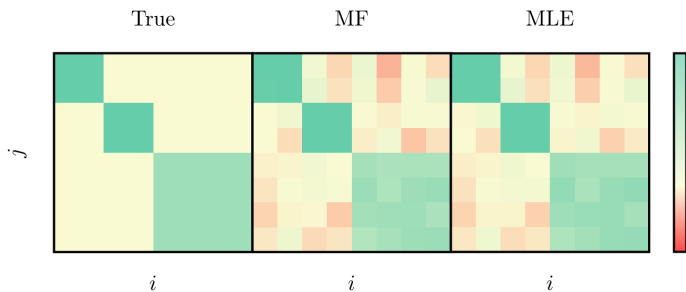
It gives an approximation of the log-likelihood

- Faster to solve
- Surprisingly sharp in large dimension

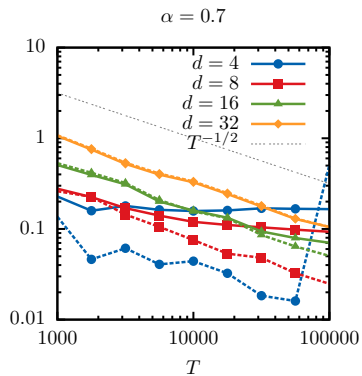
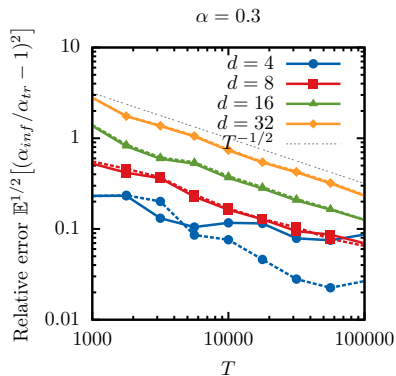
Fluctuations  $\mathbb{E}^{1/2}[(\lambda_t^1/\Lambda^1 - 1)^2]$



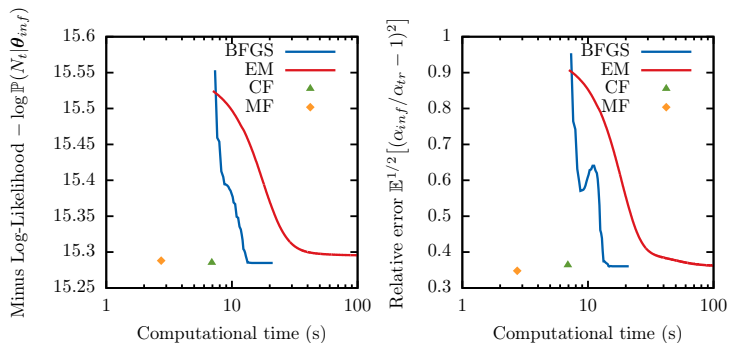
# Mean-field inference for Hawkes



# Mean-field inference for Hawkes



Training algorithm faster by several order of magnitude than state-of-the-art solvers



Thank you!