# Research Summer School on Statistics for Data Science
# S4D 2018, Caen, France
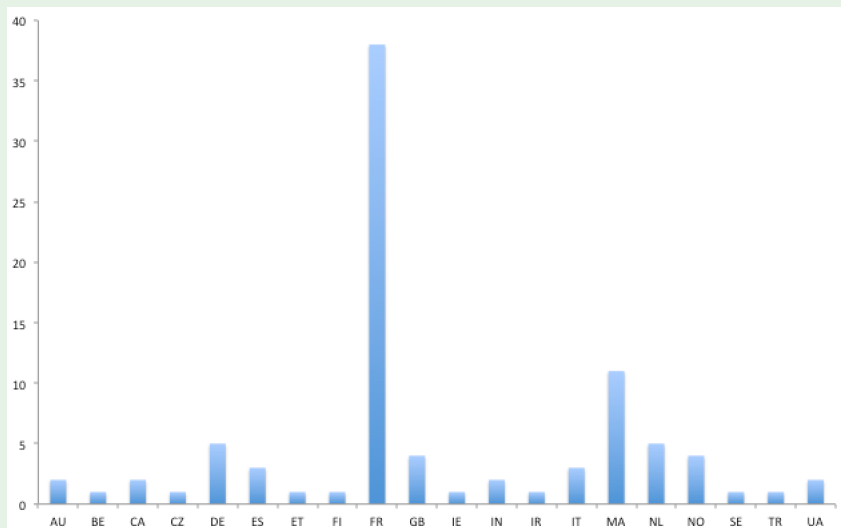
FAICEL CHAMROUKHI
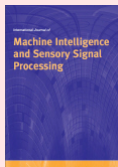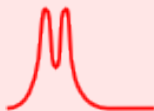
https://chamroukhi.com

## S4D key insights

- Talks covering both tutorial and advanced aspects at the interface of Statistics, Machine Learning and Optimization

  $\hookrightarrow$ the main data science fields

- Theoretical foundations and algorithmic aspects, as well as typical case studies in complex and large-scale scenarios

- We enjoy good food and nice visits to Normandie sites :)

# 86 participants from 20 countries!

Many thanks to our sponsors!

- The term "Data Science" has surged in popularity
- Data science is increasingly commonly used with "big data."
- Data science, including Big Data has recently attracted an enormous interest from the scientific community

- What does Data Science mean?
- What about Statistics in the Data Science "area" ?
- There is not yet a consensus on what precisely constitutes Data Science



- For a review, see the report of D. Donoho (2015): "50 years of Data Science"

- There is not yet a consensus on what precisely constitutes Data Science, but
- Data Science can be seen (defined ?) as[a]:
  - the study of the generalizable extraction of knowledge from data.
  - requires an integrated skill set spanning mathematics, machine learning, artificial intelligence, statistics, databases, and optimization

---

[a]Vasant Dhar (2013): Communications of the ACM, Vol. 56 No. 12: 64-73

- Data Science clearly has an interdisciplinary nature and requires substantial collaborative effort
- Databases, statistics and machine learning, and distributed systems are emerging as foundational to data science

(i) Databases: organization of data resources,

(ii) Statistics and Machine Learning: convert data into knowledge,

(iii) Distributed and Parallel Systems: computational infrastructure

Statistics play a central role in data science

- Allow to quantify the randomness component in the data

- A well-established background to deal with uncertainty (probabilistic framework) and to establish generizable methods for prediction and estimation

- allow soft decision: e.g. confidence interval in regression and posterior probabilities in classification

- help for understanding the underlying generative process

# Data science models/algorithms

New problems (big data, etc) but ... classical methods ?



## Our Core Algorithms Remain the Same

• Regression, decision trees, and cluster analysis continue to form a triad of core algorithms for most data miners. This has been consistent since the first Data Miner Survey in 2007.

Question: What algorithms / analytic methods do you TYPICALLY use? (Select all that apply)

2016 Rexer Analytics

# S4D 2018: programme

| | Mon. 18 | Tue. 19 | Wed. 20 | Thu. 21 | Fri. 22 |
|---|---|---|---|---|---|
| 08:00 | | | | | |
| 09:00 | | Probabilistic modeling for machine learning (I) | Model selection theory and considerations in large-scale scenarios (II) | Feature selection in high-dimensional problems (II) | Unsupervised learning from high-dimensional and functional data |
| 10:00 | | Coffee break | Coffee break | Coffee break | Coffee break |
| 11:00 | Welcome reception | Probabilistic modeling for machine learning (II) | Mixture models and feature selection in high-dimensional problems (I) | Theory of statistical Inference (I) | Mixture of experts for regression, clustering and classification in high-dimensional scenarios |
| 12:00 | | | | | |
| 13:00 | Lunch | Lunch | Lunch | Lunch | Lunch |
| 14:00 | Optimization for ML (I) | Model-based clustering and co-clustering in high-dimensional scenarios (I) | Le Mont Saint Michel | Majorization-Minimization (MM) Algorithms for Statistical Inference and Machine Learning Problems (II) | Conference closure |
| 15:00 | | Oral presentation | | Coffee break | |
| 16:00 | Coffee break | Coffee break | | Oral presentation | |
| | Mixed Integer Optimization for unsupervised learning. Applications to clustering and Image segmentation (II) | WW2 landing beaches | | | |
| 17:00 | | | | | |
| 18:00 | Visit of Caen | | | Posters session | |
| 19:00 | | | | | |
| 20:00 | | | | | |
| 21:00 | Dinner (at La Planche's 13 Rue Prairies Saint-Gilles, 14000 Caen) | Dinner (at Le Carlotta, 16 Quai Vendeuvre, 14000 Caen) | Dinner (at Le Dauphin, 29, Rue Gemare, 14000 Caen) | Dinner (at Le Carlotta, 16 Quai Vendeuvre, 14000 Caen) | |
| 22:00 | | | | | |

Many thanks for your participation!

Enjoy the courses!

Looking forward to seeing you next year :)