

14D008

3 ECTS

Topics in Big Data Analytics I

Overview and Objectives

Constant advances in digital sensors, Internet, mobility and storage, result in the explosion of available data that potentially carries significant value to business, science and society. This poses many challenges both technological and analytical. A wide variety of techniques have arisen with the objective of discovering hidden patterns in data. These methods are fully exploited by top technology companies such as Amazon, Netflix, Twitter or Google and define the core of their competitive advantage. This course is structured as a series of four lectures where students will be presented with the theoretical underpinning and practical implementation of some of the most groundbreaking big data applications.

Lecturers

The course is coordinated by Pau Agulló (Kernel Analytics) and classes will be delivered by Mikel Arizaleta, Emiliano Carluccio, Marçal Molins, Diego Gruber, Miquel Camprodon (Kernel Analytics)

Course Outline

- 1. Discrete optimization in Python:** We will work on a complete lifecycle of a real world optimization problem: from the theory (problem statement and mathematical encoding) to the practise (system configuration and solution programming). We will focus on MILP (Mixed Integer Linear Programming) problems, reviewing the state of the art of current solvers, and using Python libraries, which provide us interfaces to solvers making coding and integration much easier.
- 2. Deploying ML Models as APIs:** Web services have become ubiquitous in data science, both as a way to access data as well as an interface for models that simplify integration within an organization's systems. We will learn the building blocks of REST APIs, how to use them to retrieve data as well as how to deploy your own model as an API using R and OpenCPU.
- 3. Personalized recommendation on distributed architectures.** An update of last year module making it really "big data". How to prepare data and execute collaborative filter algorithms on big data architectures (Hadoop, Spark)
- 4. How to work with spatial data:** Maps, gps localizations, addresses, the spatial data of an entity give us a lot of new hidden information about the entity itself and about its neighbours. In this 4 hours class we will show how to get, clean and cross this data and how to extract the relevant information out of it.
- 5. Docker, unit tests and dev-ops for analytics:** In this 4 hours class we will show how to leverage R, Docker, GIT and other Continuous integrations tools to build from scratch a basic ML service prototype ready for deployment in a (pre)-production framework.

Required Activities

Attendance at classes, a practical project (consisting of programming exercises)

14D008

Topics in Big Data Analytics I

Evaluation

- 100% Project. Students will hand out a programming project for each topic at the end of the course.

Materials

Books:

Anand Rajaraman and Jeffrey David Ullman.
2011. Mining of Massive Datasets. Cambridge
University Press, New York, NY, USA.
Chapters 3, 8 and 9.

Christopher D. Manning, Prabhakar Raghavan,
Hinrich Schütze, Introduction to Information
Retrieval, Cambridge University Press.

Papadimitriou, C. H.; Steiglitz, K. (1998).
Combinatorial optimization: algorithms and
complexity

Walker, R. C. (1991). Introduction to Mathematical
Programming

Agile Data Science: Building Data Analytics
Applications with Hadoop, Russell Journey, O'Reilly

Others:

Bennett, James, and Stan Lanning. "The netflix
prize." Proceedings of KDD cup and workshop.
Vol. 2007. 2007.

Bell, Robert M., and Yehuda Koren, "Lessons
from the Netflix prize challenge." ACM SIGKDD
Explorations Newsletter 9.2 (2007): 75-79.

Coordinate Sysems: Coordinate Systems and Map

Projections, D.H. Maling
Carto: <https://carto.com/docs/>

[MapReduce: Simplified Data Processing on Large Clusters](#)

[An Introduction to APIs](#) by Brian Cooksey (Zapier)

[Out of the Weeds and into Product: APIs and the Future of Data Science](#) by Ken Tien (Mulesoft)

[R Packages](#) by Hadley Wickham (RStudio)

[The OpenCPU System: Towards a Universal Interface for Scientific Computing through Separation of Concerns](#) by Jeroen Ooms (OpenCPU)

CoinOr: Computational Infrastructure for Operations
Research. <https://www.coin-or.org/>
Python PuLP
library. <https://pythonhosted.org/PuLP/>
An Introduction to Linear Programming and the
SimplexAlgorithm. <https://www2.isye.gatech.edu/~s-pyros/LP/LP.html>