

Mathematics and Big Data

Code: 43478
ECTS Credits: 6

Degree	Type	Year	Semester
4313136 Modelling for Science and Engineering	OT	0	2

Contact

Name: Alejandra Cabaña Nigro
Email: AnaAlejandra.Cabana@uab.cat

Use of languages

Principal working language: english (eng)

Teachers

Albert Ruíz Cirera

External teachers

Isabel Serra

Prerequisites

Students should have basic knowledge linear algebra, statistical inference and linear models. We also assume the students have programming skills.

Previous experience with R and/or Python will be helpful.

Objectives and Contextualisation

The aim of this course is to learn and apply various mathematical and statistical methods related to the discovery of relevant patterns in data sets. Nowadays, huge amounts of data are being generated in many fields, and the goal of this course is to learn how to extract information from such data.

Skills

- Analyse, synthesise, organise and plan projects in the field of study.
- Apply logical/mathematical thinking: the analytic process that involves moving from general principles to particular cases, and the synthetic process that derives a general rule from different examples.
- Conceive and design efficient solutions, applying computational techniques in order to solve mathematical models of complex systems.
- Formulate, analyse and validate mathematical models of practical problems in different fields.
- Isolate the main difficulty in a complex problem from other, less important issues.
- Solve complex problems by applying the knowledge acquired to areas that are different to the original ones.

Learning outcomes

1. Analyse, synthesise, organise and plan projects in the field of study.
2. Apply logical/mathematical thinking: the analytic process that involves moving from general principles to particular cases, and the synthetic process that derives a general rule from different examples.
3. Identify real phenomena as models of stochastic processes and extract new information from this to interpret reality.
4. Isolate the main difficulty in a complex problem from other, less important issues.
5. Solve complex problems by applying the knowledge acquired to areas that are different to the original ones.
6. Use appropriate statistical packages and Bayesian methods solutions to solve specific problems.

Content

Statistics:

The problem of multiple testing. Linear and Generalized linear methods: LASSO, Elastic Nets. Functional Data Analysis : Observed functional data and its computational representation, descriptive statistics and dimensionality reduction, depth measures for FD, two-sample problem for FD, Functional linear models, classification techniques.

Topological data analysis:

Topology and data, quick review of linear algebra, from points to polyhedra, combinatorial topology, persistence Diagrams and software.

Statistical Learning:

Review of the state-of-the-art in statistical learning techniques.

Methodology

Lectures, supervised exercises and autonomous activities directed to perform data analysis projects based on statistical and topological tools.

Activities

Title	Hours	ECTS	Learning outcomes
Type: Directed			
Lectures	38	1.52	1, 3
Type: Supervised			
Completion of exercises	36	1.44	2, 6
Type: Autonomous			
Personal study, readings	20	0.8	3
Project	44	1.76	1, 2, 3, 4, 5, 6

Evaluation

Practical Exercises: Completion and presentation of the proposed exercises.

Final Projects: each part of the course will have its own projects, which include a written report, and maybe a public presentation.

Due dates will be announced during the course and will be strict.

Evaluation activities

Title	Weighting	Hours	ECTS	Learning outcomes
Final Projects	0.8	6	0.24	1, 2, 3, 4, 5, 6
Practical Exercises	0,2	6	0.24	2, 3, 4, 5, 6

Bibliography

Basic references

B. Efron, T. Hastie, *Computer Age Statistical Inference*, Cambridge University Press (2016) (5th Ed 2017)

G. James, D. Witten, T. Hastie and R. Tibshirani, *An Introduction to Statistical Learning (with applications in R)*. Springer, 2013.

Gunnar Carlsson, "Topology and data". Bull. AMS 46,2 (2009), 255-308.

P. Kokoszka, M. Reimherr, *Introduction to Functional Data Analysis*. CRC Press.(2017).

Ramsay, J. , B. W. Silverman, *Functional Data Analysis Springer* (2nd Ed. 2005).

Complementary references

B. Everitt and T. Hothorn, "An introduction to Applied Multivariate Analysis with R". Springer, 2011.

(B. Everitt, "An R and S+ Companion to Multivariate Analysis", Springer, 2005).

J Faraway, " Extending de Linear Model with R", Chapman & Hall, Miami, 2006.

J Faraway, "Linear Models with R", Chapman & Hall, Boca Raton, 2005.

W. Härdle and L. Simar, "Applied Multivariate Statistical Analysis". Springer. 2007.

B. Ripley, "Pattern Recognition and Neural Networks". Cambridge University Press, 2002.

L. Torgo. "Data Mining with R. Learning with Case Studies". Chapman & Hall, Miami. 2010

W Venables, B Ripley, "Modern Applied Statisticswith S-PLUS", Springer, New York.

Collins FS and Varmus H, "A new initiative on precision medicine". N Engl J Med. 2015 Feb 26;372(9):793-5 .

Jensen A.B. et al, "Temporal disease trajectories condensed from population-wide registry data covering 6.2 million patients". Nat Commun 2014 Jun 24; 5:4022.

J.D. Jobson, "Applied Multivariate Analysis". Vol I i II. Springer, 1992.

R. Johnson and D.W. Wichern, "Applied Multivariate Statistical Analysis". Pearson Education International, 2007.

P.Y.Lum et al., "Extracting insights from the shape of complex data using topology". Sci. Rep. 3, 1236; DOI:10.1038/srep01236 (2013).

A. Rencher, "Methods of Multivariate Analysis". Wiley Series in Probability and Mathematical Statistics, 2002.

D. Skillicorn, "Understanding Complex Data. Data Mining with Matrix Decomposition". Chapman&Hall, 2007.

G. Singh, F. Mémoli, G. Carlsson, "Topological methods for the analysis of High dimensional data sets and 3D object recognition". Eurographic Symp. on Point-Based Graphics, 2007

Journal of Statistical Software, <http://www.jstatsoft.org/>

Dealing with Data (2011) Special Issue. Science 11 February 2011:692-789