**UNIVERSITAT POLITÈCNICA DE CATALUNYA**
**BARCELONATECH**

**Escola Superior d'Enginyeries Industrial,
Aeroespacial i Audiovisual de Terrassa**

**MUEI, MUEA, and MASE**

# Advanced Engineering Data Analysis (AEDA)
# General Information

**Prof. Daniel Fernández**

*daniel.fernandez.martinez@upc.edu*

Daniel Fernández (**coordinator**)

daniel.fernandez.martinez@upc.edu

## My background

- B.S. Statistics . UPC. 1998
- B.S. Computer software. UPC. 2000
- M.S. in statistics and probability. Mathematics Research Centre. 2005
- Ph.D. in Statistics - Victoria University of Wellington (VUW). 2015
- Post Doctoral Fellow in Statistics. VUW. 2015
- Moore-Sloan Post Doctoral Associate in Statistics. NYU. 2016
- Professor, Department of Epidemiology and Biostatistics, SUNY, 2017
- Currently, Serra-Húnter lecturer professor. UPC.

# Class method

- Virtual campus: Atenea, https://atenea.upc.edu/

# *Class method*

- Virtual campus: Atenea, https://atenea.upc.edu/

- Mondays 17-19h. Theory (Room: TR5-0.3)

# *Class method*

- Virtual campus: Atenea, https://atenea.upc.edu/

- Mondays 17-19h. Theory (Room: TR5-0.3)

- Tuesdays 15-17h. Lab/Practice (Room: TR5-0.3)
  **Do you have laptops?**

# *Goal*

The purpose of this course is:

- To introduce the most common and well-known multivariate statistical methods to non-mathematicians.

- It is not intended to be a comprehensive course (mathematically speaking). However, it is important not to take the statistical methods a black box.

- The intention is to keep the details to a minimum while serving as a practical guide that illustrates the possibilities of multivariate statistical analysis.

- In other words, it is a course to "get you going" in a particular area of statistical methods.

# *Preliminary Knowledge*

It is assumed that you (the students) have a working knowledge of:

- Elementary statistics
  - ➢ Basic statistics: summary statistics (mean, median,…), Normal distribution, CI, Hypothesis testing,…
  - ➢ linear regression
  - ➢ EDA (visualization)

- Some facility with algebra is also required to follow the equations in certain parts of the text. Understanding the theory of multivariate methods requires some matrix algebra. However, the amount needed is not great.

# *Preliminary ideas*

- We will learn a comprehensive set of multivariate methods. The important thing you do not have to memorize the techniques and their assumptions. **You need to know what they do conceptually.**

- This course should be helpful to **train your brain** in a way that if you have data and you want to analyze it, you are can go to a book/slides/etc. and find what you want.

- Think about statistics techniques as a **toolbox**, and you use the one you consider more necessary in each moment.

- All models and techniques **are wrong**.

- Nothing can simulate the real data to perfection, but we always look for the best model (i.e., the model with less uncertainty), but **uncertainty is always there**.

- In Physics is possible, because they are laws, However, we **play with uncertainty in Statistics**. The good thing about statistics is that we can measure the error (95% of CI, for example)

- *"I think it is much more interesting to live with uncertainty than to live with answers that might be wrong"* –Richard Feynman (https://en.wikipedia.org/wiki/Richard_Feynman).

# *Tentative schedule*

**Advanced Engineering Data Analysis. Schedule**

| Week | Calendar period | Topic/Activity |
|------|-----------------|----------------|
| 1 | 21-Feb & 22-Feb | Introduction, R and RStudio, Basic Statistics & Exploratory Data Analysis |
| 2 | 28-Feb & 1-Mar | Principal Components Analysis |
| 3 | 7-Mar & 8-Mar | Linear Discriminant Analysis |
| 4 | 14-Mar & 15-Mar | Classification |
| 5 | 21-Mar & 22-Mar | Clustering |
| 6 | 28-Mar & 29-Mar | Project presentation & Quiz |

# *Class description*

- Mondays 17-19h. Theory
  - We will show the methods with examples.

- Tuesdays 15-17h. Lab Practice
  - Goal: practice.
  - We will deliver examples and lab practice (**not evaluable**).
  - Lab practice will be solved **individually** (COVID-19)
  - Questions can be posed during the two hours.
  - Important: we **will not respond** to questions regarding the lab practice after the end of the class.
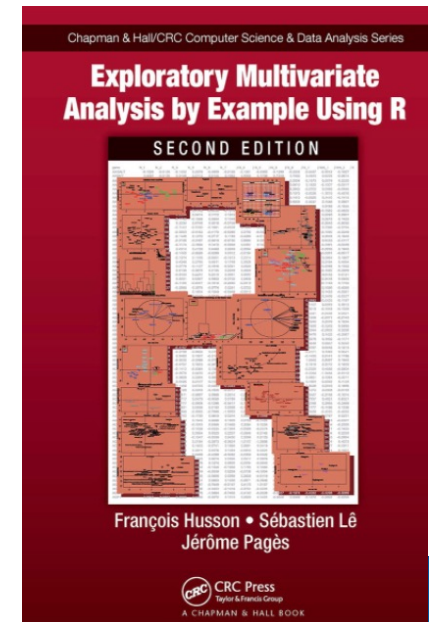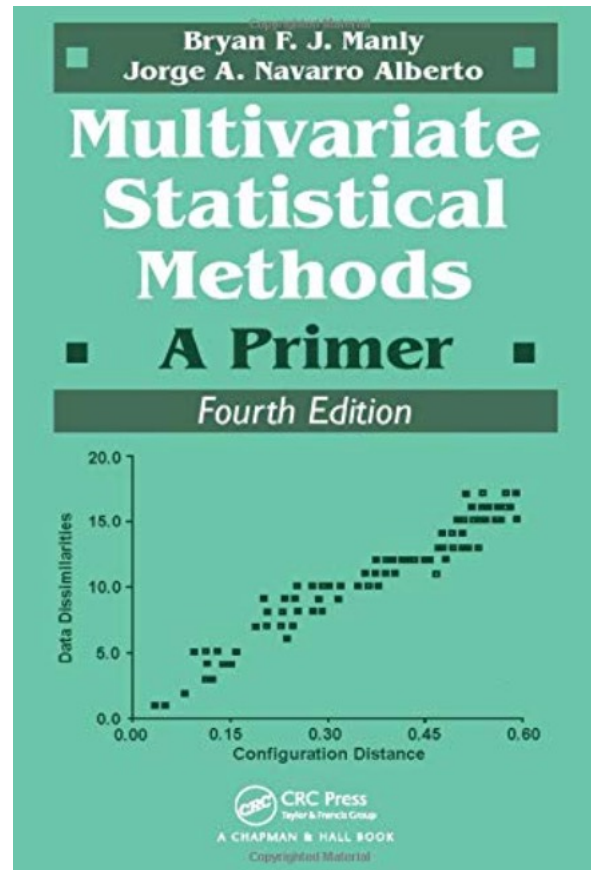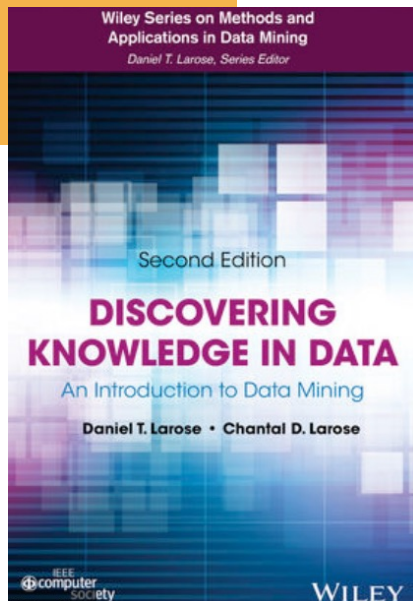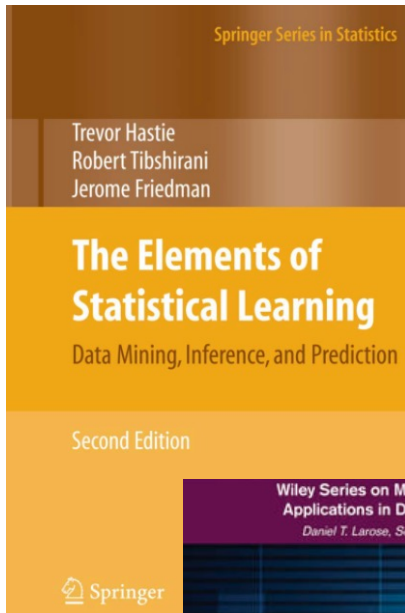
  *Note: Lab practice is not the same as the evaluable exercises you will be doing in groups (3-4 people) throughout the course*

# *Basic Bibliography*

UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONATECH
UPC
Escola Superior d'Enginyeries Industrial,
Aeroespacial i Audiovisual de Terrassa

Academic assessment will be based on the grades obtained in:

     1. Exercises – **30%**

     2. A exam (Quiz) – **30%**

     3. A project –  **40%**

UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONATECH
Escola Superior d'Enginyeries Industrial,
Aeroespacial i Audiovisual de Terrassa

- **Exercises**
  - They will be 2 assignments corresponding to the different units of teaching (throughout the course).
  - They will be done in **3-4 people**
  - Submitted as a report (**PDF format**) to Atenea.
  - The report must have **2 sheets at most (four pages)**,
  - Written in 11-12 font size.

- **Quiz**
  - It will take place the last day of the course -- **Mar 29th**
  - It can have methodological, R and practical questions.

UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONATECH
Escola Superior d'Enginyeries Industrial,
Aeroespacial i Audiovisual de Terrassa

- ## **Exercises**

### Advanced Engineering Data Analysis. Exercises

| Exercise | Topics | Weight | Handout date | Due date |
|---|---|---|---|---|
| 1 | Principal Components Analysis Linear Discriminant Analysis and extensions | 50% | 8/3/22 | 15/3/22 |
| 2 | Classification* | 50% | 15/3/22 | 22/3/22 |

*Clustering will be evaluated with the Final Project

- The project will be done on a **real data set**

- **Oral presentation** on second-to-last day of class (**Mar 28th**)

- Form a **3-4-person group** (same group as exercises. You are currently 12, that means 3 groups of 4 or 4 groups of 3, for instance)

- It must be done using 

- Basic steps of the work to do:
  - Choose a "real-world" domain and define the problem
  - Implement and test methods
  - Write a report
  - Present it orally

    Instructions: **Project_instructions.pdf** document in Atenea

- UC Irvine Machine Learning Repository
  http://archive.ics.uci.edu/ml/index.php

- Kaggle, with thousands of data sets in Business, Computer Science, Earth and Nature, Health, among other fields.

  https://www.kaggle.com/datasets

- IFCS Data Challenging Repository
  https://ifcs.boku.ac.at/repository/challenge2/

- Google public datasets https://cloud.google.com/bigquery/public-data/

- Eurostat database: https://ec.europa.eu/eurostat/data/database

# Software

## Software

- We will be using R and Rstudio.

- We recommend to start visiting the following webs:
  - ➢ The R Project for Statistical Computing: https://www.r-project.org/
  - ➢ The Comprehensive R Archive Network: https://cran.r-project.org/

- Reading the manuals included in the installation:
  https://cran.r-project.org/manuals.html

- Reading other documents:
  https://cran.r-project.org/other-docs.html

- And especially, this short introduction:
  https://cran.r-project.org/doc/contrib/Torfs+Brauer-Short-R-Intro.pdf

UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONATECH
Escola Superior d'Enginyeries Industrial,
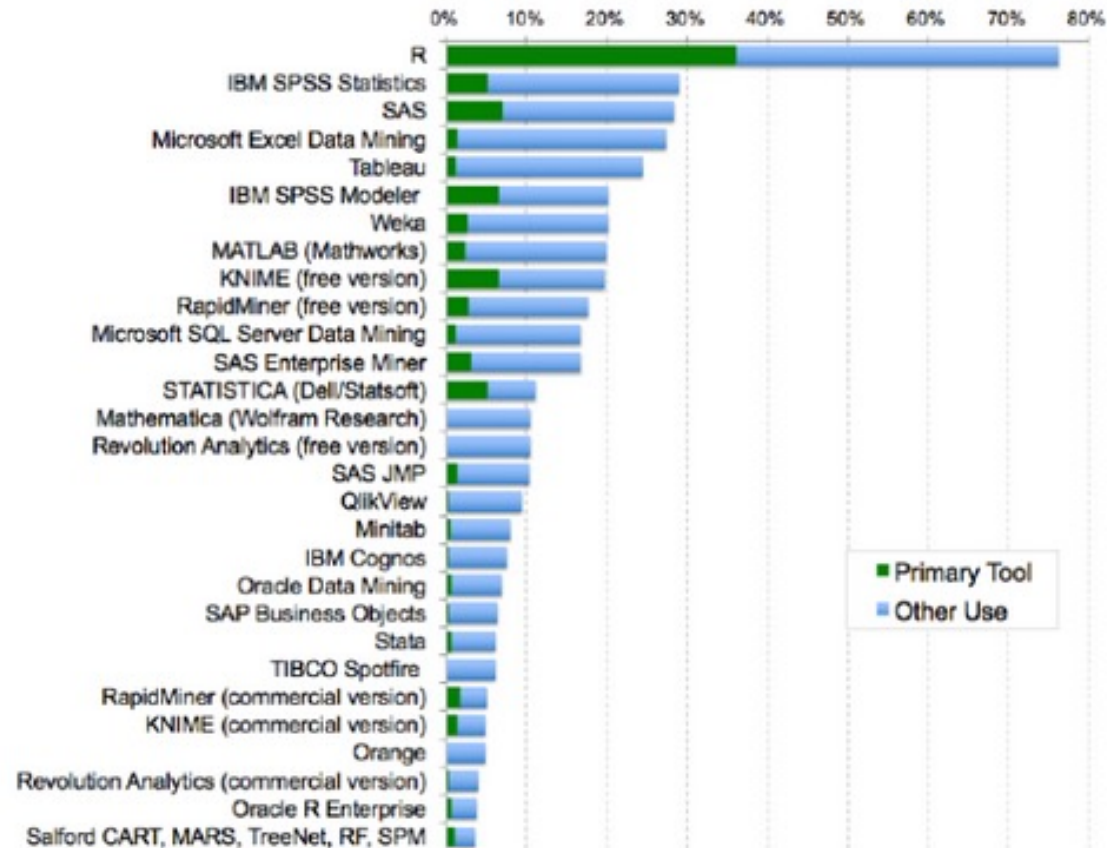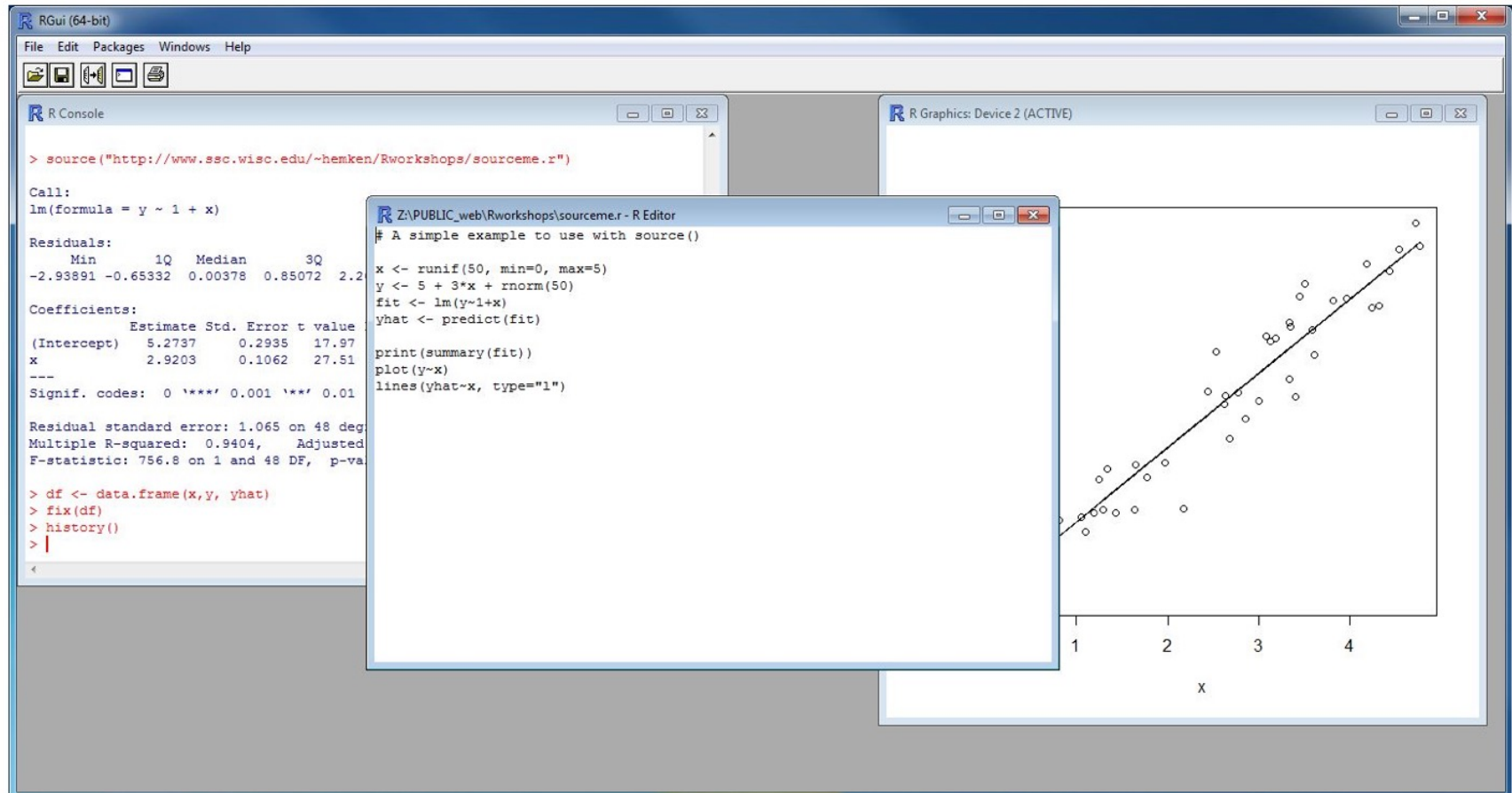Aeroespacial i Audiovisual de Terrassa

- R

It depends on the operating system, but we can consult in:

http://cran.r-project.org/bin

- Rstudio

R-Studio is one of the platforms from which you can run R as well as manage its packages, results, files, etc... https://www.rstudio.com/

It depends on the operating system, but we can consult in:

https://www.rstudio.com/products/rstudio/download/#download

# *Data analytics and R*

UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONATECH
UPC
Escola Superior d'Enginyeries Industrial,
Aeroespacial i Audiovisual de Terrassa

IMPORT

CLEANING & TRANSFORM

VISUALIZE & MODELLING

COMMUNICATE

*AEDA course. ESEIAAT. Session 1. Teaching: Daniel Fernández*

*22*

# *R is increasing a lot*

Figure 6a. Analytics tools used by respondents to the 2015 Rexer Analytics Survey. In this view, each respondent was free to check multiple tools.

R and RStudio are <u>two different programs</u>

- R is the program that calculates. It is **open-source**, **collaborative** and (**initially**) focused on **statistical computing**.
- R is also the language in which we write the commands, as it is also the **programming language**.
- R is w**idely spread** because:
  - o There are many forums on the internet where users raise / resolve their doubts and / or proposals.
  - o there are many libraries for R in continuous development and that allow to use it in a much faster and more efficient way.

# *R. Interface*

R Console interface

R and RStudio are <u>two different programs</u>

- RStudio is the **interface** where we will be working, as it offers us some **comforts**.

- It allows us to create **scripts** in a more agile way and in a much more **pleasant environment** than using R.

- It also makes it **easier** to install and uninstall libraries, load and view databases, etc.

- It offers **other possibilities** beyond R such as creating web pages, pdf's or word files with integrated R, which goes beyond this course.

UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONA**TECH**
Escola Superior d'Enginyeries Industrial,
Aeroespacial i Audiovisual de Terrassa

- Scripts
- Data viewer

- Data type
- Objects

**Programs**

**Data**

**Results**

**Graphs**

- Warnings & Errors

- Submenus

# *RStudio. Interface*

**Toolbar**
Save scripts, settings,…

**Script**
Text file where the analysis instructions are saved

**Environment**
Set of objects in memory
**History**
Set of executed Instructions



**Console**
Commands and results

**Explorer, Graphics Packages, Help, & Viewer**

*AEDA course. ESEIAAT. Session 1. Teaching: Daniel Fernández*

- We recommend to start visiting the following webs:

  o The R Project for Statistical Computing: https://www.r-project.org/

  o The Comprehensive R Archive Network: https://cran.r-project.org/

- Reading the manuals included in the installation:

https://cran.r-project.org/manuals.html

- Reading other documents: https://cran.r-project.org/other-docs.html

- And especially, these tutorials:

  https://cran.r-project.org/doc/contrib/Torfs+Brauer-Short-R-Intro.pdf

  https://www.theanalysisfactor.com/r/