



HOMEWORK 1

The goal is to get a good insight into a dataset by mean of summary statistics and visualisations. For this exercise set choose one alternative below. Regardless of your choice, your submission must comply with the guidelines at the end of this document.

DATA SELECTION

You are give the possibility to choose one set of data attached to the HW assignment:

- ALTERNATIVE 1 - ABALONE: The dataset contains the characteristics of abalones. The data can be either i) retrived from [UC Irvine Machine Learning Repository](#), or ii) within R using the commands: `library(AppliedPredictiveModeling); data(abalone)`.
- ALTERNATIVE 2 - GAPMINDER: The dataset contains an excerpt of the Gapminder data on life expectancy, GDP per capita, and population by country. The data can be either i) retrieved from the [Gapminder](#) website or ii) retrieved within R using the commands: `library(gapminder); data(gapminder)`.
- ALTERNATIVE 3 - WINE QUALITY: The dataset contains two datasets related to red and white wine samples from Portugal. The data can be retrieved from [UCI Machine learning repository](#).
- ALTERNATIVE 4 - GAS SENSOR ARRAY DRIFT: The dataset contains measurements from chemical sensors exposed to different gases at various concentration levels. The data can be retrieved from [UCI Machine learning repository](#).
- ALTERNATIVE 5 - ENVIRONMENT DATA IN MELBOURNE: Energy consumption, climate, and wastewater characteristics of Melbourne wastewater treatment plant for period of six years (2014-2019). The data can be retrieved from [Mendeley data](#).
- ALTERNATIVE 6 - CONCRETE MIXTURE DATA: Data on experiments designed to find concrete formulations that maximize compressive strength. The data used here consists of separate experiments from 17 sources with common experimental factors were combined into one “meta-experiment”. The data can be retrieved from [UCI Machine learning repository](#)
- ALTERNATIVE 7 - YOUR CHOICE: You have a set of data of your own interest. The dataset should comprise of a certain number of observations, each observation consists of a certain number of predictors and corresponding class label.

DATA ANALYSIS

Regardless of your choice, you must:

- 1 Describe your data and their features in terms of number of observations N , number of predictor variables D , number of classes L and class-distribution (that is, the number of observations for each of the classes).
- 2 Perform an unconditional mono-variate analysis of each of the D predictors. Specifically, you must plot their (unconditional) histograms and box-plots, calculate their (unconditional) mean μ_d , standard deviation σ_d and skewness γ_d , with $d = 1, \dots, D$, using all the N observations.
- 3 Perform a class-conditional mono-variate analysis of each of the predictors. Again, you must plot their (class-conditional) histograms and box-plots, calculate their (class-conditional) mean $\mu_{d|l}$, standard deviation $\sigma_{d|l}$ and skewness $\gamma_{d|l}$, with $d = 1, \dots, D$, now using only the N_l observations of class l , for each the L classes.

Item 1 leads to D histograms, D means, D standard deviations and D skewness values. Item 2 leads to $D \times L$ histograms, $D \times L$ means, $D \times L$ standard deviations and $D \times L$ skewness values. Tabulate all means, standard deviations and values of skewness, for both items. Comment on the results, highlight any remarkable fact that emerge from this exploratory analysis. Are there predictors that seem to show any discriminative power (as in, ‘are they, alone, capable to separate the classes’)?

Then, you must

- 4 Perform an unconditional bi-variate analysis of the predictors. Specifically, you must plot the scatter plots between all pairs of predictors. For each point (observation), use colours or symbols to indicate the associated class label. Investigate the existence of potential relationships between pairs of predictors and the presence of potential outliers.

Are there any relevant relationships between pairs of predictors? If yes, are these relationships linear? Quantify linear dependence between predictors using pair-wise correlation coefficients ρ_{d_i, d_j} , with $d_i, d_j = 1, \dots, D$. Either tabulate the correlation coefficients as a correlation matrix $\boldsymbol{\rho}$ with $\boldsymbol{\rho}(i, j) = \rho_{d_i, d_j}$, or show the matrix as an image. Comment on the results.

As final task, you must

- 5 Perform an unconditional multi-variate analysis of the predictors. Specifically, you must perform a principal components analysis of the predictors, for the sake of visualisation, retain only the first two principal components (those associated with the two largest eigenvalues) and plot the scatter plot of the projected observations. Again, for each projected point (observation) you must use colours or symbols to indicate the associated class label. [Remember to perform the necessary pre-processing of the data]

Are the classes well (or better) separated? Are the boundaries between classes linear? What classes show a high degree of overlap and thus are harder to separate?

GUIDELINES

Regardless of your choice of data, you must generate the following:

- Article: You must generate a report in the format of a conference paper following the template adapted from the IEEE conference proceedings¹ that is attached to the homework assignment. The paper should not be longer than 6 pages and must include the following:
 - Title: Here, you summarize your paper in one sentence. [Spend time on it and try some alternatives². As part of the preparation, this will help both you to write a clear abstract and the reader to grasp the content of the work.]
 - Abstract: Here, you introduce the main objective and overview of the work [Provide a short and informative view of the work, its scope and results].
 - Introduction: Here, you provide some context and background [Briefly, explore the literature in order to understand the data, define how and why data need to be pre-processed. Discuss application examples and provide the references.]
 - Methods: Here, you briefly describe your data set and the methods used for analysing it [Report and comment the main characteristics of the data. Plot the most representative histograms for the unconditional and class-conditional analysis (here or in the next session). Each figure or table must be discussed in the text. Describe the features and the theoretical background of the methods you use for the analysis.]
 - Results: Here, you explain and critically discuss the results of the preprocessing task [Report and comment the main results of the analysis.]
 - References: Here, you provide bibliographic references [Report the books and/or articles that you used for studying the methods and perform the analysis. Each reference reported in this section must be cited in the main text].
- Code listing: The code you used to perform the analysis. Regardless of your choice programming, your code must be executable/functioning. The code (and the relevant functions, if needed) can be either pasted at the end of the 6-page article (for instance as an appendix) or packaged together with the paper as a zip file.

The work can be done individually or in group of maximum 4 co-authors. You can chose to write your paper either in English or Portuguese³ You can base your analysis on the resources you might find on the web but you must adequately reference to them.

¹Also available at: <https://www.ieee.org/conferences/publishing/templates.html>.

²Avoid the obvious title “Homework 1: Data pre-processing”.

³In L^AT_EX, specify `\usepackage[portuguese]{babel}` in the preamble to change the language.

The work must be submitted by NOVEMBER 30, 2024. Extension on this deadline might be considered if unanimously requested at least 1 week prior the set date. Further note that delays will be penalised (<24h: 20% penalty; <48h: 40% penalty; etc.).