

Effort estimation with statistical and ML techniques

Project Management
Sergio de la Mata Moratilla
3º GII E

Index:

1. Introduction: 3

2. Summary Practice: 3

3. Estimation Models: 3

4. Evaluation Measures: 6

5. Other Techniques: 7

6. Considerations:..... 8

7. Bibliography: 9

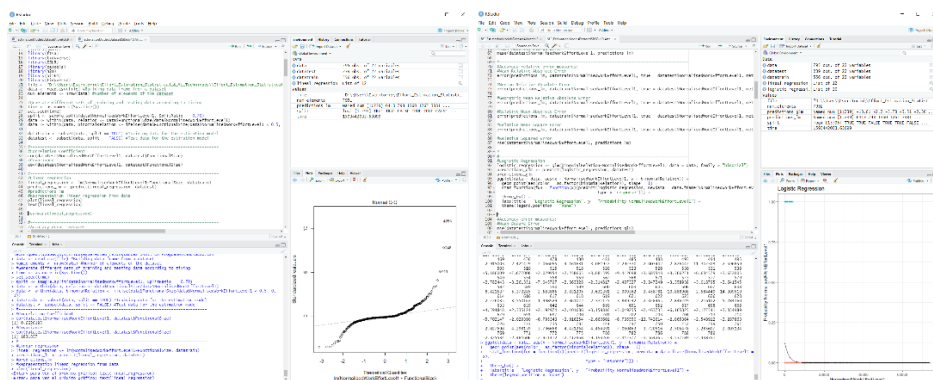
Introduction:

This document contains all the information regarding the performance of the practice. It will explain what it was made, the estimation models used, the evaluation measures used on them and some techniques used for the performance of the practice.

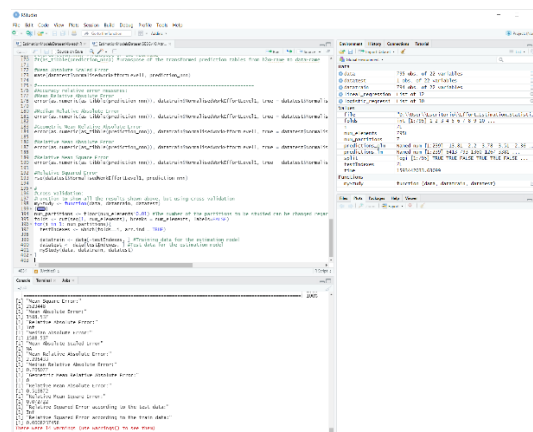
Summary Practice:

This practice asked to build different statistical or machine learning models in order to make estimations from a database of historical projects.

To do it, it was used two datasets given (one from the Github repository given and another provided directly by teacher) where it was used linear regression, logistic regression and neural network models and multiple evaluation measures in order to have different estimation from the datasets.



It has been also used another dataset in order to represent the regression tree model and the technique of cross validation to have a different point of view of estimating dividing data provided in different subsets.

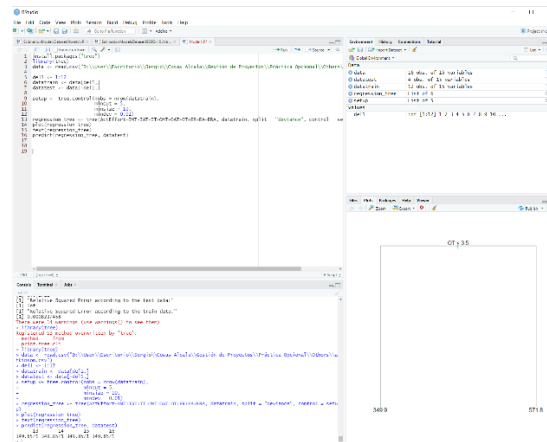


Estimation Models:

In order to make estimations and/or predictions from a set of historical data from a project, there are many estimation models which can be used considering the kind of data which is been used.

For the study of the datasets selected, were used the lineal regression, logistic regression and neural network models. Even though the main study makes use of

these estimation models, it was also seen with a different dataset the regression tree, but due to the restriction it had, it was decided not to be used with the rest of models.



A regression tree is a type of decision in which the leaves represent a numerical value and their idea is to predict the answer variable Y in function of covariables. The problem from this model is its disadvantages. These have made to not be used with the rest models as it could only be made in one of the data sets will in the order had many disadvantages its use. These were some of the disadvantages found:

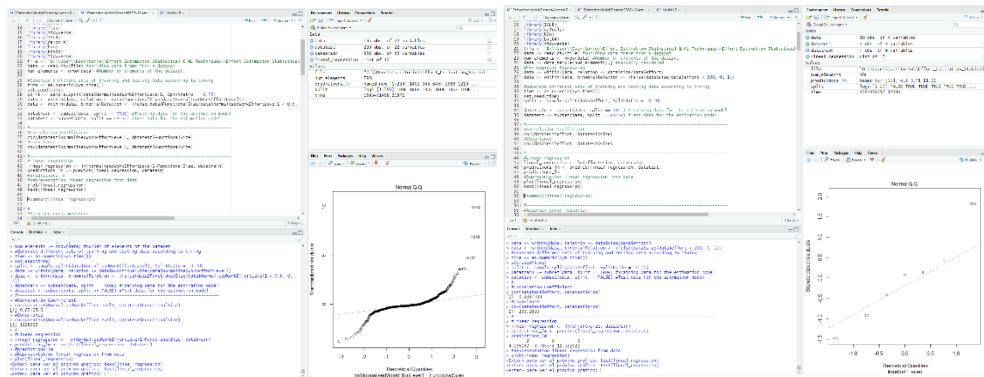
- If it is used the package “tree” as in the picture above, if the dataset provided is small, for example using 10 elements, it might not be able to be made in RStudio as it can distribute the element in a more accurate way.
- The regression tree may have also a lower number of nodes, if the series of terms provided for the prediction is low, being the lowest case only one term.
- When it is obtained a prediction for a series of terms given at first, the results obtained will only have the values which are at the nodes making very restrictive the variety of possibilities obtained and, in the cases in which there is a great variety of data, the error estimations obtained, can be high.

One of the datasets have only 10 elements with two fields each while the other have 795 elements and a huge number of fields. Therefore, the dataset with the lower number of elements is more restricted and it is not useful to be used it in this case.

The linear regression is an estimation model in which it is followed a linear approach between a dependent variable and an independent or a set of independent variables.

The problem with this estimation comes when the structure of the dataset according to its fields used for the estimation builds a different structure or the greater part of the elements are away from the estimation value which would correspond to that element in the estimation function straight line.

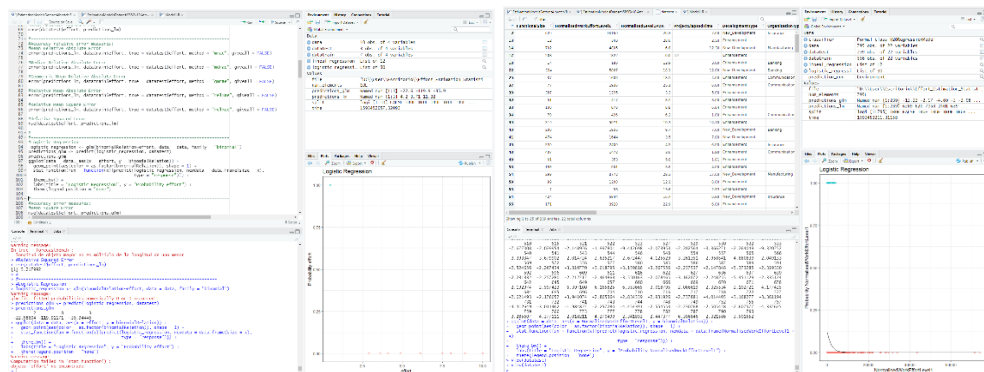
As it can be seen in the next pictures, the estimation obtained for the different values of the elements is really accurate to the different elements of the dataset.



The logistic regression is a statistical model which is used to predict the result of a variable in function from a set of independent variables. This estimation model is used when the dependent variable is binary in nature.

The main problem from this estimation is the binary nature from the dependent variable making impossible to consider intermediate variable unless it is obtained previously a mechanism to classify them using this criterion.

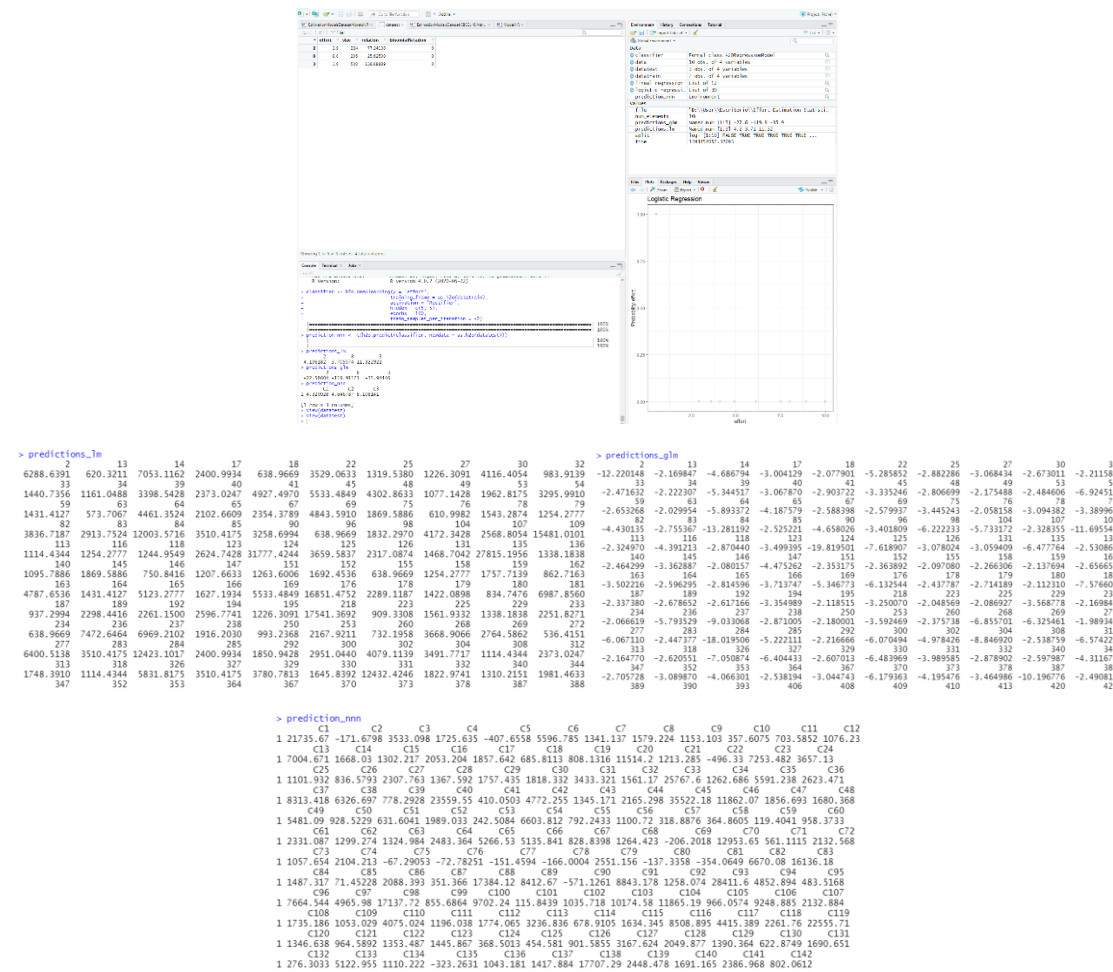
In the case of the datasets selected, it couldn't be done directly with the data given as none of the fields from their different elements had only the values of 0 and 1. To make the study, it was build a small function at the beginning from each of the R files in order to distribute according to a condition the different elements according to the value they had from the relation between the effort and the size of the project. The pictures below show the results from the predictions obtained.



Neural networks are computational model based on the biological neural connections in which a set of units are connected between them to transmit signals. The input information goes through the network obtaining different outputs passed to the next units until it is obtained a set of output data. Some of the disadvantages from this model are its hardware dependence, not quiet explained behaviour from the network, unknowledge from the proper structure of the network, difficulty of representation, ... For this last reason, it has not been represented for the different datasets.

These are the results of the predictions obtained for each of the datasets according to the different test data provided for each of the datasets in each moment using the different estimation models and the values which in really would have. It is needed to

consider that the field value predicted is the one related with effort for all the models and both datasets.

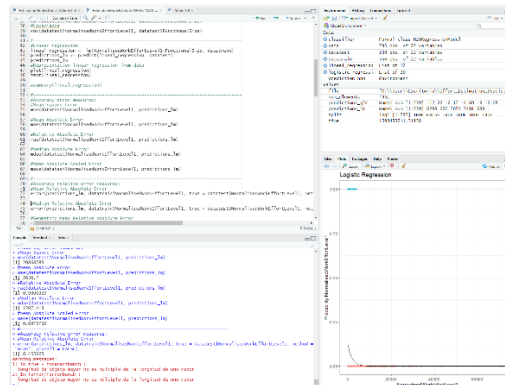


Evaluation Measures:

In order to evaluate how the different estimation projects were more or less accurate to the data given from each dataset, it was used a set of evaluation measures. These would be the ones used:

- Correlation Coefficient.
- Covariance.
- Mean Square Error.
- Mean Absolute Error.
- Relative Absolute Error.
- Median Absolute Error.
- Mean Absolute Scaled Error.
- Mean Relative Absolute Error.
- Median Relative Absolute Error.
- Geometric Mean Relative Absolute Error.
- Relative Mean Absolute Error.
- Relative Mean Square Error.

According to the kind of measure is provided from the ones above, it can be seen when they are executed different values like “Inf” (referring to infinite), “NA” (referring to not available, which happens when there are not enough data to obtain the result) and integer/float numbers and it can be joined with a set of warnings in relation with how the train and test datasets are built.



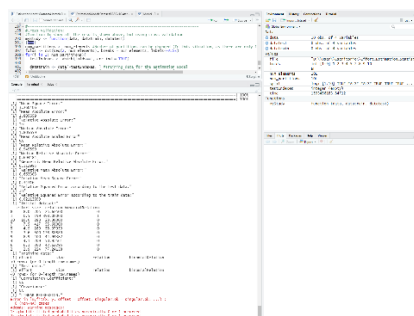
All the measures except the correlation coefficient and the covariance, can change according to the different kind of prediction used if the train and test data are maintained for each of them while correlation and covariance only depend on the dependent and independent variables established to be studied.

Other Techniques:

With all the estimation models and evaluation measures used, it was also used the cross validation which is used to evaluate the different results from different statistical analysis which guarantees the independency between the train and the test data.

The cross validation used for the two datasets is differentiated on how many partitions each of them are represented. In the case of the dataset with a higher number of elements, it was only shown a few of the possible analysis which could be obtained having only one element as train data. In the case of the other dataset, the number of analysis obtain are the ten possible resulting analysis if only it is selected one of its elements as part of the training data.

This technique uses all the operations used to build each of the different models and evaluation measures used and is repeated for each analysis made. In this case, it is necessary to print the different results obtained as the command is not directly executed.



Considerations:

Regarding this practice, it is needed to consider that:

- It was made using RStudio IDE in order to use R language.
- All the packages used for both R files it was needed to be the first time installed with the command `'install.packages("name_package")'`
- To execute any of the operations, it is needed to do the next sequence of keys: Ctrl + Enter.
- The packages used for a practice need to be loaded each time, RStudio is initialized. The command needed is `"library(name_package)"`.
- At the linear regression it is needed sometimes to execute again the command given for its representation to have a correct picture of data. This mainly happens with the linear regression.
- The practice is divided in different sections which can be done at the order the user wants, but it is compulsory to have executed the operations established before the section for correlation coefficient and covariance.
- It is needed to consider that the section related to the cross validation it is recommended to be executed after or without executing the rest of sections, as the variables `"data_train"` and `"data_test"` will be changed from the one given initially making to obtain the same results as when it is made alone.
- The operations at the function used for the cross validation are the same as the ones used for the rest of sections, but with the difference that the results obtained each operation are printed at screen and it is not needed to execute them.
- The idea of having both representations at the R document is to show a specific study having specific set of train and test data using the approximations used for this kind of studies (70-80% of data are for training and 20-30% of data are for testing) distributing them in a random way, and using a cross validation (in this case, the set of train data could only be made of one element).
- The number of validations which can be made with cross validation can be changed at any moment on the source code depending on the specifications of the user.
- The project contains three datasets and three R files. Two of each are used for the study of the linear regression, the logistic regression and the neural networks models while the other show an example of use of the regression tree.

Bibliography:

- Correlations in R: <https://www.statmethods.net/stats/correlations.html>
- Read table in R: <https://www.rdocumentation.org/packages/utils/versions/3.6.2/topics/read.table>
- Data structures in R: <http://ocw.uc3m.es/estadistica/aprendizaje-del-software-estadistico-r-un-entorno-para-simulacion-y-computacion-estadistica/estructuras-de-datos-en-r>
- Introduction to R: <https://cran.r-project.org/doc/contrib/R-intro-1.1.0-espanol.1.pdf>
- Package mlr: <https://cran.r-project.org/web/packages/mlr/mlr.pdf>
- Model Training and Tuning: <https://topepo.github.io/caret/model-training-and-tuning.html>
- Confidence Interval: <https://fhernanb.github.io/Manual-de-R/ic.html>
- Measures of Evaluation in Software Engineering: <http://danroddgar.github.io/DASE/evaluationSE.html>
- Mean Squared Error in Linear Regression in R: <https://stats.stackexchange.com/questions/107643/how-to-get-the-value-of-mean-squared-error-in-a-linear-regression-in-r>
- Package ftsa, error function: <https://www.rdocumentation.org/packages/ftsa/versions/3.4/topics/error>
- Documentation package Metrics: <https://cran.r-project.org/web/packages/Metrics/Metrics.pdf>
- Forecast function in R: <https://www.rdocumentation.org/packages/forecast/versions/8.12/topics/forecast>
- Documentation package ftsa: <https://cran.r-project.org/web/packages/ftsa/ftsa.pdf>
- Forecast error measure in R: <http://finzi.psych.upenn.edu/library/ftsa/html/error.html>
- Forecast error measure in R: <https://rdr.io/cran/ftsa/man/error.html>
- Simple and multiple logistic regression in R: https://rpubs.com/Joaquin_AR/229736
- Introduction Artificial Neural Networks in R: <https://rpubs.com/rdelgado/402754>
- Aggregation and Restructuring data in R: <https://www.r-statistics.com/tag/transpose/>
- Conversion H2OFrame to a data frame: <https://stackoverflow.com/questions/43189340/how-to-convert-a-column-in-h2oframe-to-a-python-list>
- Median in R need numeric data: <https://stackoverflow.com/questions/13204008/median-in-r-need-numeric-data>

- Cross validation in R: <https://www.diegocalvo.es/validacion-cruzada-en-r/>
- Conditionals in R: <https://stackoverflow.com/questions/54825426/how-to-use-conditional-statement-and-return-value-for-a-function-in-r>
- Add calculated fields to data in R: <https://www.dummies.com/programming/r/how-to-add-calculated-fields-to-data-in-r/>
- Logistic Regression in R: <http://idaejin.github.io/courses/R/2019/euskaltel/regresion-logistica.html>
- Simple and multiple logistic regression: https://rpubs.com/Cristina_Gil/Regresion_Logistica
- Example logistic regression with R: <https://rpubs.com/emilopezcano/logit>
- Simple and multiple logistic regression in R: https://rpubs.com/Joaquin_AR/229736
- Creation of random numbers from time stamp: <https://decisionstats.com/2013/09/28/using-r-for-random-number-creation-from-time-stamps-rstats/>
- Introduction to cross validation in R: <https://rpubs.com/rdelgado/405322>
- Minimum sample size to make a correlation test in R: https://www.researchgate.net/post/What_is_the_minimum_sample_size_to_run_Pearsons_R
- Neural Networks: https://es.wikipedia.org/wiki/Red_neuronal_artificial
- Regression Trees: <https://www.sciencedirect.com/topics/computer-science/regression-tree#:~:text=Analogously%20to%20decision%20trees%2C%20regression,by%20discretizing%20the%20target%20variable.>
- Linear Regression: https://en.wikipedia.org/wiki/Linear_regression#:~:text=In%20statistics%2C%20linear%20regression%20is,is%20called%20simple%20linear%20regression.