

# Welcome to Wine Project's documentation!

For this project, we participate in a competition on Kaggle:  
<https://www.kaggle.com/datasets/uciml/red-wine-quality-cortez-et-al-2009>.

This competition aims at a classification task: try to understand whether a wine is good or not from chemical wine characteristics.

For the aim, three modules are built: “main,” which contains the start of the project; “preprocessing,” in which the dataset is prepared for the classification part; “analysis,” which includes algorithms used for prediction. Every subsection is explained in them.

Contents:

- [Main module](#)
- [preprocessing module](#)
- [analysis module](#)

## Indices and tables

- [Index](#)
- [Module Index](#)
- [Search Page](#)

## Contacts

[s.aranciofebbo@studenti.unipi.it](mailto:s.aranciofebbo@studenti.unipi.it)  
[m.secondo@studenti.unipi.it](mailto:m.secondo@studenti.unipi.it)

## Footnotes

# Main module

This module contains the main function.

`main.classification(X_train, X_test, y_train, y_test)`

This function execute various learning models used for the classification task. Then take the results of each algorithm and call `comparison()`.

- Parameters::**
- **X\_train** (*non-target train*,) –
  - **X\_test** (*non-target for prediction test*,) –
  - **y\_train** (*target for prediction train*,) –
  - **y\_test** (*target for prediction test*) –

`main.comparison(models, scores)`

- Parameters::**
- **models** (*classification model analyzed*,) –
  - **scores** (*It is a list that memorizes [f1score, recall, precision] for each model in the sequence obtained using models listed in “models.”*) –

**Return type::** None

`main.feature_Importance(df)`

`main.visualization(df)`

## Footnotes

# preprocessing module

This module implements many functions to transform datasets with preprocessing tasks.

`preprocessing.preprocessing(db)`

In this function are done the significant tasks of preprocessing step are. Are defined the targets of interest. The dataset is divided into a training set and test set(20%), normalized with `MinMaxScaler()`. The training set is balanced with `SMOTE(random_state=42)` algorithm.

Other internal functions are dedicated to plotting some helpful representations of the dataset. If the users want them, they must write “`plt. show()`”.

<b>Parameters::</b>	<b>db</b> ( <i>DataFrame the competition DataFrame</i> ) –
<b>Yields::</b>	<ul style="list-style-type: none"><li>• <b>df</b> (<i>DataFrame structured according to our needs,</i></li><li>• <b>X_train</b> (<i>non-target train,</i>)</li><li>• <b>X_test</b> (<i>target for prediction test,</i>)</li><li>• <b>y_train</b> (<i>non-target for prediction train,</i>)</li><li>• <b>y_test</b> (<i>target for prediction test,</i>)</li><li>• <b>X</b> (<i>attributes</i>)</li></ul>

`preprocessing.remove_outlier(db)`

This function aims to remove all outliers using the Box Plot Diagram. The purpose of this diagram is to identify outliers and discard it from the data series.

<b>Parameters::</b>	<b>db</b> ( <i>DataFrame</i> ) – the competition DataFrame
<b>Yields::</b>	<b>df</b> ( <i>Database cleaned of outliers</i> )

## Footnotes

# analysis module

This module contains all algorithms used for the classification task.

`analysis.DecisionTC(X_train, X_test, y_train, y_test)`

This function is used to visualize the results of the grid searches for the Decision Tree algorithm.

**Parameters::**

- **X\_train** (*non-target train*,) –
- **X\_test** (*non-target for prediction test*,) –
- **y\_train** (*target for prediction train*,) –
- **y\_test** (*target for prediction test*) –

**Yields::** **score** (*it is the result of the `csv_algo(model, grid, y_test, y_pred)` function. A list with [] – macro\_f1, macro\_precision, macro\_recall] computed from `report()`.*

`analysis.KNN_CV(X_train, X_test, y_train, y_test)`

This function is used to visualize the results of the grid searches for the k nearest neighbor algorithm. We used GridSearchCV to test Hyper-parameters with cross-validation algorithms kFold.

In kFold, the data is shuffled once at the start, and then divided into the number of desired splits.

The only Hyper-parameters are used are the `n_neighbors` parameters which contains the interval from 2 to k neighbors, and they test among all these values which one gives the best results

**Parameters::**

- **X\_train** (*non-target train*,) –
- **X\_test** (*non-target for prediction test*,) –
- **y\_train** (*target for prediction train*,) –
- **y\_test** (*target for prediction test*) –

**Yields::**

**score** (*it is the result of the csv\_algo(model, grid, y\_test, y\_pred) function. A list with [] – macro\_f1, macro\_precision, macro\_recall] computed from report()*).

analysis.NaiveBayes(*X\_train, X\_test, y\_train, y\_test*)

This function is used to visualize the results of the grid searches for the Naive Bayes algorithm. As hyper-parameters, we used var\_smoothing.

It is a stability calculation to widen (or smooth) the curve and therefore account for more samples that are further away from the distribution mean.

**Parameters::**

- **X\_train** (*non-target train,*) –
- **X\_test** (*target for prediction test,*) –
- **y\_train** (*target for prediction train,*) –
- **y\_test** (*target for prediction test*) –

**Yields::** **score** (*it is the result of the csv\_algo(model, grid, y\_test, y\_pred) function. A list with [] – macro\_f1, macro\_precision, macro\_recall] computed from report()*).

analysis.RandomForestClassifier\_CV(*X\_train, X\_test, y\_train, y\_test*)

This function is used to visualize the results of the grid searches for the random forest algorithm.

**Parameters::**

- **X\_train** (*non-target train,*) –
- **X\_test** (*non-target for prediction test,*) –
- **y\_train** (*target for prediction train,*) –
- **y\_test** (*target for prediction test*) –

**Yields::** **score** (*it is the result of the csv\_algo(model, grid, y\_test, y\_pred) function. A list with [macro\_f1, macro\_precision, macro\_recall] computed from report()*).

analysis.csv\_algo(*model, grid, y\_test, y\_pred*)

This function prints the algorithm's results in a CSV file.

- Parameters::**
- **model** (*classification model analyzed,*) –
  - **grid** (*results of GridsearchCV from the model,*) –
  - **y\_test** (*target for prediction train,*) –
  - **y\_pred** (*the target used for prediction*) –
- Yields::** **Scores** (*It is a list that contains the f1 average, precision average, and recall average.*)

analysis.logisticRegression(*X\_train, X\_test, y\_train, y\_test*)

This function is used to visualize the results of the grid searches for the Logistic Regression algorithm. As hyper-parameter we used C .It's a penalty term, meant to disincentivize and regulate against Overfitting.

- Parameters::**
- **X\_train** (*non-target train,*) –
  - **X\_test** (*target for prediction test,*) –
  - **y\_train** (*non-target for prediction train,*) –
  - **y\_test** (*target for prediction test,*) –
- Returns::** **score** – macro\_f1, macro\_precision, macro\_recall] computed from report().
- Return type::** it is the result of the csv\_algo(model, grid, y\_test, y\_pred) function. A list with [

## Footnotes

# Python Module Index

[a](#) | [m](#) | [p](#)

**a**

[analysis](#)

**m**

[main](#)

**p**

[preprocessing](#)



# Index

[A](#) | [C](#) | [D](#) | [F](#) | [K](#) | [L](#) | [M](#) | [N](#) | [P](#) | [R](#) | [V](#)

## A

analysis  
[module](#)

## C

[classification\(\).\(in module main\)](#)

[comparison\(\).\(in module main\)](#)  
[csv\\_algo\(\).\(in module analysis\)](#)

## D

[DecisionTC\(\).\(in module analysis\)](#)

## F

[feature\\_Importance\(\).\(in module main\)](#)

## K

[KNN\\_CV\(\).\(in module analysis\)](#)

## L

[logisticRegression\(\).\(in module analysis\)](#)

## M

main

[module](#)

module

[analysis](#)

[main](#)

[preprocessing](#)

## N

[NaiveBayes\(\).\(in module analysis\)](#)

## P

preprocessing

[module](#)

[preprocessing\(\).\(in module preprocessing\)](#)

## R

[RandomForestClassifier CV\(\).\(in module analysis\)](#)

[remove\\_outlier\(\).\(in module preprocessing\)](#)

## V

[visualization\(\).\(in module main\)](#)