



Intelligent Systems

Laboratory activity 2018-2019

Project title: Insight into the grades of Romanian students
Tool: Scikit-Learn

Name: Csaba-Laszlo Gabor
Group: 30434
Email: csabagabor97@gmail.com



Contents

1	Overview	3
1.1	Installing the tool	3
1.2	Overview of the tool	3
2	Main Features	4
2.1	Machine learning Algorithms	4
2.2	Other Algorithms	4
3	Algorithm - Details	5
3.1	General Description	5
3.2	Algorithm in Scikit-Learn	6
4	examples	8
4.1	Provided Examples	8
4.2	Own Examples	10
5	Problem specification	12
5.1	General specification	12
5.2	Data	12
5.3	Scraping the Data	13
5.4	Merging the Data	16
6	Related work	20
7	Preliminary Results	21
7.1	Predicting the target value <code>romana final</code>	21
7.2	Predicting the target value <code>rezultatul final</code>	23
8	Implementation details	26
8.1	Achieving higher accuracy	26
8.2	Optimizations in predicting <code>romana final</code>	32
9	Graphs and experiments	37
10	Related work and documentation	40

Chapter 1

Overview

1.1 Installing the tool

1. Make sure **Python** is installed on the computer(check it with *python -version*). If it is not, install it from the official Python website([8]).
2. Make sure **pip** is installed on the computer. Pip comes pre-installed with new Python releases.
3. To install the tool as a standalone version, just type `pip install scikit-learn` in the terminal
4. These tools: **numpy**, **pandas**, **matplotlib**, **seaborn** need to be installed to process data, to plot functions and to use mathematical functions so for every tool just type `pip install {name of the tool}`
5. Also it is highly recommended to install **Jupyter Notebook** which is a code editor in which separate pieces of code can be executed, commented, plots can be made etc. To install it, type: `pip install jupyter` in the terminal

1.2 Overview of the tool

Scikit-learn is an open-source machine learning library for the Python programming language. It includes several classification and regression algorithms including support vector machines, random forests etc.

Machine learning is all about predicting some labels in unseen data like predicting the height of new children if their age, gender and health condition is known(in this case height is a label while age, gender and health condition are features which help to predict labels). **Scikit-learn** helps in this by automizing some algorithms so the end programmer does not have to know how these algorithms work under the hood (of course it is recommended to have a basic understanding of them - but the programmer does not have to implement these functions from scratch) so he can start working on a problem from the beginning.

Chapter 2

Main Features

2.1 Machine learning Algorithms

There are 2 big types of algorithms:

1. **Classification:** it means that the label we want to predict can have limited possible values. For example classifying if an apple is red or green is a classification problem.
2. **Regression:** regression methods are used to predict a continuous value which can have an unlimited number of values like predicting the cost of a house (this can be 1000 dollars but it can also be 1000.0001 dollars...).

For both of these two categories, **Scikit-learn** comes with predefined algorithms.

1. **Classification:** `DecisionTreeClassifier()`, `RandomForestClassifier()`, `KNeighborsClassifier()`, `SVC()` etc.
2. **Regression:** `LinearRegression()`, `KNeighborsRegressor()`, `RandomForestRegressor()` etc.

These algorithms behave differently for different datasets, so there isn't a best choice between them. Another constraint which needs to be considered is time. There are some really good algorithms like the **Random Forest** but requires more time to run than other simpler algorithms such as **LinearRegression**.

2.2 Other Algorithms

We don't need just machine learning algorithms. We also need algorithms for other tasks such as:

1. splitting the dataset into training and test datasets with **`train_test_split()`**
2. cross-validation with **`cross_val_score()`**
3. testing the accuracy of our predictions with **`accuracy_score()`**

Scikit-Learn comes in handy because it already contains algorithms for these subtasks, so no need to develop our own ones.

Chapter 3

Algorithm - Details

3.1 General Description

The algorithm I've chosen is **Linear Regression**. Probably it is the easiest one to understand and the most intuitive. Imagine that we are given a number of features and we have to predict a label. In this case we want to assign a weight to each feature which shows how much the label depends on it.

$$y' = b + w_1x_1 + w_2x_2 + w_3x_3$$

The picture shows that y' is the label we want to predict, x_1, x_2, x_3 are the features and w_1, w_2, w_3 are the weights. b is called bias.

In the first phase we assign a random value to each weight. Next, based on those values we calculate the predicted label values for all entries in the dataset and compare the results with the correct results. The average of the absolute difference between the correct results and our results is the **training error**. In other words this is our **loss**. Our goal is to minimize this loss. Of course there are several types of errors like:

1. **Mean Absolute Error**: this was discussed previously. This error is calculated easily but has one problem: for this error calculation it does not matter how far the predicted values are from the real values. This is when the Mean squared error helps.

$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n}$$

2. **Mean squared error**: which calculates the power of 2 of the loss difference. It favors closer values over further values from the real values. **This method is usually used.**

$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

3. **Mean Bias Error:** Same as Mean Absolute Error only that we don't take the absolute value of the difference. It is not used often.

$$MBE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)}{n}$$

There are other types of errors which are not discussed here.

We have a loss function now which depends on the weights. Our goal is to find the minimum of this function. The algorithm calculates the new values for the weights using the **gradient descent method**. The general idea of this method is to calculate the partial derivatives which indicate in which direction we have to move to minimize the loss. We can make a big step in that direction or a smaller step. This is called the **learning rate**. A small learning rate will give us the correct solution the same way the bigger one calculates the correct solution but it needs more calculations and more time. On the other hand a too big learning rate can make the algorithm fail to converge. So it won't find the minimum error (the minimum of the error function). Having a very good learning rate is hard to achieve.

3.2 Algorithm in Scikit-Learn

Scikit-Learn comes with a built-in function for Linear regression. This function initializes a `LinearRegression` object:

`LinearRegression(fit_intercept=True, normalize=False, copy_X=True, n_jobs=None)`

1. The **fit_intercept** is a *boolean* parameter. It specifies if we want to include an intercept (bias) to the model. By default it is set to *True*.
2. The **normalize** is a *boolean* parameter. If set, it normalizes the data set.
3. The **copy_X** is a *boolean* parameter. If set, it copies the original **X** passed to the function. If it is *False* then the original **X** may change after running the algorithm.
4. The **n_jobs** is an *int* parameter. It is the number of jobs to use for the computation (to use more cores in a CPU). By default it is set to *None*.

After initialization, the **fit()** method needs to be called on the training set. This function trains the model.

`fit(X, y, sample_weight=None)`

X is the set of features and **y** is the set of labels. With **sample_weight** we can add individual weights for each sample.

After training our model, we have to verify its accuracy. To do this, we run the **predict()** function.

predict(X)

X is the feature set in the test set. This function returns the set of the predicted labels.

To check the mean square error of our model on the test set, we can use the `mean_squared_error(y_test, y_pred,...)` function where `y_pred` is the set returned by the `predict(X)` function and `y_test` is the original label set for the test set obtained after splitting our data set(it contains the correct results).

We can also check the accuracy using the `accuracy_score` function.

Simple usage:

```
regressor = LinearRegression()
regressor.fit(X_train, y_train) #training the algorithm
y_pred = regressor.predict(X_test)

#To retrieve the intercept:
print(regressor.intercept_)

#For retrieving the slope:
print(regressor.coef_)

print('Mean Squared Error:', mean_squared_error(y_test, y_pred))
```

Chapter 4

examples

4.1 Provided Examples

The next example is taken from [9] and [7]. It uses the Heart Disease UCI dataset from Kaggle[6]. The data set includes features like *age*, *cholesterol level of the blood* or *gender*. The example uses the **Random Forest Classifier** which is a decision tree, but replicated multiple times with a randomly taken set of features and the final result is the average of the partial results in the trees. The example predicts if a person has a heart disease or not.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import cross_val_score
from sklearn.metrics import accuracy_score
from sklearn.metrics import mean_squared_error

#GET THE DATA
df = pd.read_csv('heart.csv')

#get the first 5 rows of data
df.head()

#information for the values
df.info()

# We must clean the NaN values , since we cannot train
# our model with unknown values.
# Luckily , there is nothing to clean.
df.isna().sum()

# First of all let's see how many zeros and ones do we have...
negative_target = len(df[df.target == 0])
```



```

positive_target = len(df[df.target == 1])
sns.countplot(x = "target", data = df, palette = "pastel")
plt.xlabel("Target (0 = no, 1= yes)")
plt.ylabel("count")
plt.show()

# It is a classification problem, and more precisely a BINARY
  CLASSIFICATION.
# we divide into two subsets: the training data and the testing data
# X are the explanatory variables, Y is the response variable
  ('target'):
X = df.drop('target', axis=1) # everything except target.
Y = df['target']              # only target.

# Since the amount of data is not extremely large, we will use a
  small test_size (0.10-0.15).
X_train, X_test, Y_train, Y_test = train_test_split(X, Y,
  test_size=0.15)

# Firstly, let's take Random Forest Classifier
rfc = RandomForestClassifier(n_estimators=100)

#train the model
rfc.fit(X_train, Y_train)

#predict unseen data
rfc_predictions = rfc.predict(X_test)
test_error = accuracy_score(rfc_predictions, Y_test)
test_error

# get importances from RF
importances = rfc.feature_importances_

# then sort them descending
indices = np.argsort(importances)

# get the features from the original data set
features = df.columns[0:13]

# plot them with a horizontal bar chart
plt.figure(1)
plt.title('Feature Importances')
plt.barh(range(len(indices)), importances[indices], color='b',
  align='center')
plt.yticks(range(len(indices)), features[indices])
plt.xlabel('Relative Importance')

```

After running the program we get an accuracy score of 89%. The example also makes a plot which shows the relative importance of every feature, the most important being **thelach** feature which is *the maximum heart rate achieved*.

4.2 Own Examples

We are going to predict the final bacculaureate grades of students if we know their grades from 8th Grade(*Evaluate Nationala*). So we'll have only one feature(the grade from 8th Grade) and one label(final bacculaureate grade). This time we will run a simple Linear Regression.

After importing the necessary modules, we read the data and check some of its attributes such as size.

```
df = pd.read_csv('note.csv')
df.keys()
df.shape
df.head()
df['romana final']
df['nota la limba romana']
```

Now we split the data in training and test set and train our model:

```
X = df['nota la limba romana'].values.reshape(-1,1)
Y = df['romana final'].values.reshape(-1,1)
X_train, X_test, Y_train, Y_test = train_test_split(X, Y,
    test_size=0.20)
regressor = LinearRegression()
regressor.fit(X_train, Y_train) #training the algorithm
```

Now we predict unseen data and show the error and accuracy score:

```
Y_pred = regressor.predict(X_test)
mean_squared_error(Y_test, Y_pred)
mean_absolute_error(Y_test, Y_pred)
regressor.score(X_test, Y_test)
res = pd.DataFrame({'Note evaluate': X_test.ravel(),
    'Predicted': Y_pred.flatten(), 'Actual': Y_test.flatten()})
res
```

We get a mean squared error of 2.18 and an absolute error of 1.15. The accuracy score is 45%. Showing the intercept and the coefficient:

```
print(regressor.intercept_)
print(regressor.coef_)
```

We get 0.794 for the coefficient and 0.88 for the intercept. This is an interesting result and shows that the results from the Bacculaureate are worse than the ones from the previous exam. Of course this was expected because the Bacculaureate is a harder exam.

Before judging the results, let's run another example. This time we replace every value in the feature column with the average in that column.

```

avg = df[' nota la limba romana'].mean()
df[' nota la limba romana'] = avg
X = df[' nota la limba romana'].values.reshape(-1,1)
Y = df['romana final'].values.reshape(-1,1)
X_train, X_test, Y_train, Y_test = train_test_split(X, Y,
    test_size=0.20)
regressor = LinearRegression()
regressor.fit(X_train, Y_train) #training the algorithm
Y_pred = regressor.predict(X_test)
mean_squared_error(Y_test, Y_pred)
mean_absolute_error(Y_test, Y_pred)

```

This time we get a mean squared error of 3.86 and an absolute error of 1.63. This shows that, despite the fact that in our previous example, we achieved a low accuracy score, there is a high correlation between this feature and this label.

Chapter 5

Problem specification

5.1 General specification

My project consists of predicting the grades of students in Romania at the Baccalaureate. Also other predictions will be carried out such as predicting the failure of students, the chance to have an exam in a minority language based on the county, predicting the grades for optional subjects for those who had an exam at the same subject 4 years earlier and for those who didn't have etc. Plots will be drawn to visualize the data better. The aim of this problem is not just to have a high accuracy but to carry out as many predictions as we can and to analyze those predictions.

5.2 Data

The data from the 2012's Baccalaureate is used. This data is taken from an Open sourced source [11]. It contains two hundred thousand rows. But we also need data for the grades from 4 years earlier at the *Evaluare Nationala* which will be scraped from the well-known <http://web.archive.org> domain. The extracted data needs to be merged with the first one to form a single unit. The merge will be done on the Name(*nume*) column. We want to have a single data unit for a single person with his/her grades from both exams. The second dataset also contains the average grades for the 4 years in school from grade 5 to grade 8. The final dataset has the following columns: *nume*, *judet8*, *scoala*, *media la admitere*, *media teze nationale*, *media de absolvire*, *nota la limba romana*, *nota la matematica*, *optional8*, *optional8 nota*, *limba materna8*, *limba materna8 nota*, *Unnamed: 13*, *pozitia in ierarhie pe judet*, *pozitia in ierarhie pe tara*, *unitatea de invatamant*, *judetul*, *promotie anterioara*, *forma invatamant*, *specializare*, *romana oral*, *romana nota*, *romana contestatie*, *romana final*, *limba materna*, *materna oral*, *materna nota*, *materna contestatie*, *materna final*, *limba moderna*, *limba moderna nota*, *disciplina profil*, *disciplina profil nota*, *disciplina profil contestatie*, *disciplina profil final*, *optional*, *optional nota*, *optional contestatie*, *optional final*, *competente digitale*, *media*, *rezultatul final*

Some columns which were redundant like the first column which contained information which remained from the scraping process were removed. The data also contained some duplicates because there were students with identical names(in this case the merging process could have merged 2 **different** students with identical names - every such student should have been checked manually and added to the data set but because there were only few of them every duplicate was eliminated).

Some info about the type of the data for every column:

nume	10237 non-null object
judet8	10237 non-null object
scoala	10237 non-null object
media la admitere	10237 non-null object
media teze nationale	10237 non-null float64
media de absolvire	10237 non-null float64
nota la limba romana	10237 non-null float64
nota la matematica	10237 non-null float64
optional8	10237 non-null object
optional8 nota	10237 non-null object
limba materna8	10237 non-null object
limba materna8 nota	10237 non-null object
pozitia in ierarhie pe judet	10237 non-null int64
pozitia in ierarhie pe tara	10237 non-null int64
unitatea de invatamant	10237 non-null object
judetul	10237 non-null object
promotie anterioara	10237 non-null object
forma invatamant	10237 non-null object
specializare	10237 non-null object
romana oral	10199 non-null object
romana nota	10237 non-null float64
romana contestatie	1108 non-null float64
romana final	10237 non-null float64
limba materna	944 non-null object
materna oral	943 non-null object
materna nota	944 non-null float64
materna contestatie	64 non-null float64
materna final	944 non-null float64
limba moderna	10237 non-null object
limba moderna nota	9176 non-null object
disciplina profil	10237 non-null object
disciplina profil nota	10237 non-null float64
disciplina profil contestatie	1205 non-null float64
disciplina profil final	10237 non-null float64
optional	10237 non-null object
optional nota	10237 non-null float64
optional contestatie	751 non-null float64
optional final	10237 non-null float64
competente digitale	10184 non-null object
media	10237 non-null float64
rezultatul final	10237 non-null object

We can see that *Pandas* recognized most of the columns with grades as float values and the columns with strings as object values. The other columns for which the type was inferred incorrectly will be corrected manually. Before we run any algorithm on string columns, we have to encode them with the power of *scikit*. This is a rather complex data set with 42 columns and more than 10000 rows (when merging the 2 datasets only a portion of the names were common in both datasets so the final size of the dataset is smaller than the initial ones - some students didn't finish high school, others didn't go to the Baccalaureate exam - the reasons can be many). Also, the column names are self-explanatory.

5.3 Scraping the Data

For scraping, the following *JavaScript* code was used. This code can be just copied in the console of any modern browser after injecting/importing *jQuery*.

After selecting the county in [1], the following piece of code extracts the links which take us to the pages of schools.

```
a = $("tdc > .lnk");

var hrefs = new Array();
a.each(function(){
    hrefs.push($(this).attr('href'));
})
```

```
$("body").html("");
```

```
hrefs.forEach(function(item) {
    $.get(window.location.href+item, function(data){
        var t = $("<div id=" + item + "></div>");
        t.prepend(data);
        $('body').append(t);
    });
});
```

The following code extracts the links which take us to the tables where the students are shown in each school and loads the tables from the saved links into different *div* elements.

```
var hrefs = new Array();
$('body').children('div').each(function () {
    var id = $(this).attr('id');
    a = $(".tdc > .lnk", this);

    a.each(function(){
        hrefs.push( id + $(this).attr('href') );
    })
});
```

```
$("body").html("");
```

```
var t = $('<table width="88%" align="center" cellpadding="0"
    class="mainTable" id="table"><tbody
    id="tbody"></tbody></table>');
$('body').append(t);
```

```
hrefs.forEach(function(item) {
    $.get(window.location.href+item, function(data){
        var t = $("<div></div>");
        t.prepend(data);
        $('body').append(t);
    });
});
```

This code contains 2 functions which were taken from[4] and are used to save the data in a table into a .csv file. Also it merges the small tables into a big one.

```

function download_csv(csv, filename) {
    var csvFile;
    var downloadLink;

    // CSV FILE
    csvFile = new Blob([csv], {type: "text/csv"});

    // Download link
    downloadLink = document.createElement("a");

    // File name
    downloadLink.download = filename;

    // We have to create a link to the file
    downloadLink.href = window.URL.createObjectURL(csvFile);

    // Make sure that the link is not displayed
    downloadLink.style.display = "none";

    // Add the link to your DOM
    document.body.appendChild(downloadLink);

    // Lanzamos
    downloadLink.click();
}

function export_table_to_csv(html, filename) {
    var csv = [];
    var rows = html.querySelectorAll("table tr");

    for (var i = 0; i < rows.length; i++) {
        var row = [], cols = rows[i].querySelectorAll("td, th");

        for (var j = 0; j < cols.length; j++)
            row.push(cols[j].innerText);

        csv.push(row.join(","));
    }

    // Download CSV
    download_csv(csv.join("\n"), filename);
}

```

```

bodies = $('#mainTable > tbody ');

var t = $('<table width="88%" align="center" cellpadding="0"
    id="BIGtable"><tbody id="tbody"></tbody></table>');
$('#body').append(t);

bodies.each(function() {
    $(this).children('tr').not(':first').each(function() {
        t.append($(this));
    })
})

```

This final code saves the table in a .csv file.

```

var html = document.getElementById("BIGtable");
export_table_to_csv(html, "table.csv");

```

Note that these snippets must be run after each other, one at a time, with some time between them because the loading of data is done using *Ajax* and it needs time to load the data.

5.4 Merging the Data

Loading the data and transforming the values in the *nume* column to lowercase and stripping the whitespace characters. Then merging the result into a single *DataFrame*.

```

df1 = pd.read_csv('eval.csv')
df2 = pd.read_csv('bac.csv')
df1.nume = df1.nume.str.lower()
df2.nume = df2.nume.str.lower()
df1.nume = df1.nume.str.strip()
df2.nume = df2.nume.str.strip()
df = pd.merge(df1, df2, how='inner', on=['nume'])

```

Duplicates in the *nume* column were identified and eliminated.

```

sum(df['nume'].duplicated())
df.drop_duplicates(subset='nume', keep=False, inplace=True)

```


Columns without values were identified.

```
df.isna().sum()
```

nume	0
judet	0
scoala	0
media la admitere	14
media teze nationale	0
media de absolvire	0
nota la limba romana	0
nota la matematica	0
optional8	0
optional8 nota	0
limba materna8	0
limba materna8 nota	0
pozitia in ierarhie pe judet	0
pozitia in ierarhie pe tara	0
unitatea de invatamant	0
judetul	0
promotie anterioara	0
forma invatamant	0
specializare	0
romana oral	48
romana nota	0
romana contestatie	9506
romana final	0
limba materna	9681
materna oral	9683
materna nota	9681
materna contestatie	10601
materna final	9681
limba moderna	0
limba moderna nota	1093
disciplina profil	0
disciplina profil nota	0
disciplina profil contestatie	9407
disciplina profil final	0
optional	0
optional nota	0
optional contestatie	9882
optional final	0
competente digitale	62
media	0
rezultatul final	0

It is unacceptable to have rows in which *media la admitere* column does not have a value so these rows(14 rows in total) will be eliminated.

```
df.dropna(subset=['media la admitere'], how='all', inplace = True)
```

Also, after the scraping process, some columns have invalid data such as the *media la admitere* which has *string* in some columns. These rows will be deleted. Also, other rows with NaN values will be deleted except those from columns like *limba materna8 nota* - where the majority of rows contain NaN values anyway.

```
df[df.columns[3]] = df[df.columns[3]].apply(pd.to_numeric,
      errors='coerce').fillna(0).astype(float)
df = df[df['media la admitere'] != 0]

df['optional8 nota'] = df['optional8 nota'].astype(float)
df = df.dropna(axis=0, how='any', subset=['competente digitale'])
df = df.dropna(axis=0, how='any', subset=['romana oral'])
```

Finally we save the corrected data set:

```
df.to_csv("note.csv", index=False)
```

The final dataset contains 10112 rows and 42 columns and the column types are identified correctly (except the column *limba materna8 nota* because only some students had this exam and '-' appears in some places, but these lines will be removed when this column is included in an analysis).

nume	10172 non-null object
judet8	10172 non-null object
scoala	10172 non-null object
media la admitere	10172 non-null float64
media teze nationale	10172 non-null float64
media de absolvire	10172 non-null float64
nota la limba romana	10172 non-null float64
nota la matematica	10172 non-null float64
optional8	10172 non-null object
optional8 nota	10172 non-null float64
limba materna8	10172 non-null object
limba materna8 nota	10172 non-null object
pozitia in ierarhie pe judet	10172 non-null int64
pozitia in ierarhie pe tara	10172 non-null int64
unitatea de invatamant	10172 non-null object
judetul	10172 non-null object
promotie anterioara	10172 non-null object
forma invatamant	10172 non-null object
specializare	10172 non-null object
romana oral	10134 non-null object
romana nota	10101 non-null float64
romana contestatie	1103 non-null float64
romana final	10101 non-null float64
limba materna	944 non-null object
materna oral	943 non-null object
materna nota	934 non-null float64
materna contestatie	64 non-null float64
materna final	934 non-null float64
limba moderna	10172 non-null object
limba moderna nota	9111 non-null object
disciplina profil	10172 non-null object
disciplina profil nota	10088 non-null float64
disciplina profil contestatie	1196 non-null float64
disciplina profil final	10088 non-null float64
optional	10172 non-null object
optional nota	10088 non-null float64
optional contestatie	747 non-null float64
optional final	10088 non-null float64
competente digitale	10119 non-null object
media	10111 non-null float64
rezultatul final	10172 non-null object

Chapter 6

Related work

As far I know nobody has published an analysis on the Baccalaureate grades in Romania but I've found related work from other countries which predict the grades of students. These are [10], [2] and [3].

[10] works on the prediction of the grades of Portuguese students in 2 subjects: Mathematics and Portuguese language. It works on a dataset which includes several features like sex, age, family size, the job of the father, the job of the mother, extra-curricular activities etc. This data set isn't similar to mine in any aspect because it just contains information about the environment of the student and wants to see in which environments students achieve better. My data set is more concrete than this because it has some concrete past data about students such as their past grades. The majority of Kernels which try to solve this problem use decision trees which seems a good choice in this case.

[2] is the most similar to my problem because it is also about the Baccalaureate grades(in another country). But the data set only contains 8 features, and the problem with these 8 features is the same as with [10] - they are too general. There is only a single Kernel working on this data set which tries to visualize the data with plots.

[3] is the most interesting one because it tries to solve the general problem of predicting grades, not just a local one. It tries using *Bayesian Linear Regression*. The data set it uses is a broad one with 33 features[13]. The good part about this tutorial is that it uses the same tools I'll be using. After following this tutorial I realized that making plots of the data, making plots which show the relation between features is a really important thing because it helps in choosing the best features for the current analysis. Overall, this related work is worth reading.

Chapter 7

Preliminary Results

Two separate experiments were carried out. Code is not pasted here but it is attached and every step is commented to better understand the reason of the steps that were done.

7.1 Predicting the target value romana final

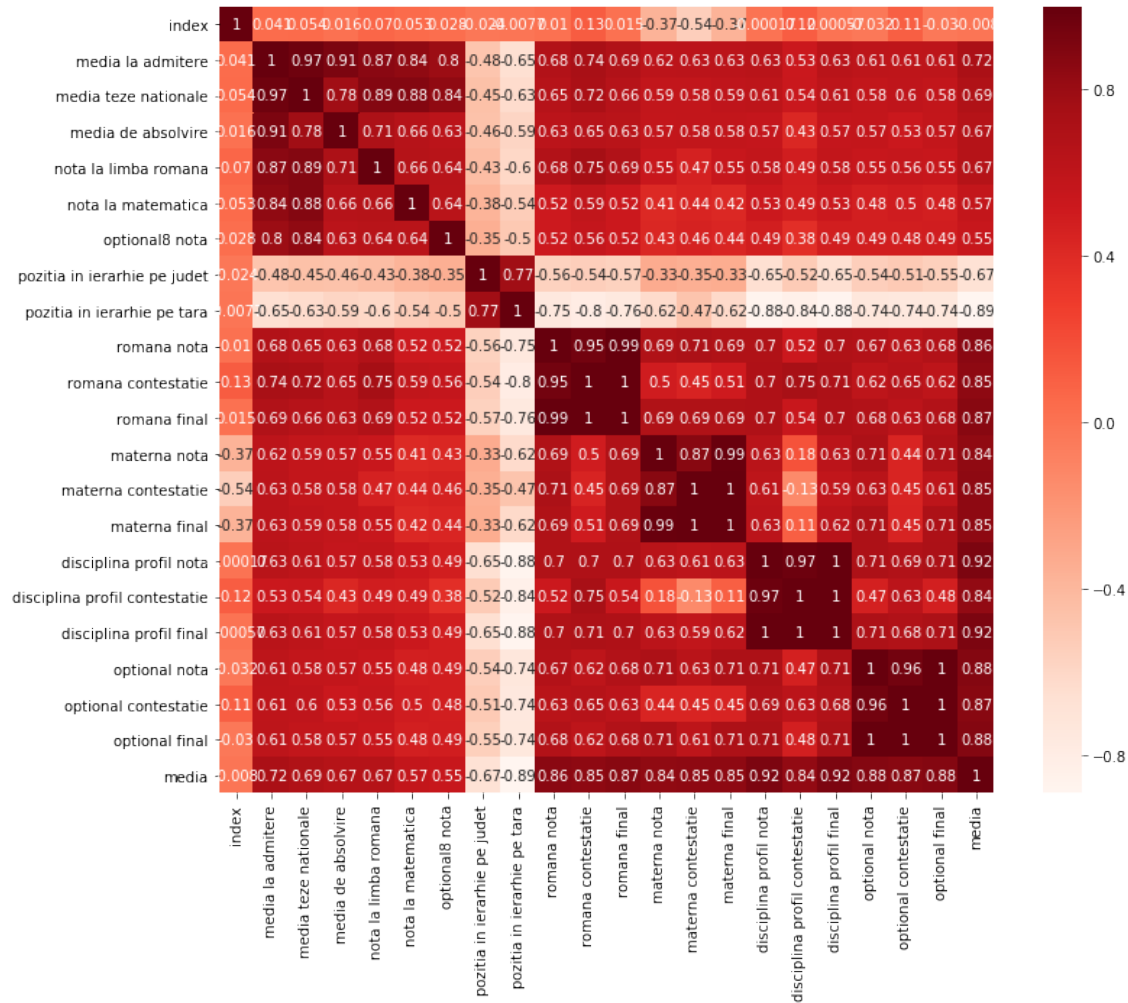
`romana final` is the grade for the final Baccalaureate grade in Romanian language. It is a continuous value ranging from 0 to 10. So, this is a regression problem. Some information about the target value:

```
count    10112.000000
mean      6.893004
std       2.040470
min       0.000000
25%       5.450000
50%       7.200000
75%       8.600000
max       10.000000
Name: romana final, dtype: float64
```

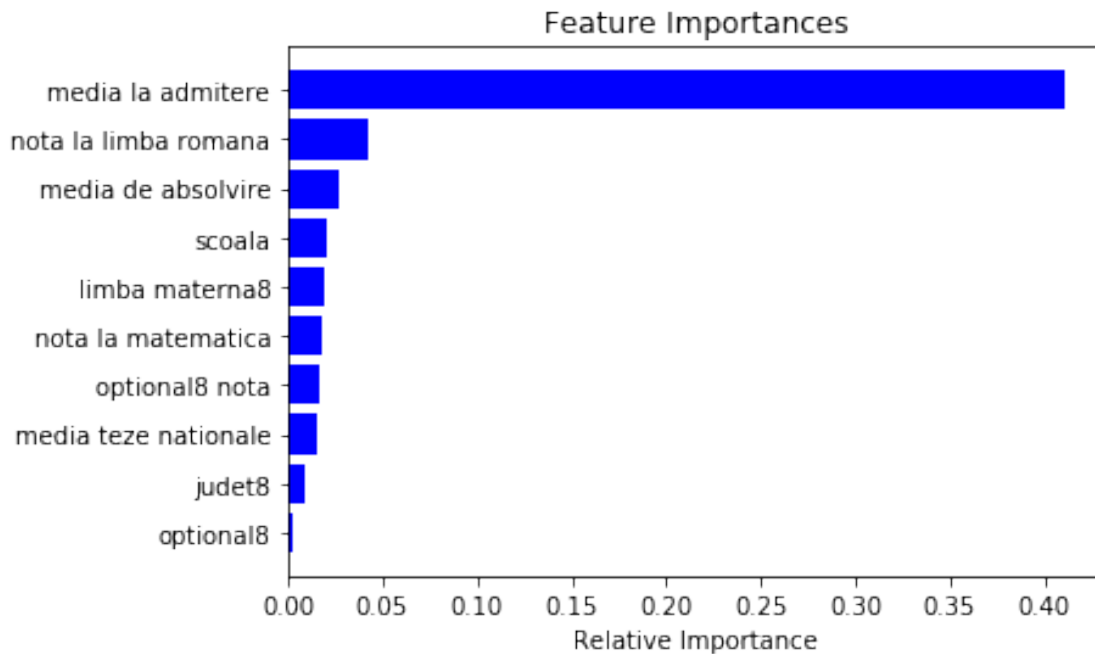
This was a tough experiment because a lot of feature engineering needed to be done, namely:

1. Removing highly correlated columns, because there are 42 columns in total with many categorical features and we must speed up the algorithms somehow. Besides this, some algorithms don't work well with highly correlated features such as *Linear Regression*. Several features were excluded because of correlation such as `'..nota'` and `'..contestatie'` for every type of grade because the `'...final'` column contains the final grade which is of our interest. Yes, some observations could have been made regarding who goes to dispute the grade(*contestare*), but because the columns contain sparse data I excluded them to speed up the algorithm. We can see that `'media teze nationale'` and `'media la admitere'` are highly correlated(because one is used in the calculation of the other one). Also, `'nota la limba romana'` and `'optional8 nota'` and `'nota la matematica'` are highly correlated with both the `'media la admitere'` and `'media teze nationale'` These are obvious because some of the columns are just the average of other columns, so these will be removed. But there is a neat trick here: we won't include the `'limba materna8 nota'` feature because it has many empty values, but instead we include both the `'media teze nationale'` and the

grades for other subjects because in this way 'media teze nationale' will contain in itself the value for 'limba materna8 nota'.



2. Transforming categorical features. **Random Forest** can work well with *Label Encoded* features and it is also quicker. One hot encoded features were used initially with Random Forest to see the feature importance for every subcategory. In this way, we got some interesting results such as *Harghita* county in itself is a very good predictor for the final target value. *Nume* feature was replaced with firstnames, and students can contain multiple firstnames so they were one-hot encoded to allow multiple firstnames for a single person. Also, 'limba moderna nota' was a categorical feature but was replaced with a single number which is the average of the grades got at the exam with $A1 = 1$ and $B2 = 4$. For those who have taken another language exam, their grades were replaced with the maximum grade of 4.0.



After running a *Random Forest Regressor* we got an accuracy of 70%. Cross validation showed a score of 67% which is still a high value. Several other algorithms were tried. Only one was with success(**GradientBoostingRegressor**) with an accuracy score of 72%.

The next step will be to remove outliers because probably they are causing overfitting in our model.

7.2 Predicting the target value rezultatul final

rezultatul final is the column which shows if a student passed the final Baccalaureate. So, it makes sense to not include grades from the Baccalaureate subjects in the features because we'll get a high accuracy, but without any general knowledge. What we'd like to do instead is to predict this value **ONLY** knowing the grades from *Evaluate Nationala* and eventually some categorical values from the Baccalaureate such as the name of the optional subject(but **NOT** the grade itself, just the name) or the type of the High School class(*Mathematics*, *Computer-Science*, *etc.*).

In this experiment, feature engineering wasn't so heavy, only some categorical features had to be transformed into numerical ones like 'specializare', 'limba moderna', 'disciplina profil' etc.

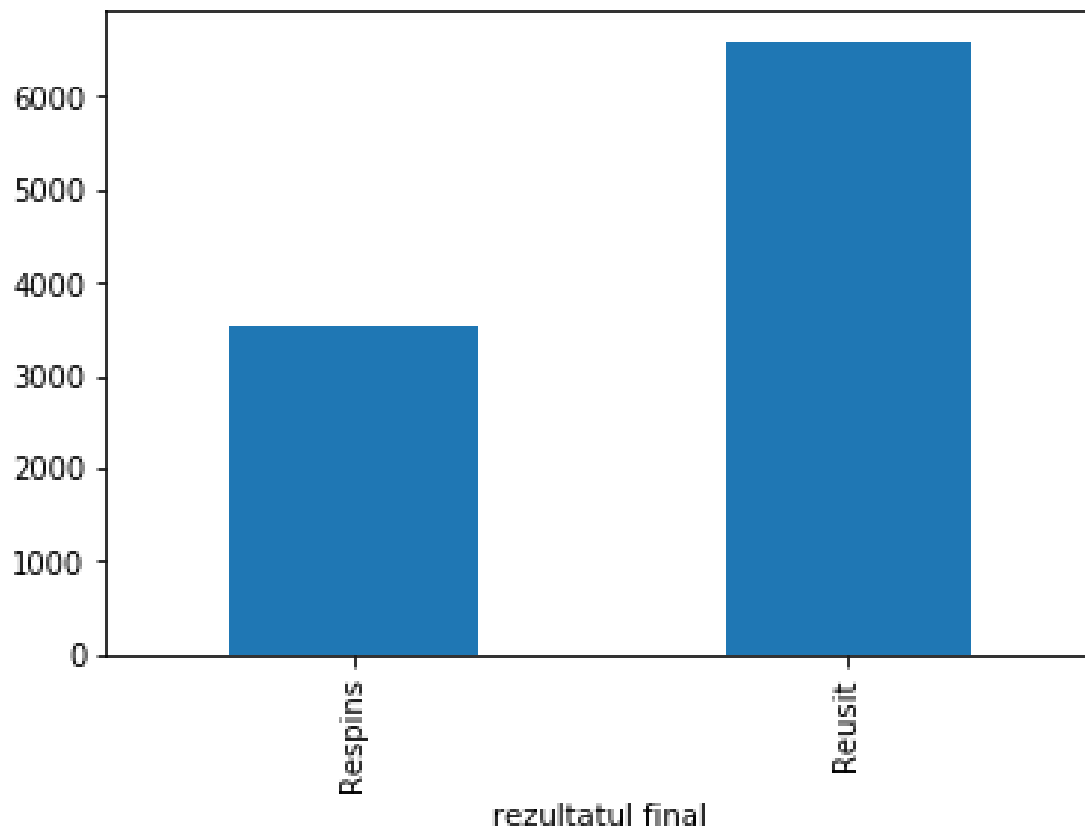
After running a *Random Forest Classifier* we got an accuracy of 83%. For classification problems, it is important to calculate the confusion matrix which shows if there are many false negatives or false positives(this usually happens because the data is unbalanced). This wasn't the case here:

```
#in a classification problem they are really important
print confusion_matrix(Y_test, Y_pred)

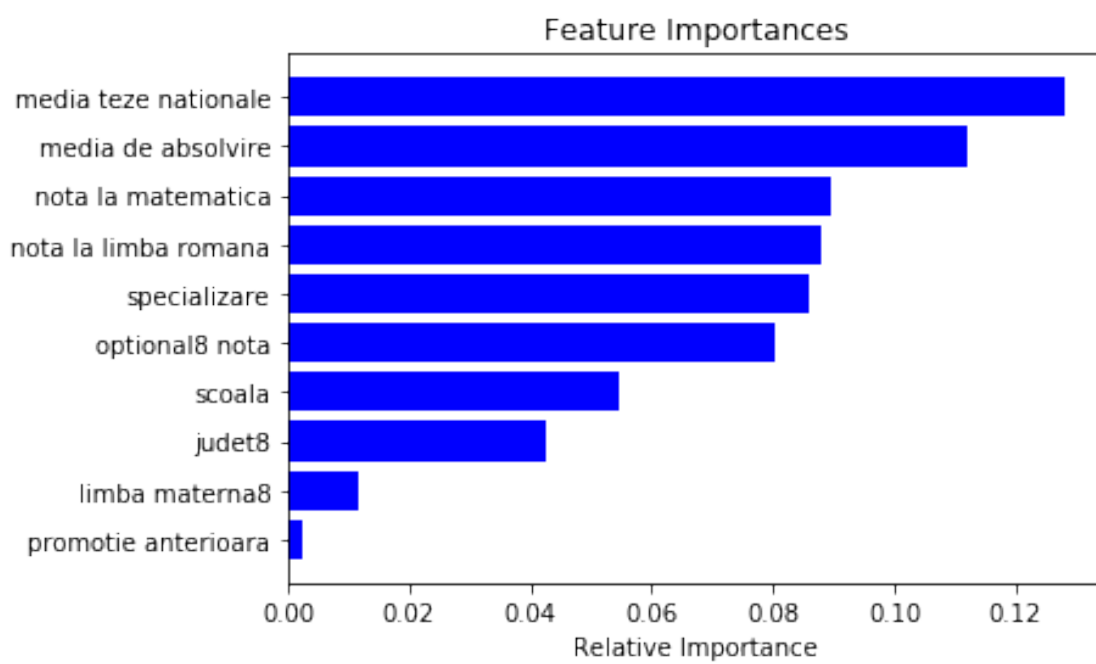
print precision_score(Y_test, Y_pred, average='macro')
print recall_score(Y_test, Y_pred, average='macro')
```

```
[[ 535  171]
 [ 168 1149]]
0.8157393637656796
0.8151138629514671
```

We have a high recall and precision score, which is really nice. In fact, our data is not unbalanced indeed:



These are the feature importances. Feature importances are of great interest because we can see where we have to do more feature engineering, which features can be removed easily without having an impact on accuracy (but making the training phase faster).



Chapter 8

Implementation details

8.1 Achieving higher accuracy

One idea was to extract the gender from the first names. This was done using an online api([5]). A .csv file was uploaded to the API and the api added a gender column and accuracy column to this column. There were 25 names where the API couldn't guess the gender, I completed those values myself.

```
firstName_gender = pd.read_csv('all_names_gender.csv')
firstName_gender.isnull().sum()
```

There were 25 missing values, 'ga_gender' column was completed by me:

```
Unnamed: 0      0
Unnamed: 0.1    0
name           0
ga_first_name   25
ga_gender       0
ga_accuracy     25
ga_samples      25
dtype: int64
```

This code was used to include gender information as a feature:

```
#add gender
genders = []

for ind, name in enumerate(X['name']):
    all_first_names = name.split('-')
    male = 0
    female = 0

    for first in all_first_names:
        if not firstName_gender[firstName_gender['name'] ==
            first]['ga_gender'].empty:
```

```

gender = firstName_gender[firstName_gender['name'] ==
    first]['ga_gender'].values[0]
if gender == 'male':
    male +=1
else:
    female +=1

if male > female:
    genders.append(0)
else:
    genders.append(1)

```

```
X_dummy['gender'] = genders
```

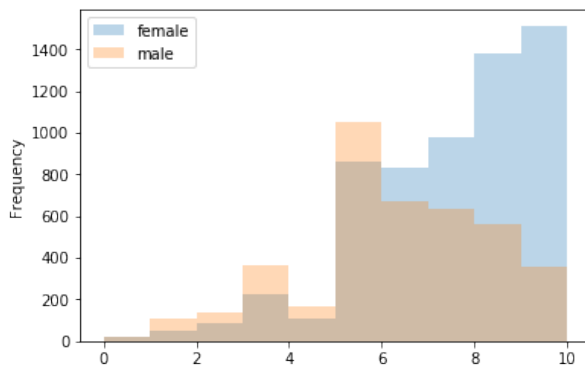
Plot to see the correspondence between gender and target value(INDEED: there is a high correlation: girls are more likely to get a better grade):

```

full_X_Y = pd.concat([X_dummy, Y, df['romana final']], axis = 1)
full_X_Y['gender'] = full_X_Y['gender'].map({1: 'female', 0: 'male'})

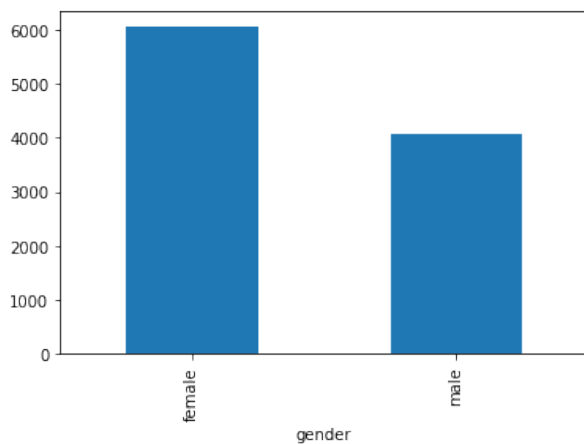
full_X_Y.groupby('gender')['romana final'].plot(kind='hist',
    legend=True, alpha=0.3)

```



Also, we can see that there are more girls than boys:

```
full_X_Y.groupby(full_X_Y['gender']).size().plot(kind='bar')
```



After running the same Random Forrest Classifier, we get an accuracy which is 1% higher. Indeed, gender is a good feature.

Taking a look at the confusion matrix:

```
In [45]: #wow, great, just by including a gender column, our accuracy improved by 1 percent
#now check the recall and precision
print confusion_matrix(Y_test, Y_pred)

print precision_score(Y_test, Y_pred, average='macro')
print recall_score(Y_test, Y_pred, average='macro')

[[ 539  167]
 [ 149 1168]]
0.829168299799669
0.8251600878466598
```

We can see that the precision score and recall score is approximately the same. There are 167 false positives and 149 false negatives.

Now try to do another feature engineering: make another feature from the first word of the "scoala" feature and one from the 'unitatea de invatamant':

```
scoala = []
unitate = []

for ind, row in X.iterrows():
    scoala.append(X['scoala'][ind])
    unitate.append(X['unitatea de invatamant'][ind])

scoala2 = []
for sc in scoala:
    full = sc.split(" ")
    scoala2.append(full[0])

unitate2 = []
for sc in unitate:
    full = sc.split(" ")
    unitate2.append(full[0])

X['scoala_tip'] = scoala2
```

```
X['unitate_tip'] = unitate2
```

Looking at the new feature values:

```
In [50]: X['scoala_tip'].unique()
Out[50]: array(['SC.CU', 'LICEUL', 'LIC.TEORETIC', 'COLEGIUL', 'SCOALA', 'GRI',
               'GRUP', 'GR.SC.', 'SC.I-VIII', 'GRUPUL', 'GIMNAZIUL', 'FUNDATIA',
               'SC.', 'CENTRUL', 'COLEGIL'], dtype=object)

In [51]: X['unitate_tip'].unique()
Out[51]: array(['Colegiul', 'Liceul', 'Grupul', 'Grup', 'Centrul', 'Seminarul',
               'Colegiu', '.liceul', 'Scoala'], dtype=object)

In [ ]: #we can see that there are similar names, we won't replace them with the more generic one
         #because they may hold important information in them
```

Running the algorithm again does not increase the accuracy score. In fact, it decreases slightly to 83,6%. Confusion matrix is almost the same. Ok, so this step didn't help us, the accuracy score is not higher now.

Let's try PCA to reduce the number of features - to speed up the Random Forest algorithm so we can use more trees in it:

```
from sklearn.decomposition import PCA
```

```
#reduce feature size to 75%
pca = PCA(.75)
principalComponents = pca.fit_transform(X_dummy)
```

```
X_dummy_pca = pca.transform(X_dummy)
```

Great, even reducing our feature size to 75% didn't decrease the accuracy score. And the running time is a lot better now.

```
print accuracy_score(Y_pred, Y_test)

print confusion_matrix(Y_test, Y_pred)

print precision_score(Y_test, Y_pred, average='macro')
print recall_score(Y_test, Y_pred, average='macro')
```

```
0.8388531883341572
[[ 548  158]
 [ 168 1149]]
0.8222377998999799
0.8243206618183226
```

Now try to increment n_estimators to 500:

```
0.8358872960949085
[[ 535  171]
 [ 161 1156]]
0.8199080329842614
0.8177714180008218
```

Results are not better, and running time is awful...

Now it's time to experiment with all the algorithms available:

	rec	prec	algorithm
0	0.819108	0.815486	AdaBoostClassifier
1	0.797446	0.787974	BaggingClassifier
2	0.823028	0.818589	ExtraTreesClassifier
3	0.825686	0.822584	GradientBoostingClassifier
4	0.818823	0.808511	RandomForestClassifier
5	0.807659	0.820276	GaussianProcessClassifier
6	0.823795	0.828954	LogisticRegressionCV
7	0.755645	0.816438	PassiveAggressiveClassifier
8	0.815523	0.821478	RidgeClassifierCV
9	0.606401	0.766607	SGDClassifier
10	0.725984	0.705378	Perceptron
11	0.827271	0.814310	BernoulliNB
12	0.652669	0.652972	GaussianNB
13	0.807426	0.811656	KNeighborsClassifier
14	0.778236	0.788949	SVC
15	0.801833	0.804071	NuSVC
16	0.803834	0.808453	LinearSVC
17	0.775622	0.776538	DecisionTreeClassifier

We can see that tree based algorithms work very well. Also, despite its simplicity, Logistic Regression is working as well as the other more complex ones.

OK, now try to use Stacking, which means running some algorithms and including their predictions as features and running another set of algorithms on top of this new dataset. I will use the **vecstack** library.

```
from xgboost import XGBClassifier
from vecstack import stacking

models = [
#     KNeighborsClassifier(n_neighbors=8,
#                           n_jobs=-1),
#     GradientBoostingClassifier(n_estimators=100),
#     KNeighborsClassifier(n_neighbors = 20),
```

```

AdaBoostClassifier() ,
LogisticRegression(max_iter = 100) ,
    ExtraTreesClassifier(n_estimators=100) ,

RandomForestClassifier(random_state=0, n_jobs=-1,
                        n_estimators=100, max_depth=3) ,

XGBClassifier(n_jobs=-1, n_estimators=100, max_depth=3)

]

```

```

S_train , S_test = stacking(models ,
                             X_train , Y_train , X_test ,
                             regression=False ,
                             mode='oof_pred_bag' ,
                             needs_proba=False ,
                             save_dir=None ,
                             metric=accuracy_score ,
                             n_folds=4 ,
                             stratified=True ,
                             shuffle=True ,
                             random_state=0 ,
                             verbose=2)

```

For the second round, I'll be using the powerful **xgboost** library.

```

model = XGBClassifier(random_state=0, n_jobs=-1, learning_rate=0.1,
                      n_estimators=100, max_depth=3)

model = model.fit(S_train , Y_train)
Y_pred = model.predict(S_test)
print('Final prediction score: ', accuracy_score(Y_test , Y_pred))

```

```

In [83]: #ok, so this is our final, tuned accuracy score
         #Looking at the recall and precision scores
         print confusion_matrix(Y_test, Y_pred)

         print precision_score(Y_test, Y_pred, average='macro')
         print recall_score(Y_test, Y_pred, average='macro')

[[ 585  121]
 [ 203 1114]]
0.8222050391500031
0.8372368525772154

In [ ]: #all in all, this is a good result, because we are predicting if a student passes the Baccalaureate without knowing any grade
        #in advance, in a real life scenario everything can happen that causes a student to fail

In [87]: #cross validating our result
         scores = cross_val_score(model, X_dummy_pca, Y, cv=3)
         print scores
         print scores.mean()

[0.77402135 0.79020772 0.81246291]
0.7922306584861895

```

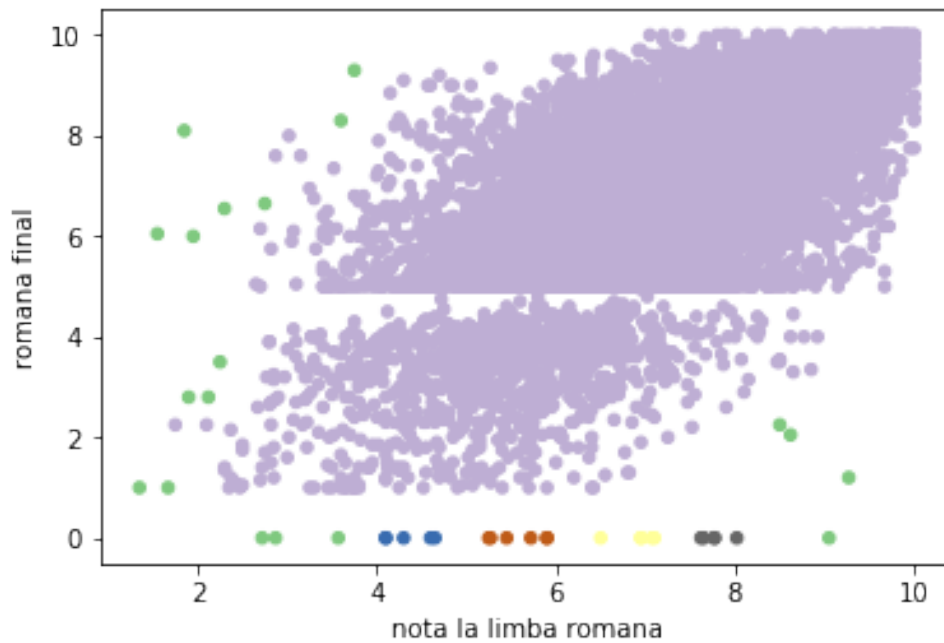
Seeing the result, we can conclude that stacking does not work as predicted. Cross validation shows an accuracy of 79% which is acceptable.

Neural networks were tried with an iteration count of 100, 3 hidden layers of size 10. The result is an accuracy of 79%. One problem with neuronal networks is that they are slow. So not much hyper-parameter tuning can be achieved. Another test was carried out using 20 as the hidden layer size and using only 2 hidden layers. This resulted in an accuracy of 80%.

Also AdaBoostClassifier was tried using a DecisionTreeClassifier as a base learner. The accuracy that resulted was 82.5%. Also, it was tried with a RandomForestClassifier base learner resulting in an even lower accuracy score.

8.2 Optimizations in predicting romana final

Outliers were detected and removed from the training set to learn the important patterns and do not overfit. Taking a look at the plot which shows the correlation between the final Romanian grade and the one from 8th grade. The well-known **DBSCAN** algorithm was used to identify outliers. Outliers are grouped into categories and each category has its own colour. This algorithm also has a parameter **eps** which tells the algorithm how far each point has to be from other points to be considered an outlier.

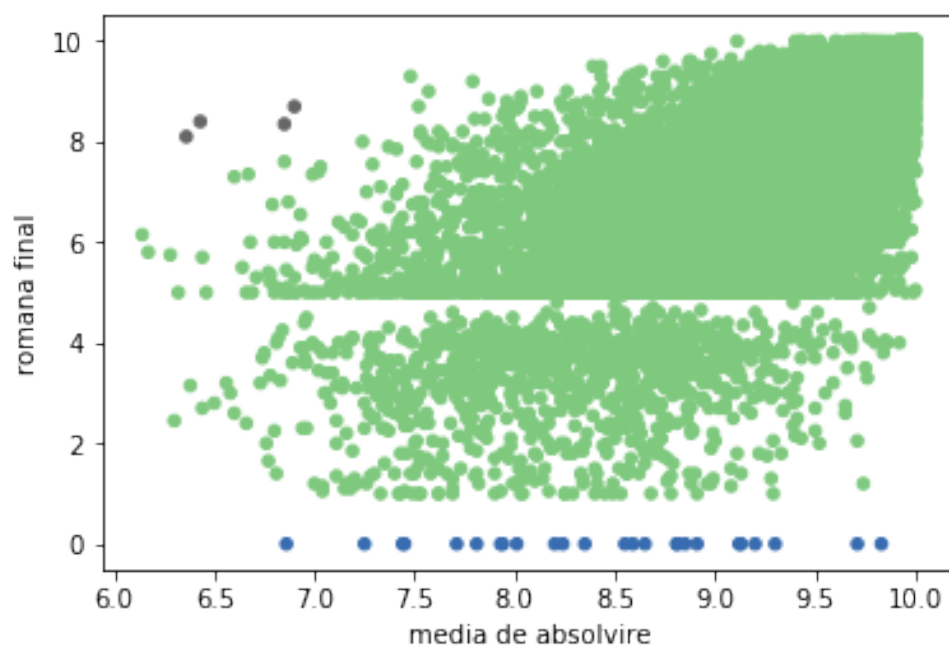
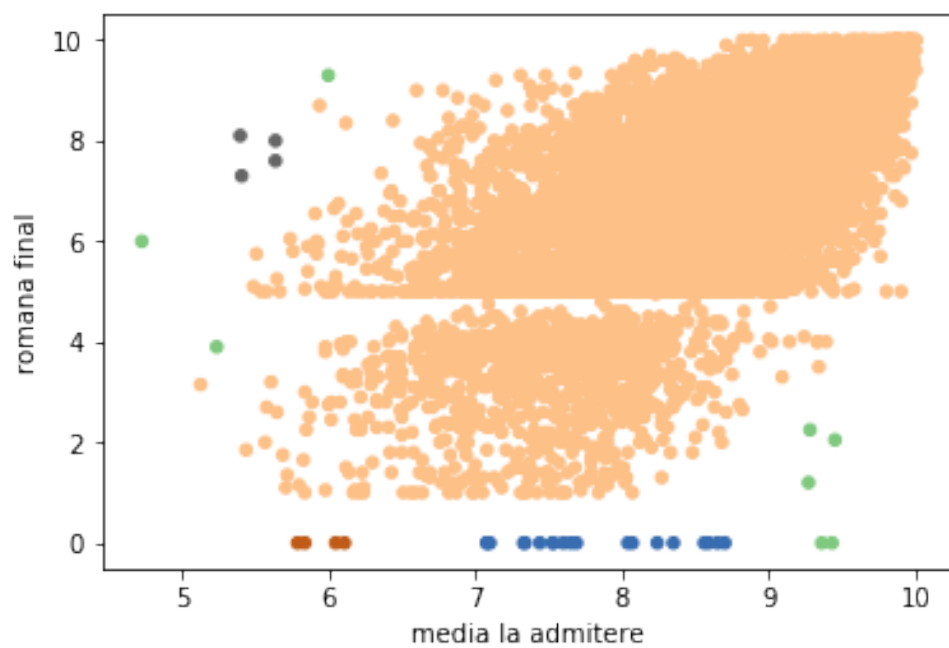


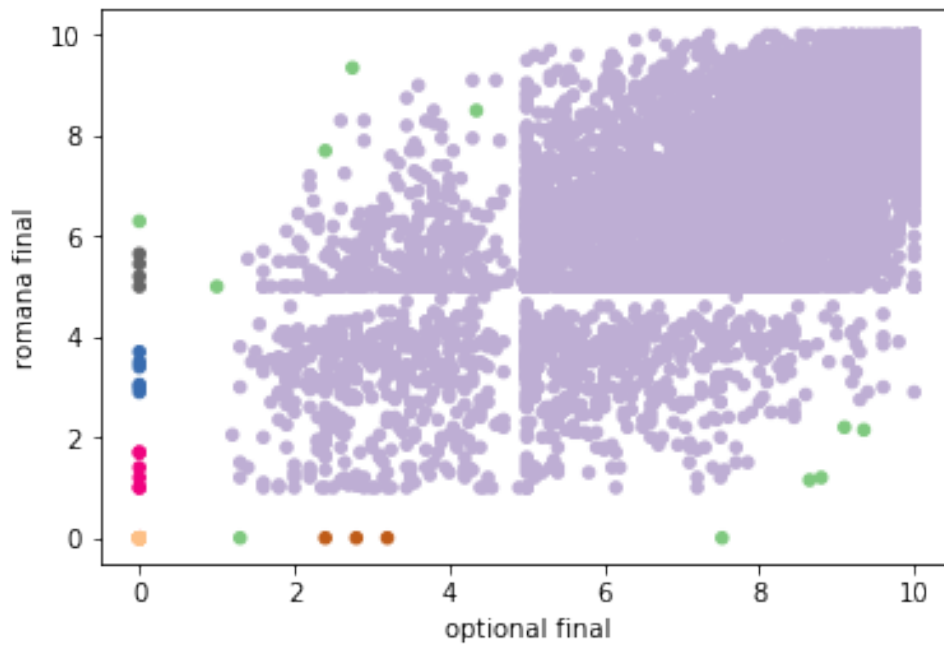
In fact, we can see that those point are outliers, so need to remove them.

```
outlier_list_index_all = set([])

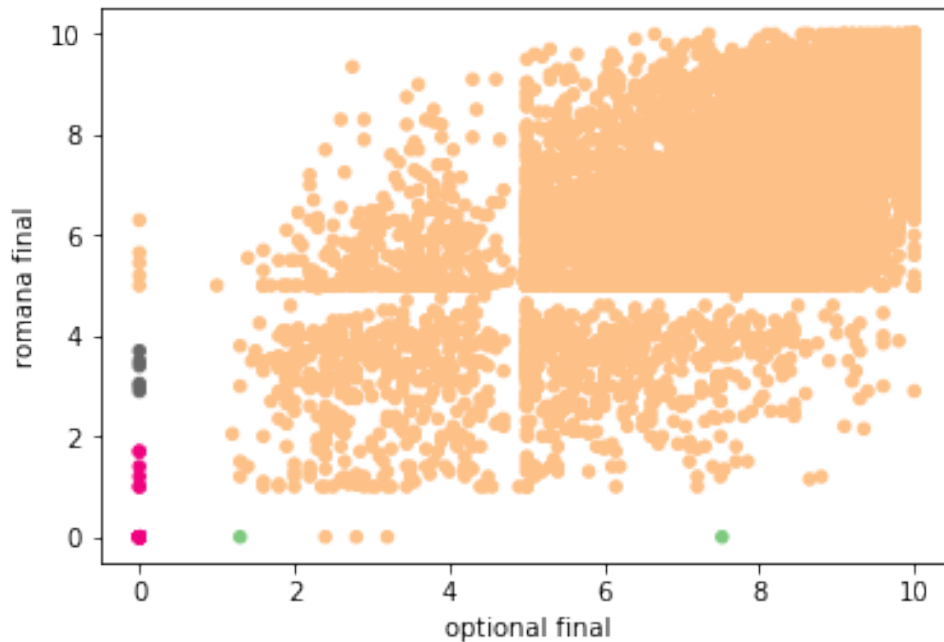
for ind, v in enumerate(clusters):
    if v != 0: # outlier
        outlier_list_index_all.add(ind)
```

Before removing outliers, we want to be sure that they can be removed - so manual checking is important (checking the plot):





In the last plot we can see that some points are not really outliers, so we will try again but using another value for the margin **eps**. A value of 1 was chosen instead of the default 0.5.



We can now remove these ones because they are outliers. After having an index of every outlier, we can remove them from the initial dataset. It is important to check outliers only in the train set because we would like to predict every kind of data in the test set even if it is an outlier.

```
#remove the outliers from the training set
```

```
X_train2 = X_train2.drop(outlier_list_index_all)
X_train2.shape
```

We get a size of (8030, 1048) instead of (8089, 1048) so a total of 59 outliers were removed. Running the Random forest regressor again yields an even lower accuracy(67%), so removing outliers didn't help.

Chapter 9

Graphs and experiments

The most important first names in predicting the final Romanian grade are as follows (the number denotes the relative importance):

```
In [34]: #most important names  
for (a,b) in feature_importances:  
    if "nume" in a:  
        print a,b
```

```
nume_maria 0.00087  
nume_andrei 0.00067  
nume_alexandru 0.00058  
nume_florin 0.00055  
nume_daniel 0.0005  
nume_andreea 0.00046  
nume_rebeca 0.00046  
nume_mihai 0.00045  
nume_ionut 0.00044  
nume_mihaela 0.00042  
nume_alexandra 0.00041  
nume_elena 0.00041  
nume_cristian 0.00039  
nume_cristina 0.00036  
nume_razvan 0.00035  
nume_... 0.00033
```

The interesting thing about these names is that these students do not get high or low grades, their grades do not have significance over the importance as shown here:

```

nume_maria avg. grade = 7.60212765957
nume_andrei avg. grade = 6.59431279621
nume_florin avg. grade = 6.00275590551
nume_alexandru avg. grade = 6.56743772242
nume_mihaela avg. grade = 7.63743455497
nume_daniel avg. grade = 5.99934640523
nume_elena avg. grade = 7.65062761506
nume_mihai avg. grade = 6.31598360656
nume_andreea avg. grade = 7.84133738602
nume_ionut avg. grade = 6.11931216931
nume_cristian avg. grade = 6.42536764706
nume_alexandra avg. grade = 7.83611111111
nume_cristina avg. grade = 7.58917525773
nume_razvan avg. grade = 6.7686440678

```

The most important counties:

```

In [35]: #most important counties
         for (a,b) in feature_importances:
             if "judet8" in a:
                 print a,b

```

```

judet8_HARGHITA 0.00104
judet8_BIHOR 0.00084
judet8_TIMIS 0.0007
judet8_DOLJ 0.00063
judet8_SALAJ 0.00062
judet8_MARAMURES 0.00058
judet8_SUCEAVA 0.00052
judet8_SATU-MARE 0.0005
judet8_ALBA 0.00049
judet8_MURES 0.00047
judet8_BACAU 0.00042
judet8_GALATI 0.00037
judet8_HUNEDOARA 0.00036
judet8_CONSTANTA 0.00033
judet8_BOTOSANI 0.00026
judet8_BRASOV 0.00025
judet8_IASI 0.00024
judet8_TELEORMAN 0.00024
judet8_VASLUI 0.00024
judet8_BISTRITA-NASAUD 0.00021
judet8_NEAMT 0.00021
judet8_GIURGIU 0.0002
judet8_VRANCEA 0.00019
judet8_ARGES 0.00017
judet8_BRAILA 0.00015
judet8_CLUJ 0.00015
judet8_BUZAU 0.00013
judet8_CARAS-SEVERIN 0.00013

```

The most important schools:

```
In [36]: #most important schools
for (a,b) in feature_importances:
    if "scoala" in a:
        print a,b

scoala_SC,CU CLS.I-VIII MARGHITA 0.00028
scoala_GIMNAZIUL "LIVIU REBREANU" TG MURES 0.00027
scoala_LICEUL TEORETIC "GEORGE MOROIANU" SACELE 0.00027
scoala_LICEUL TEORETIC "HENRI COANDA" CRAIOVA 0.00026
scoala_SCOALA CU CLASELE I-VIII "PETRU RARES" PROBOTA 0.00026
scoala_SCOALA CU CLASELE I-VIII NR 1 BOSANCI 0.00025
scoala_SCOALA CU CLASELE I-VIII NR. 2 HUSI 0.00024
scoala_SCOALA CU CLASELE I-VIII "GEORGE BARITIU" TURDA 0.00022
scoala_SCOALA CU CLASELE I-VIII LIPIA COMUNA MEREI 0.00022
scoala_SCOALA CU CLASELE I-VIII MOLDOVENESTI 0.00022
scoala_SCOALA CU CLASELE I-VIII "AVRAM IANCU" ALBA IULIA 0.00019
scoala_SCOALA CU CLASELE I-VIII "JEAN BART" CONSTANTA 0.00019
scoala_SCOALA CU CLASELE I-VIII NR 8 "MIHAIL KOGALNICEANU" DORHOI 0.00019
scoala_SCOALA DE ARTE SI MESERII MARGINEA 0.00019
scoala_SC,CU CLS.I-VIII CONSTANTIN SERBAN ALESU 0.00017
scoala_LICEUL CU PROGRAM SPORTIV SEBES 0.00016
scoala_SC,CU CLS.I-VIII ARANY JANOS SALONTA 0.00016
scoala_SCOALA CU CLASELE I - VIII FILIASI 0.00016
scoala_SCOALA GENERALA "MIRON CRISTEA" TOPLITA 0.00016
scoala_SCOALA GENERALA "TAMAS ARON" LUDINT 0.00016
```

And the most important features when using LabelEncoder which does not create a new feature from every categorical value:

	importance
media la admitere	0.402439
disciplina profil final	0.121583
optional final	0.114280
nota la limba romana	0.048751
limba moderna nota	0.045490
media de absolvire	0.026489
romana oral	0.022494
scoala	0.020394
unitatea de invatamant	0.018966
nota la matematica	0.018673
limba materna8	0.018140
optional8 nota	0.016102
media teze nationale	0.016068
limba materna	0.012797
specializare	0.009672
judet8	0.009073

Chapter 10

Related work and documentation

[3] uses Bayesian Networks to predict the grades for students. Before this, it uses plots to see correlation between data. I've used the same approaches to identify correlations namely pair plots, plotting the deviation curve of the final grade. The only difference is that this related work has more data to work with including absences from high school, father education, mother education etc. I don't have these information so in my case developing a good model is harder. I've tried using Bayesian Networks, but without success. Also, my dataset is rather large and the model training phase took too long.

The kernels working on [10] use RandomForestRegressor, ExtraTreesRegressor and Gradient-BoostingRegressor to predict the final grades. The thing is that even this dataset has features such as age of student, if the student aims for higher education etc. I don't have these data so my accuracy suffers. These kernel achieved an accuracy score of around 77%. My data set is totally different but this shows that predicting grades for students is a tough process because during their study years, they can stop learning, they can start learning harder or everything can happen to them....This problem is not an objective prediction problem like predicting if an image shows a dog or a cat.

There are two kernels in [2] and both contain a lot of information analysis and feature engineering. The one thing which I've realized is that feature engineering is the most important aspect in this problem. It is more important than choosing good algorithms. Because I am using a dataset with many features, some patterns may lie underneath which are hard to observe but can increase the accuracy by 1 or 2 percents. This was the case when I had the idea to generate another feature that contains the gender by analyzing the first names of students. Also by including the first names rather than the full names I achieved a small but significant accuracy boost. Also, another good idea was to make a numerical value from the grades got at the language exam because in this way this single numerical value has more significance than the letters.

There is one thing in common in these related kernels: they don't throw out any information which may seem as unnecessary. They try to modify these information or where information is missing, they try to replace missing values with meaningful ones. For example, one of the kernels gave high importance to the feature which showed if a student wanted to continue with higher education or not. In the end, it resulted that students who **want** to score higher on the exam(to enter university) will really score higher eventually. Also this same kernel concluded that students who are not in a relationship score higher. Or students who go out more, score less. These features may seem unimportant at first, but they carry important information([12]).

Bibliography

- [1] Admitere 2008-grades. <http://web.archive.org/web/20081201094019/http://admitere.edu.ro/2008/staticRepI/j/>.
- [2] Bagrut grades in israeli high schools (2013-2016). <https://www.kaggle.com/emachlev/bagrut-israel>.
- [3] Bayesian linear regression in python: Using machine learning to predict student grades. <https://towardsdatascience.com/bayesian-linear-regression-in-python-using-machine-learning-to-predict-student-grades-part-1-7d0ad817fca5>.
- [4] Export html table to csv. <https://jsfiddle.net/gengns/j1jm2tjx/>.
- [5] Gender api. <https://gender-api.com/>.
- [6] Heart disease uci. <https://www.kaggle.com/ronitf/heart-disease-uci>.
- [7] Kaggle example. <https://www.kaggle.com/triomni/heart-disease-prediction-by-a-newbie-to-newbies>.
- [8] Official Python website. <https://www.python.org/downloads/>.
- [9] Plot the importances. <https://stackoverflow.com/questions/17057139/how-to-find-key-trees-features-from-a-trained-random-forest>.
- [10] Predict the final grade of portugese high school students. <https://www.kaggle.com/dipam7/student-grade-prediction>.
- [11] Rolisoft : Bacalaureat-data. <https://github.com/RoliSoft/Bacalaureat-Data>.
- [12] Student grades kernel. <https://www.kaggle.com/dipam7/introduction-to-eda-and-machine-learning>.
- [13] Student performance data set. <https://archive.ics.uci.edu/ml/datasets/student+performance>.

