

# Wine Quality Prediction

Francesco Di Salvo

Politecnico di Torino

Student id: s282417

s282414@studenti.polito.it

**Abstract**—In this report we propose a possible approach for predicting the quality score associated with a wine review. In particular, this approach consists on encoding the categorical features, managing outliers, null values, duplicated and a textual field. Once completed the preprocessing step, two regression models have been compared and they both outperformed a naive solution, obtaining a satisfactory result.

## I. PROBLEM OVERVIEW

The proposed competition is a *Regression Problem* based on a collection of 150 930 *wine reviews*, divided into *development set* (120 744) and *evaluation set* (35 716). The aim of the competition is to predict the *quality* defined for a given review, by considering 8 features defined in Table 1.

All the categorical features present a wide cardinality domain, in a range between 49 (country) and 27 801 (designation). Then, both development and evaluation set contain a considerable number of *null values*. These two aspects must be taken in consideration for the processing phase. In addition to these categorical features, there is the *description* field, that contains the textual description of the review. We discovered that there are just 85 005 *unique descriptions* over a total number of 120 744. Thanks to this insight, we discovered that the development set contains 35 739 duplicates.

Another key aspect regards the predictor (*quality*). It is normally distributed on the interval [0,100] (Figure 1), whereas the median (50<sup>th</sup> percentile) lays on 46. In order to detect the outliers, it has been used the 1.5(*IQR*) *Rule*, that allows to define the minimum and maximum thresholds beyond which we consider the data points as outliers.

$$\begin{aligned} \min &= Q_1 - 1.5(Q_3 - Q_1) = 12.5 \\ \max &= Q_3 + 1.5(Q_3 - Q_1) = 80.5 \end{aligned} \quad (1)$$

where  $Q_1$  and  $Q_3$  are the first and the third quantile, respectively. This technique allowed to discover 4 972 outliers. In order to verify if these values must be considered as outliers or not, we inspected two random descriptions with quality 0 and 100 and this is what we obtained :

- 0 : "Clean as anyone should reasonably expect given the almost unheard-of price. [...] Drink it now; cook with it; make a sangria with it. Just enjoy it."
- 100 : "The nose on this single-vineyard wine from a strong, often overlooked appellation is tight and minerally before showing a slightly tropical kiwi element [...]."

The difference is not noticeable at all, so we assumed that these "rare" values are probably not reflecting the real quality.

TABLE I: Summary

Feature	Development set		Evaluation set	
	Cardinality	Null values	Cardinality	Null values
country	49	0.01%	40	0%
designation	27 801	30.24%	11 952	30.53%
province	445	0.01%	342	0%
region_1	1 207	16.57%	949	16.74%
region_2	19	59.64%	19	59.53%
variety	603	0%	457	0%
winery	14 105	0%	8 962	0%
description	-	0%	-	0%
quality	-	0%	-	-

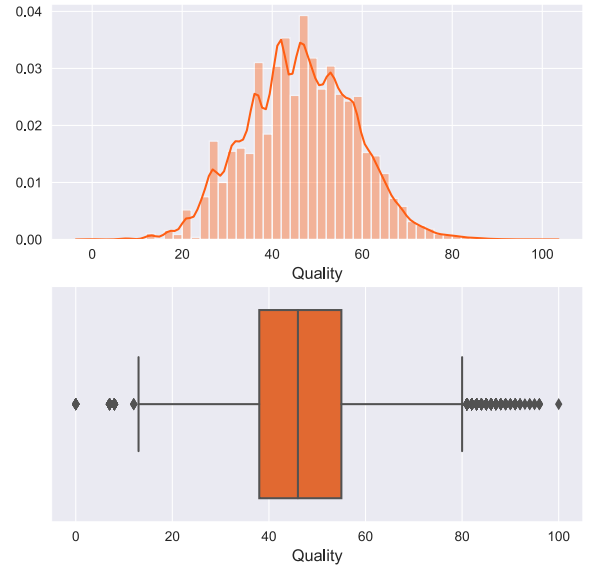


Fig. 1: Quality distributions and outliers detection

It is also interesting that the quality associated to most categorical values has a considerable variance. For example in Figure 2 it is possible to observe the top 15 countries (based on the median quality) and most of them covers a wide range, even if some of them are stretched by the outliers or they have just one review (US\_France).

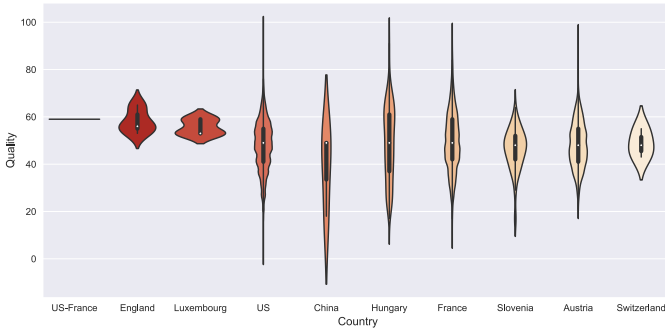


Fig. 2: Top 10 countries by median quality

In addition to the geographical features, there are other two features strongly related to the wine itself, and they are *winery* and *designation*. Even if the latter has around 30% of null values, it may be still be relevant for the prediction. On the flip side, this strong representativeness is related to a huge cardinality domain, that will affect the overall performances.

## II. PROPOSED APPROACH

### A. Data preprocessing

It has been demonstrated that duplicated values on a regression problem may bias the estimated coefficient and standard errors [1]. So, we decided to drop all the duplicated values on our dataset (35 716 samples).

Then, we encoded the categorical features. This process is extremely important, especially when it needs to encode considerable categorical domains. So, after some preliminary analysis and thanks to the Cedric and Seger's contribute [2], we decided to use the *OneHotEncoding* technique, even if it does involve a remarkable slowdown in performances. This is particularly valuable when there is not a strong "ordinal" relationship among the values, because it creates a binary column for each category, as follows :

country quality				country_France		country_US	quality
0	France	45.0	➡	0	1	0	45.0
1	US	31.0		1	0	1	31.0
2	US	35.0		2	0	1	35.0

Then, as it can be clearly seen in Table 1 that *region\_2* contains around 60% of null values, we decided to drop it. Then, since *country* and *province* contains just the 0.004% , they have been replaced with the most common attribute on the respective columns. The last two columns with null values are *designation* and *region\_1*. They both have a reasonable amount of null values, and because of that, they have been replaced with the sentinel value "other", otherwise they would have given too much importance to the most frequent ones.

Since before we assumed that some target values should be considered as outliers, we analyzed the behaviour of the model. So, we tested the model with outliers and also without them.

We can consider them as hyperparameters for the preprocessing step:

- $\alpha$  : removing outliers
- $\beta$  : not removing outliers

We will discuss it more in detail on the *hyperparameter tuning* section. Lastly, the *description* feature required a refined preprocessing, in fact we removed it from the main data structure and we manipulated it separately. In particular we:

- removed non alphabetic characters,
- converted each word in *lowercase*,
- removed *stop-words* (most common words),
- extracted the root from each word (*stemming*).

Then, once we obtained the cleaned descriptions, we used *CounterVectorizer* and *TfidfTransformer* provided by Scikit-Learn. The former converts the descriptions to a matrix of token counts used by *TfidfTransformer* for calculating the *text frequency - inverse document frequency* (tf-idf):

$$TF-IDF(t) = freq(t, d) \times \log \frac{m}{freq(t, D)} \quad (2)$$

where  $t$  is the token,  $d$  is the current document, and  $D$  is the collection of  $m$  documents. The idea of the TF-IDF is to stand out tokens that occurs frequently in a single document but rarely on the entire collections of documents. For the Counter Vectorizer we considered the 4-grams (subsequence of 4 tokens) that are present just on less than 50% of the documents (presumably a positive or negative half).

### B. Model selection

The following algorithms have been tested:

- *Linear Regression* : it is a model that assumes a linear relationship among its features. The prediction is made by computing a weighted sum of the input features and the intercept.
- *Stochastic Gradient Descent Regressor (SGD)* : it picks a random instance instance in the training set at every step and compute the gradients based only on that single instance. Due to its stochasticity, the cost function bounces up and down and slowly converges to the minimum. This randomness is good for escaping from the local optima, but on the other hand it cannot settle at the minimum. This issue can be partially solved by considering the behaviour of the well known simulated annealing algorithm.

For both regressor models, we found out the best configurations through the grid search approach, as it is possible to see it in the next section.

### C. Hyperparameters tuning

As we mentioned before, we tuned two different sets of hyperparameters, one for the preprocessing task ( $\alpha$ ,  $\beta$ ) and one for the prediction models (Table 2).

Firstly, we analyzed the performances of ( $\alpha$ ,  $\beta$ ) on both regressors with their naive configuration by means of the  $R^2$  score. Once we found the apparently best preprocessing, we can finally tune the models. The best configuration can be

found by using the *Grid Search Cross Validation* technique that works as follows :

- 1) Generates all the possible candidates that needed to be trained.
- 2) Train and test each configuration with a Cross Validation approach ( $k = 3$  in our case).
- 3) Select the best configuration.

TABLE II: Hyperparameters Tuning

Model	Hyperparameters	Values
LinearRegression	<i>fit_intercept</i>	{True, False}
	<i>normalize</i>	{True, False}
SGDRegressor	<i>loss</i>	{'squared_loss', 'huber'}
	<i>penalty</i>	{'l1', 'l2', None}
	<i>alpha</i>	{1e-5, 1e-4, ..., 0.1, 1}
	<i>eta0</i>	{0.01, 0.1}

### III. RESULTS

The comparison among the two preprocessing proposed ( $\alpha$  = with outliers,  $\beta$  = without outliers) produced the results reported in Figure 3. It seems that removing the outliers mostly affected the SGDRegressor’s performances, whereas the Linear Regression is stable in both cases. So, we have choosen to tune both Linear Regression and SGD without removing the outliers  $\beta$ . The best configuration for Linear Regression was {fit\_intercept=True, normalize=False} and it obtained an average  $R^2$  of 0.7444. Then, the best one for SGD was {loss='squared\_loss', penalty='l1', alpha=1e-05, eta0=0.1} with 0.745. Actually, the Linear Regression outperformed SGD on the public score with 0.879 against 0.856. These results have been compared with a naive solution defined by the encoding of the cateogorical features and without considering duplicated, outliers and the description field. This one obtained 0.828 on the public score.

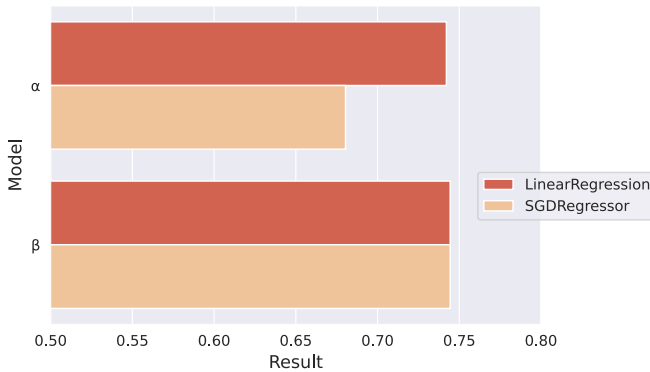


Fig. 3: Preprocessing evaluation

### IV. DISCUSSION

The current approach obtains an overall satisfactory result, outperforming the proposed baseline. In particular, removing the duplicates and managing the textual field has strongly influenced the performances. Further considerations can be

made by studying how the model behaves to a filter for each cardinality, that tries to iteratively remove the least frequent categorical attributes for wide categorical domains, such as designation and winery. Lastly, it would helpful to reduce the overall size of the final data structure with the fairly well known technique : *Principal Component Analysis*. For the proposed approach it was not possible to test it, due to the lack of strong performances of the used workstation.

### REFERENCES

- [1] Sarracino, Mikucka, “Bias and efficiency loss in regression estimates due to duplicated observations: a Monte Carlo simulation”
- [2] Seger, Cedric , “An investigation of categorical variable encoding techniques in machine learning: binary versus one-hot and feature hashing”