

DSL Report

Wine Quality Regression

Armando La Rocca
Politecnico di Torino
s279401
s279401@studenti.polito.it

Abstract—This report contains a possible approach to the regression task on the wine quality dataset. Firstly, the problem is briefly explained. In the second section the proposed approach is presented and described. The last two sections contains a presentation of the results obtained and some considerations also on possible improvements.

I. PROBLEM OVERVIEW

The following report present a possible solution to the competition related to the estimation of wine quality. Each wine is characterized by 8 categorical features: four related to the geographical provenience and three related to the winery, variety and designation. The last feature is a textual description of the wine. The problem is a regression task and the aim is to predict the quality of the wine that is represented by a value from 0 to 100.

The proposed dataset is divided in 2 splits:

- *development set* containing 120738 labelled samples.
- *evaluation set* containing 30186 samples.

First of all the development set is analyzed in order to extract useful information.

In the dataset are encoded a lot of wines of different types and provenience. However it is difficult that a dataset represents all the existing (or at least a large part of the) wines produced in the world. So, to increase the generality, data should be encoded in a way that allow to work also with categories that are not present in the development dataset or with missing values. Missing values, in fact, are quite frequent in the dataset, above all in certain features as shown in Fig 1.

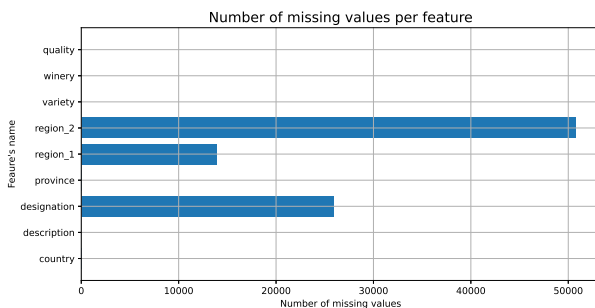


Fig. 1. Number of Missing values for each feature.

The most interesting feature is probably the description. For each wine is provided a textual description of the wine,

however is not easy to analyze this kind of reviews. First of all because the descriptions are very similar even with a great difference of quality. An example is reported below :

'the nose on this single-vineyard wine from a strong, often overlooked appellation is tight and minerally before showing a slightly tropical kiwi element. brightly acidic on the lively palate, flavors range from key lime and meyer lemon to pear skins and apple flesh.' → *quality* = 100

'clean as anyone should reasonably expect given the almost unheard-of price. it offers solid cherry and blackberry flavors, while the mouthfeel is round and totally inoffensive. drink it now; cook with it; make a sangria with it. just enjoy it.' → *quality* = 0

In fact, between the two presented wines there is an high difference of quality but from the descriptions is not easy to catch that.

The choice of the encoding and a good management of the description features are crucial steps to obtain good performances.

II. PROPOSED APPROACH

A. Pre-processing

First of all, a data cleaning step is performed on the development set. Duplicates are founded and dropped out from the dataset because represent not useful and redundant information. After this operation the development set accounts 85028 samples.

Then all the feature values are transformed in lower case to avoid mapping mistakes generated by different cases letters. Another transformation is related to the designation category in which same wine categories appear in different languages. In Figure 2 can be seen as the feature *reserve* appears also as *riserva* and *reservé*. To exploit that the wine category is the same all the different names are mapped with the same label. In Table 1 are presented the transformations made to solve this problem.

The dataset is also characterized by missing values, above all in three features as represented in Figure 1. Considering the real case, it is possible that a wine is not characterized by a specific designation or that the specific region provenience is not provided. However, the quality evaluation should be possible also without this data. For this reason the missing

TABLE I
DESIGNATION MAPPING

New designation	Old designation
<i>reserve</i>	<i>reserva, réserve</i>
<i>grand reserve</i>	<i>gran reserva</i>
<i>special reserve</i>	<i>reserva especial</i>
<i>red wine</i>	<i>red, rosso, tinto</i>
<i>rosé</i>	<i>rosé of, rosato, rosado</i>

values of the features *region_1*, *region_2* and designation are filled with the string *other*.

In contrast, the samples with missing values in *country* features are dropped out (just 3 samples not distinguishable in Figure1).

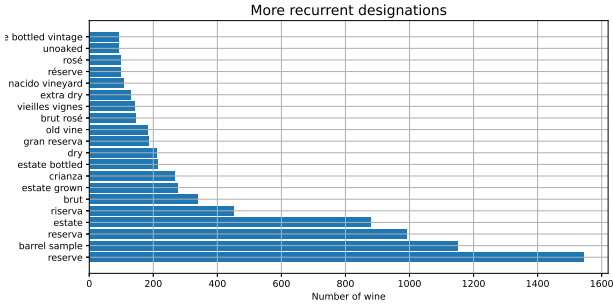


Fig. 2. Number of wines per designation

As said before, all the features are categorical except the description that is a textual review. Categorical features are transformed with the one-hot technique. This method produce for each feature a number of elements equal to the number of different values assumed by the features. Then put to 1 the bit associated to the value assumed by the feature for that sample and to 0 the others. Moreover, this approach allows to encode values that are not in the development set putting all the elements to 0. This is a good point for this analysis because also new or unknown wines can be encoded. However this encoding increases consistently the number of features of the dataset even if the data is quite sparse.

The preprocessing step for the description feature is more complex. First of all, each description is tokenized and the tokens containing numbers or punctuation are removed. Then, the tokens are lemmatized in order to avoid repetitions caused by inflected forms. Finally, all the obtained tokens are transformed to lower case and stop words are removed. This cleaning step is performed to exclude tokens that are not useful as numbers or parts of the text without specific meaning.

At the end to extract features from the obtained tokenized descriptions are used two different approaches :

- **Approach 1 (A1)** The descriptions are divided in three groups in according to their quality. From each group is created a vocabulary excluding the tokens that appear over 70% of documents and less that 5%. Then are selected only the words that are not in common with all the obtained vocabularies. This approach aim to find the

words that most characterize wines belonging to a quality ranking of high, medium and low quality (as represented in Figure 3). At the end, features are constructed considering the count of the tokens of the obtained vocabulary in each sample description.

- **Approach 2 (A2)** The tokens are filtered excluding the words that appears in more than 50% of descriptions samples and less than in 10 reviews. Then features are constructed considering the term frequency of each token in the descriptions.

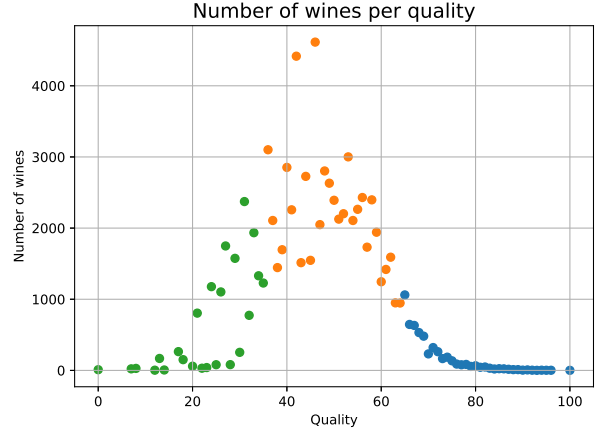


Fig. 3. Number of wines per quality. In green are represented wines associated with low quality, in orange medium and in blue high.

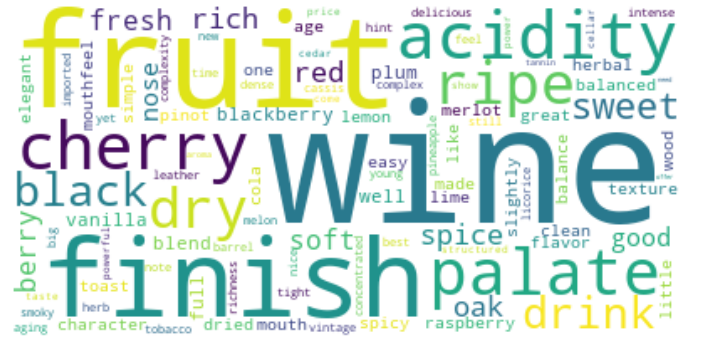


Fig. 4. Most recurrent words with approach A1

B. Model selection

For both the approaches proposed in the previous paragraph the number of features is high and data are quite sparse. This aspect can affect some algorithms for curse of dimensionality. For this reasons the models chosen for this tasks are:

- **Random Forest [1]:** this algorithm generates and train a given number of decision trees using the bagging method. Each tree works using one feature at time chosen between a subset of the original features. So the algorithm should be less affected by the high dimension of the data.

TABLE III
MULTI LAYER PERCEPTRON GRIDSEARCH

Hyperparameters		Features approaches		
nn	lr	ND	AI	A2
100	0.01	0.7473	0.6776	0.7551
100	0.001	0.7369	0.7239	0.7504
500	0.01	0.7499	0.6641	0.7575
500	0.001	0.7417	0.7300	0.7524

are two main aspects that can be better analyzed to try improve the results.

- In the preprocessing two possible description's encoding are proposed. Text data are quite complex data and different other approaches can be tried. A first, very basic, possibility is to investigate the range of word excluded for high or low frequency in the sample descriptions. Also analyze group of words, as bigram and trigram, is an interesting ways. Word embeddings, actually, is considered one of the main technique to deal with text data and can be interesting understand if it can improve the results.
- Neural networks are characterized by a lot of different parameters that can affect the performances in a significant way. For this reason, a better study on the hyperparameters as other solvers(i.e. *adam*) or loss should be done. Also try more complex and deep models can be a way to improve the results.

REFERENCES

- [1] Scikit-learn Multi layer regressor
https://scikit-learn.org/stable/modules/generated/sklearn.neural_network.MLPRegressor.html
- [2] Scikit-learn Random Forest regressor
<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
- [3] Natural language toolkit
<https://www.nltk.org>