# Wine Quality from content

René Angeles

March 16, 2020

*He who knows nothing, loves nothing... But he who understands also loves, notices, sees ... The more knowledge is inherent in a thing, the greater the love. Anyone who imagines that all fruits ripen at the same time as the strawberries knows nothing about grapes.*

—Paracelsus

## 1 Intro

We succinctly report findings for the quality prediction of wines using [1]. Our goals correspond to the data challenge at [2].

## 2 Exploratory data analysis

The EDA performed (see `EDA.ipynb` or `EDA.html`) shows that

- Acidity predictors (fixed, volatile & citic), residual sugar, chlorides, sulfur dioxide (total & free), density , sulphates all have a heavy right tail. Naively, these tails are outliers. A log/log1p transformation bring these closer to a bell shape

- The distributions for the red and white wine predictor are either shifted or different.

- The correlations coefficients larger/smaller than 0.5-0.8. Hence the used of PLS methods is justified. Correlations patterns are different between white/red wines.

Thus motivated, training was performed

- With and without log transformed heavy tailed predictors. It was found that in all cases log transformation does not improves MAE.

- With and without potential outliers where removed using a $\alpha \times IQR$ rule with $\alpha \in \{1.5, 1.8, 2, 2.2\}$. Yet, the removal of such potential outlier was never found to improve MAE.

- On the red and white wine datasets separately and jointly.

No missing values/duplicates were found. To-do: Multivariate distributions could shed light relation between predictors. It is obvious that pH cannot be independent of acidity predictors, and similarly for other variables. Domain knowledge could actually help to tell how are features related.

# 3 Modeling

PLS is expected to perform better when predictors are center and scale Ref. [3, 4]. However, in our case we performed analysis with and without these transformation, finding that these do not increase the MAE score.

Respectively, the files `Baseline.R`, `Pls_polr.R`, `Plsr_ensembles.R` and `GLMs_ensemble.R` are standalone models containing the feature engineering options presented above.

| Main models reported | | |
|---|---|---|
| Model. | Train MAE | Test MAE |
| Baseline: plsR, 5fold-tuned, independent models for red/white datasets. Predictions rounded. | 0.4363 | 0.5026 |
| Pls-polr, 5fold-tuned separately for red/white sets. | 0.5074 | 0.5181 |
| Rounded prediction from a plsr stack of 3 plsr models trained into 4 separate train-sets (5fold tuned). Combined datasets. | 0.514522 | 0.5285054 |
| Rounded PLSR stack of 3 PLSR models trained with resampling with replacement, 5fold tuned. Combined datasets. | 0.5074 | 0.5181 |
| pls-polr blend of 3 models (pls-polr 5fold-tuned, pls-gaussian, psl). Only white wine was considered. | 0.5080 | 0.5556 |

# 4 Other insights:

- The PSLR baseline stayed undefeated but it is worth noting the stacked ensemble has a very similar MAE for training and validation sets. This suggest the latter has low variance.

Figure 1: Predictor importance of baseline

- Due to the limited time, I didn't tune the hyperparameters for the models appearing in the blended ensemble, this is possibly why it performed poorly.

- When performing kfold tuning, the number of predictors kept varied between 2 and 6. This means that between 10 to 5 predictors suffer from collinearity problems.

- Our models, work comparatively better when applied separately to the white and red datasets. The stack model was fix on the combined dataset, MAE would reduce addressing white and red wines separately.

# 5    Predictor importance

A detailed analysis of several statistics using `plsrglm` is possible but beyond the scope of this work. Though, for completeness, let us check at least the largest coefficients of the sucesfull baseline model, Fig. 1, which suggest that volatile acidity, density and clorides have the largest coefficient in the regresion.

# References

[1] Paulo Cortez *https://archive.ics.uci.edu/ml/datasets/Wine+Quality* .

[2] Analytic Flavour System *https://www.gastrograph.com/blogs/gastronexus/interviewing-data-science-interns.html* .

[3] M. Kuhn K. Johnson *Applied Predictive Modeling* .

[4] Kee Siong Ng *A Simple Explanation of Partial Least Squares* .