
CAPSTONE PROJECT

— MODEL TO PREDICT QUALITY OF —
WINE

INTRODUCTION

- Wine quality depends on the vinification process, geographical origin of grapes and varietal composition of grape must.
- The assessment of wine quality and authenticity can be considered as a social demand since consumers require high quality products.
- The wine chemistry and several features influencing the complexity of chemical composition turns its quality analysis highly challenging.
- Effective tools for wine classification and detection of adulterants is required.

OBJECTIVE

- The aim is to create a prediction model in Python language for the quality analysis of wine
- A large dataset is considered which shows the quality of wine for several chemical components such as pH,density,alcohol etc.
- Multiple classification algorithms will be applied on the dataset after data preparation.
- The accuracy will be analysed for each algorithm and optimised to choose a final model for the quality prediction of wine.

DATA GATHERING

IMPORTING LIBRARIES

- The initial step before importing dataset is to import various libraries in Python.
- To name a few, pandas is used in converting csv format to dataframe and for exploratory data analysis
- Numpy is used for mathematical operations.
- Matplotlib and seaborn is used for data visualisation.
- Scikit learn features various algorithms, data preprocessing and model evaluation matrices.

READ DATASET

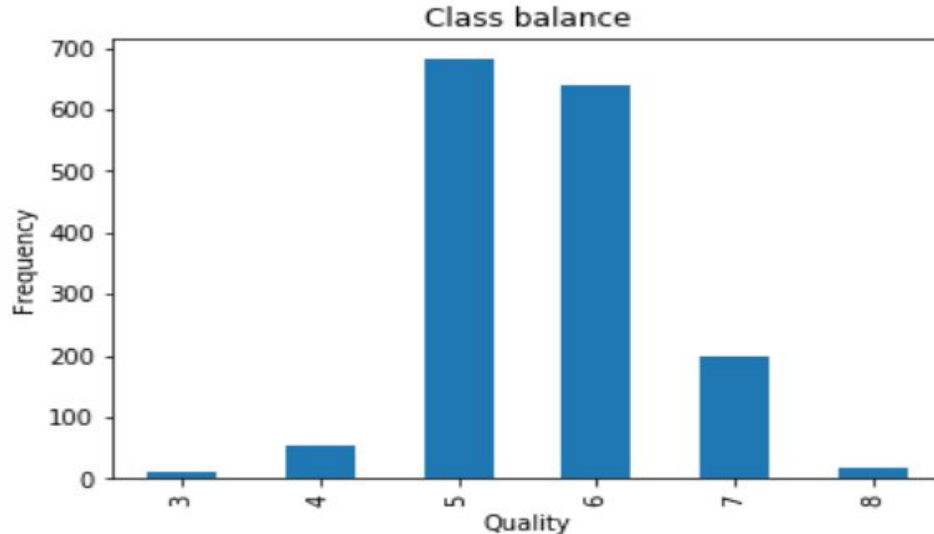
- The dataset which is in csv format is extracted to a dataframe 'wine'.
- The independent variables are 'fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar', 'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density', 'pH', 'sulphates' and 'alcohol'. All are continuous variables.
- The target variable is 'quality', which is a categorical variable and rates the wine quality on a scale of 1 to 10.

EXPLORATORY DATA ANALYSIS

UNDERSTANDING DATA

- Checking null values : To find missing values which should be treated. Our dataset has no null values.
- Dimension : The dataset has total 1599 rows and 12 columns.
- Statistical summary : describe() is used to analyse mean , median , count and five-point summary. Variation in median and mean was observed which hints the presence of outliers.
- Target variable is extracted to 'Y' and the predictor variables to 'X'

- Heat map: To visualise correlation between variables
- Visualising the balance of classes : A bar graph was plotted to see the contribution towards each output classes.



UNDERSAMPLING

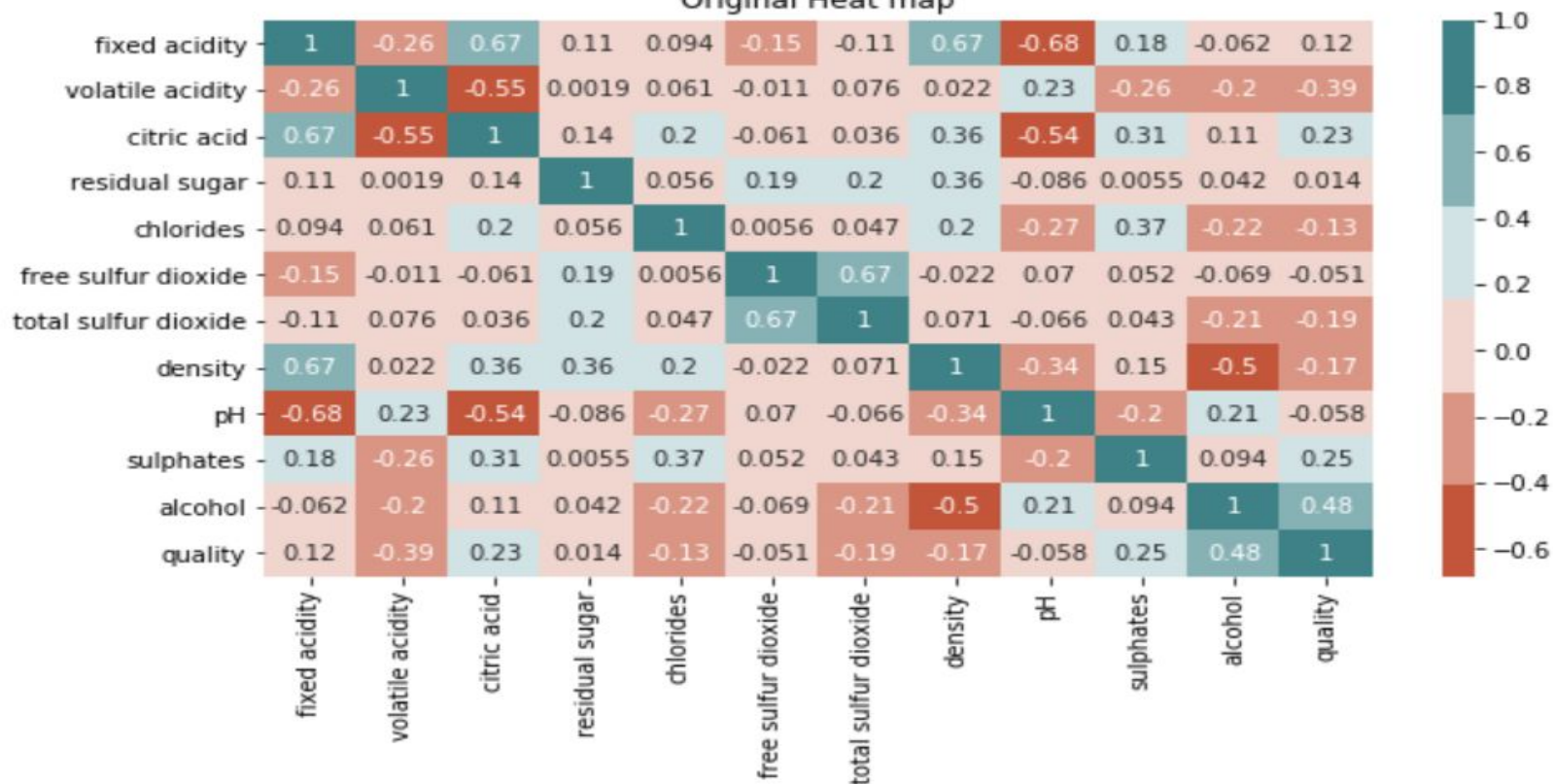
- Here, we observe data imbalance. A model fed with such a data could give biased decision and so less accuracy in prediction. It also affects the correlation between features. To avoid this, we do resampling.
- Re-sampling: Undersampling and oversampling.
 1. Undersampling : Removes some observations from majority class.
 2. Oversampling : Adds data to generate output from minority class-SMOTE

- From the bar graph 'Class balance', majority classes are 5,6 and minority classes are 3,4,7,8.
- SMOTE : To oversample the data in minority classes to the class with most number. Alternatively, each class could be assigned with required number of samples. This resulted in strengthening the relationship between variables, ie. , an increase in correlation coefficients.
- This transformation should be performed on training data only.

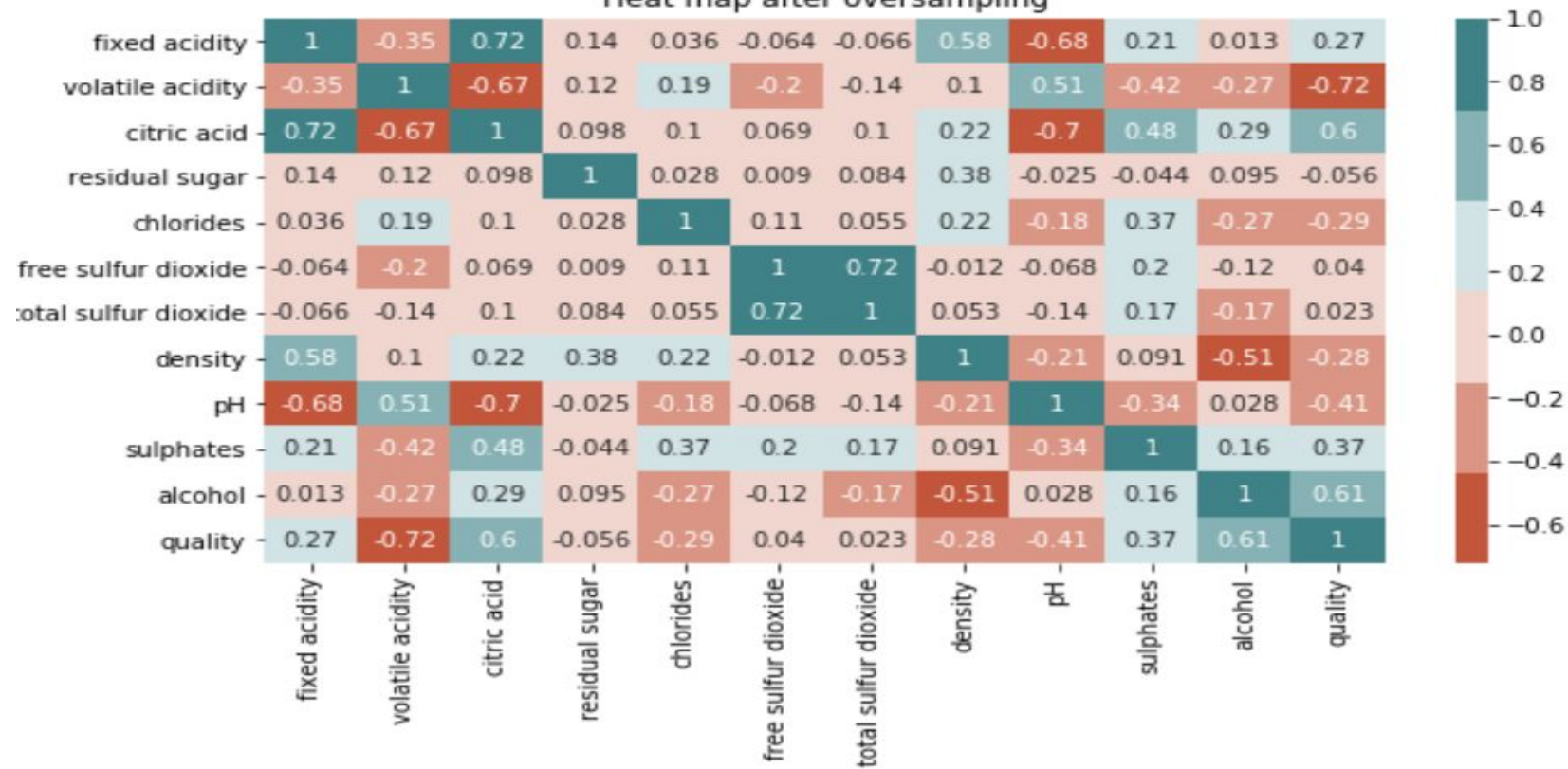
SPLITTING DATASET

- First the target variable and predictor variables are separated.
- The entire dataset is splitted into training set and testing set-
x_train,y_train,x_test,y_test at 25% split.
- Random_state : The parameter that is used to reproduce the result.

Original Heat map



Heat map after oversampling



FEATURE SELECTION

- Heatmap : The heatmap is plotted for the dataset to find the strength of correlation between variables. Some analysis:
 1. Relatively strong positive correlation : 'fixed acidity' and 'citric acid', 'fixed acidity' and 'density', 'free sulfur dioxide' and 'total sulfur dioxide'.
 2. Strong negative correlation : 'pH' and 'fixed acidity', 'pH' and 'citric acid'
- 'Total sulfur dioxide' has strong positive correlation with 'free sulfur dioxide', and its contribution towards 'quality' variable is very less. So we drop that column.

DETECTING OUTLIERS

- There are various methods to detect and remove outliers. These are:
 1. Visualising through boxplot
 2. Visualising through scatterplot
 3. Z-score method
 4. Interquartile Range method(IQR)
- Here both box plot and IQR method is chosen , where upper and lower threshold is found out and any data points not in the range are removed.

BUILDING MACHINE LEARNING MODEL

DECISION TREE ALGORITHM

- Different models were built for different values of max_depth and splitting criteria('gini' and 'entropy').
- For each model, training score and testing score are evaluated to choose the better model out of them, to find a balance between better accuracy and reducing overfitting problem.
- The best out of the lot is decision tree with max_depth=4 and criteria=gini. Accuracy=0.61
- For each of the algorithms, SMOTE was avoided since it lowered accuracy.

NAIVE BAYES ALGORITHM

- It is a supervised learning algorithm which applies Bayes' Theorem with 'naive' assumptions such as every feature is independent and every feature has equal contribution to the outcome.
- The basic logic is the comparison of probabilities of a data point falling into one of the output classes.
- Here, GaussianNB() classifier is chosen where the likelihood of features is assumed to be Gaussian. GaussianNB() showed better accuracy and F1 score than ComplementNB() ,which is used for imbalanced datasets.

ENSEMBLE ALGORITHMS

- Various Ensemble methods are chosen for the same processed data.
- `AdaBoostClassifier()` : This uses Adaptive boosting technique, which assigns higher weight on weak learners and lower weights on strong learners. Accuracy=0.5675 and F1 score= 0.524
- `GradientBoostingClassifier()` : This boosting method supports both binary and multi-class classification. This works by optimising loss functions. Accuracy=0.615 and F1 score= 0.611

- RandomForestClassifier() : Uses bagging method and uses different decision trees having random subset of dataset(sampling with replacement).Accuracy= 0.64, F1 score=0.62
- It is clear that Random Forest showed better result.

RANDOM FOREST

- Hyperparameter tuning for Random Forest is done using GridSearchCV() function. It gave parameters to be : criterion='gini', max_depth= 6, max_features= 'auto'.
- These parameters are fit into the model.
- OOB error rate: It is the misclassification rate and is used to find the optimum value of n_estimators (no of decision trees) by analysing the graph. Optimum value of n_estimators is found to be 300.
- The model is fitted and trained using these parameters.

MODEL EVALUATION

- The model is evaluated using Accuracy, Precision, Recall, F1 score.
- F1 score: It is the harmonic mean of Recall and precision. This metric is chosen for imbalanced multi class classification.
- The results are :
 - Accuracy=0.6375,
 - Precision=0.60,
 - Recall=0.6375,
 - F1 score=0.617

CONCLUSION

DECISION

- Although oversampling increased correlation between variables, it was found to decrease the accuracy and F1 scores. In the final model, it was avoided.
- With removal of column and outliers, the evaluation scores decreased slightly, but the training score decreased had this step been omitted. So, we followed the conventional method.
- The final results doesn't improve due to the highly imbalanced dataset.

THANK YOU !