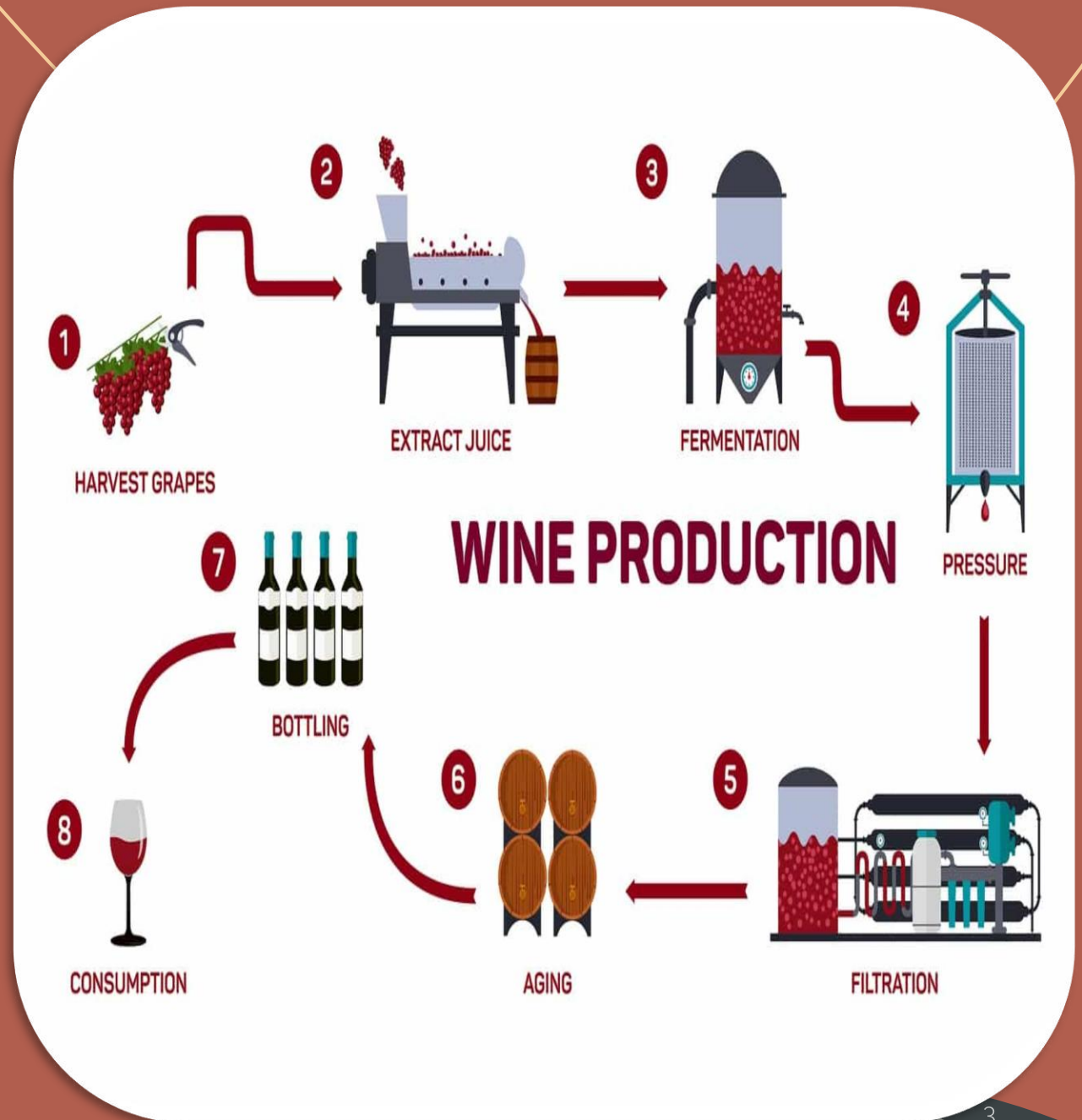# WINE QUALITY PREDICTION

Francis Osei

# Agenda

1. Introduction

2. Data Preparation

3. Machine Learning Models

4. Results and Discussion

5. Summary

# Introduction

- Wine is an alcoholic drink typically made from fermented grapes.

- It is very difficult to assess the quality of wine just by reading the label.

- Many quality wines are made to mature over a long period of time before finally reaching their best.



WINE PRODUCTION

1 HARVEST GRAPES
2 EXTRACT JUICE
3 FERMENTATION
4 PRESSURE
5 FILTRATION
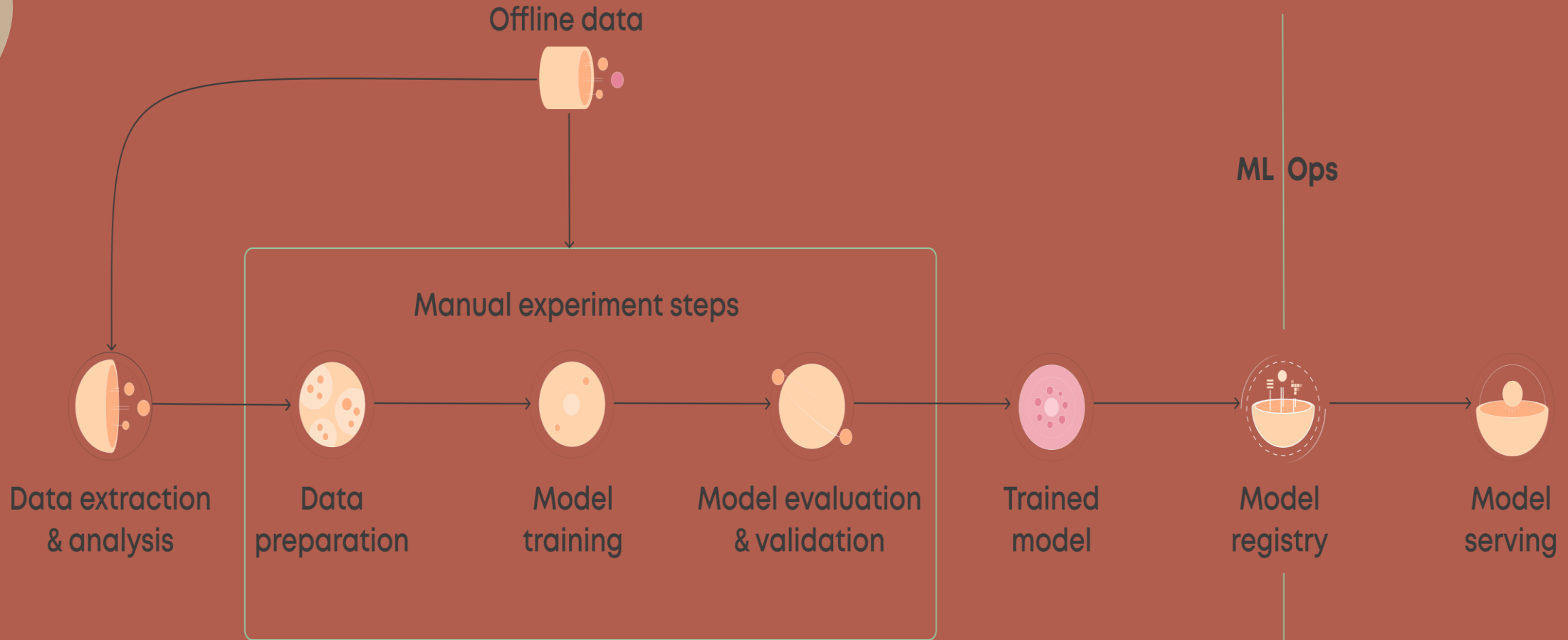6 AGING
7 BOTTLING
8 CONSUMPTION

# Problem Statement

- Price is only an indication of quality when similar wines are being compared.

- Determining wine quality is a complex task for companies due to various factors involved in assessing and evaluating wine.

# Objective

- To discover the quality of a wine (red and white variants of the Portuguese "Vinho Verde" wine) considering only its chemical properties.

- To identify the qualities of an excellent wine that are most indicative.
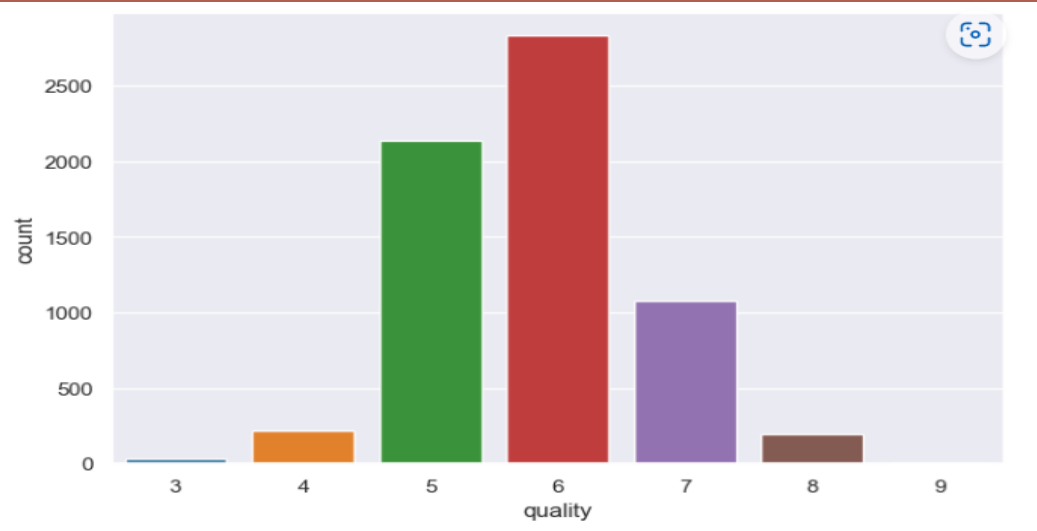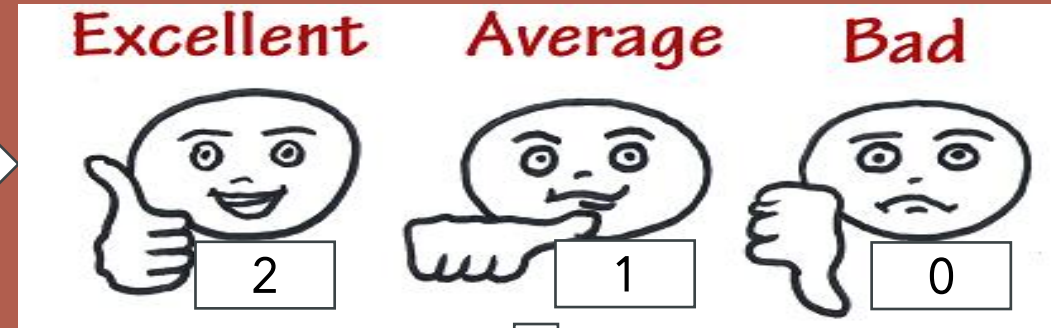
# Machine Learning Pipeline



Offline data

ML Ops

Manual experiment steps

Data extraction & analysis

Data preparation

Model training

Model evaluation & validation

Trained model

Model registry

Model serving

# Data Description

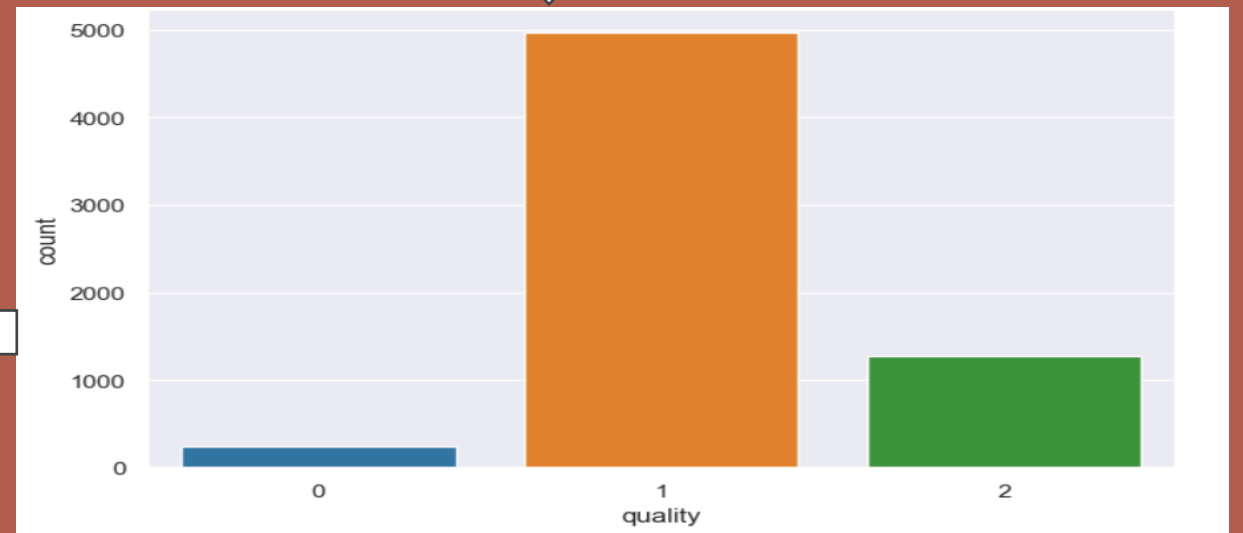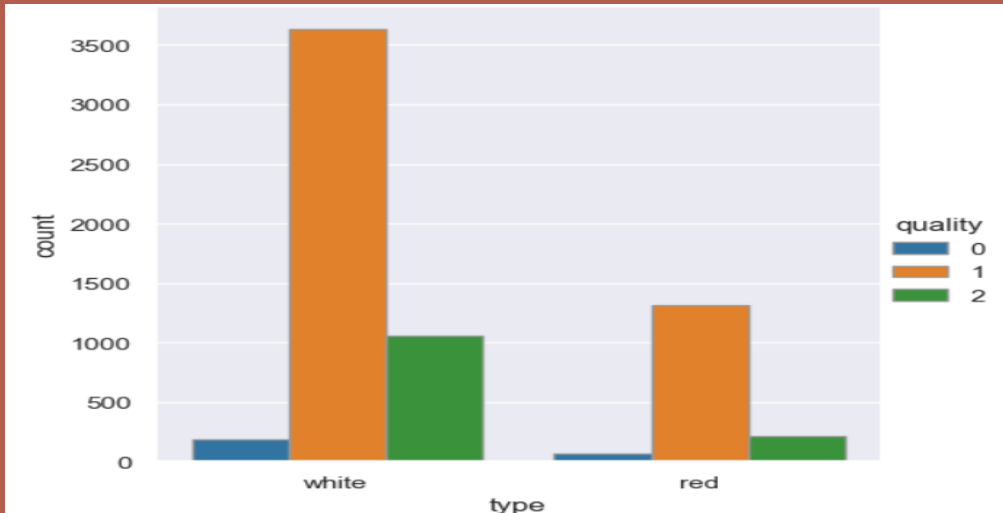| Variables | Description |
| --- | --- |
| Type | Red/White wine |
| Fixed acidity | The acids that naturally occur in the grapes used to ferment the wine and carry over into the wine. |
| Volatile acidity | Acids that evaporate at low temperatures |
| Citric acid | Citric acid is used as an acid supplement which boosts the acidity of the wine. |
| Residual sugar | The amount of sugar remaining after fermentation stops. |
| Chlorides | The amount of chloride salts (sodium chloride) present in the wine |
| Free sulfur dioxide | The free form of SO exists in equilibrium between molecular $SO_2$ (as a dissolved gas) and bisulfite ion |
| Total sulfur dioxide | The amount of free and bound forms of $SO_2$ |
| Density | The density of wine juice depending on the percent alcohol and sugar content |
| pH | A measure of the acidity of wine |
| Sulphates | Amount of potassium sulphate as a wine additive which can contribute to sulfur dioxide gas ($SO_2$) levels |
| Alcohol | How much alcohol is contained in each volume of wine |
| Quality | score between 0 (very bad) and 10 (very excellent) by wine experts |

# Data Preparation (Data Visualization)
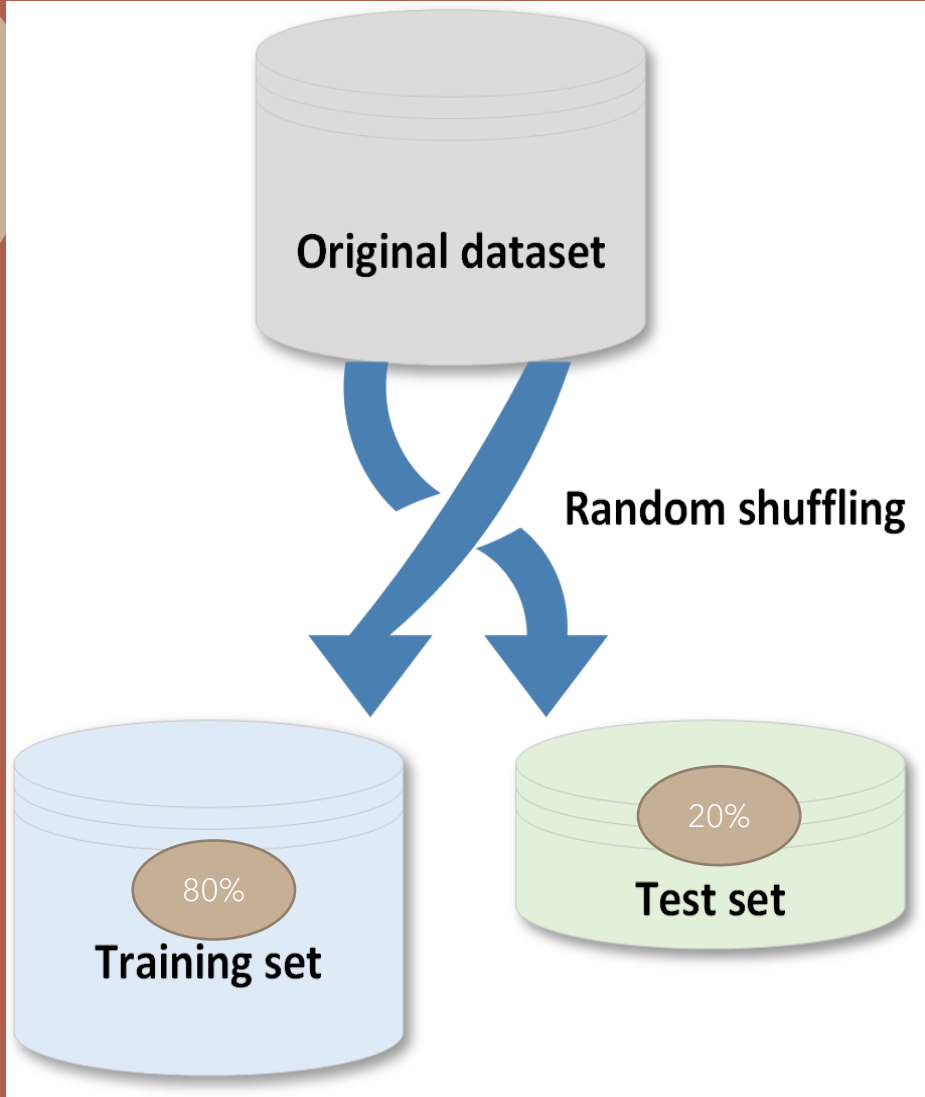


(A)

(B)

Excellent    Average    Bad

2    1    0

(C)

(D)

# Data Preparation (Missing Values)

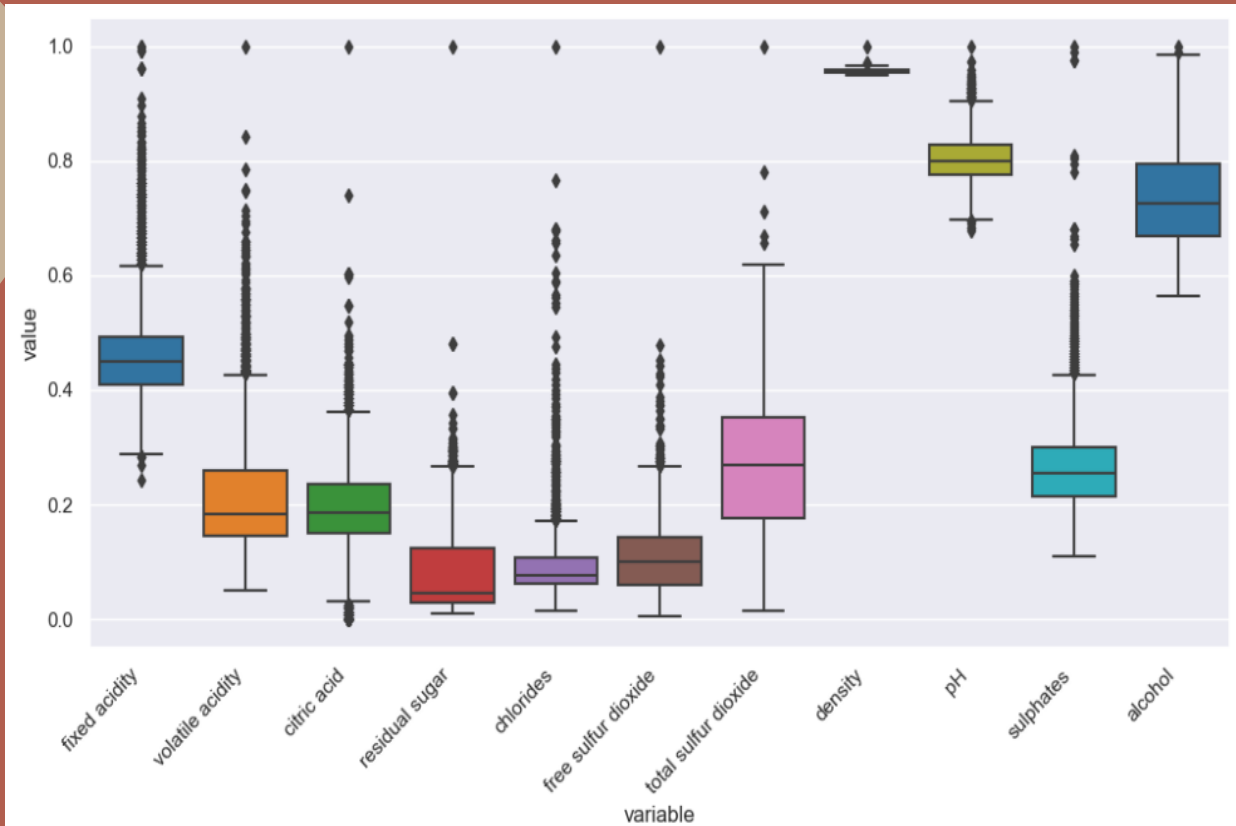| | Total | Percent |
|---|---|---|
| fixed acidity | 10 | 76.923077 |
| pH | 9 | 69.230769 |
| volatile acidity | 8 | 61.538462 |
| sulphates | 4 | 30.769231 |
| citric acid | 3 | 23.076923 |
| residual sugar | 2 | 15.384615 |
| chlorides | 2 | 15.384615 |
| type | 0 | 0.000000 |
| free sulfur dioxide | 0 | 0.000000 |
| total sulfur dioxide | 0 | 0.000000 |
| density | 0 | 0.000000 |
| alcohol | 0 | 0.000000 |
| quality | 0 | 0.000000 |

- Missing data can significantly affect the conclusions that can be drawn from the data.

- 0.005% of information is missing
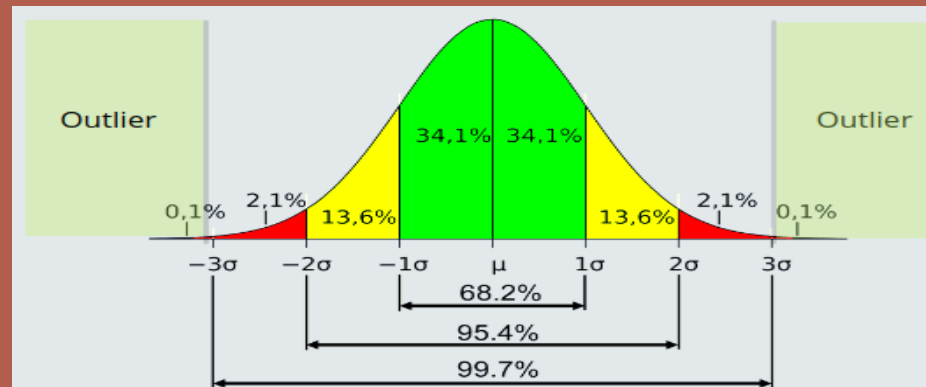
# Data Preparation (Data Splitting)



- Data splitting is an essential step in machine learning and data analysis.

- Training datasets comprise samples used to fit models under construction.

- A test dataset is a separate sample to provide an unbiased final evaluation of a model fit.
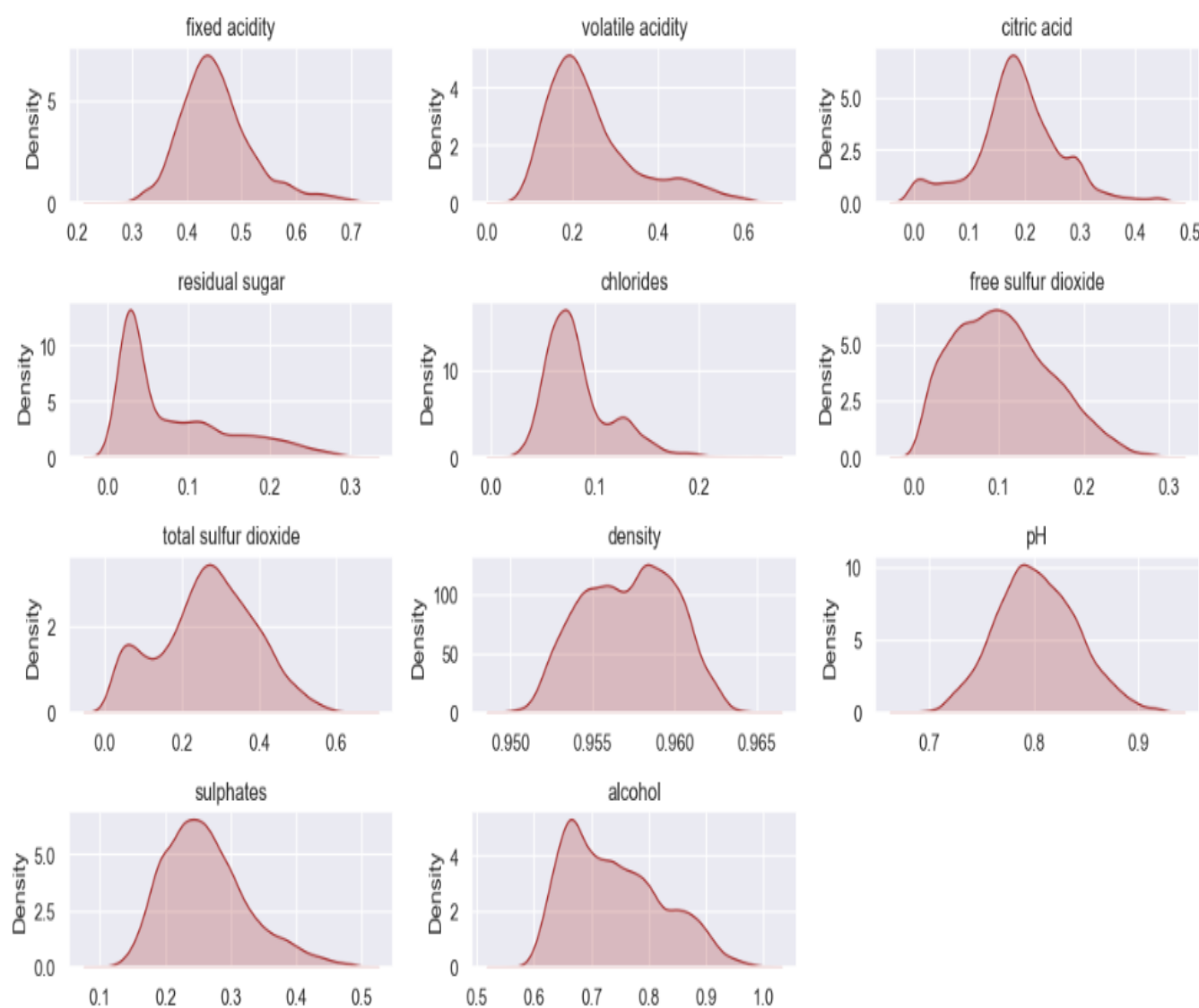
- Outliers in the data may causes problems during model fitting

- Outliers may inflate the error metrics which give higher weights to large errors
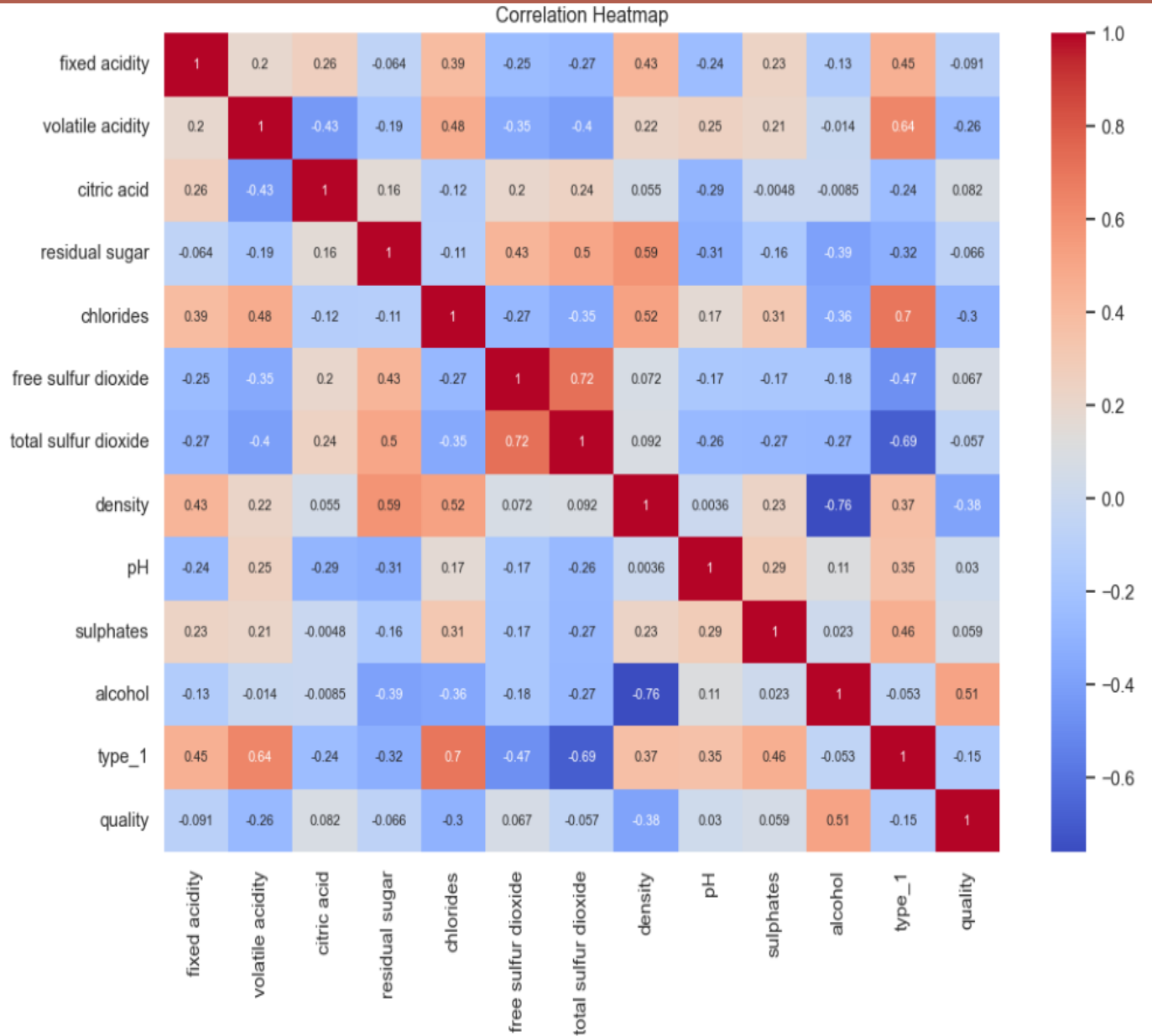
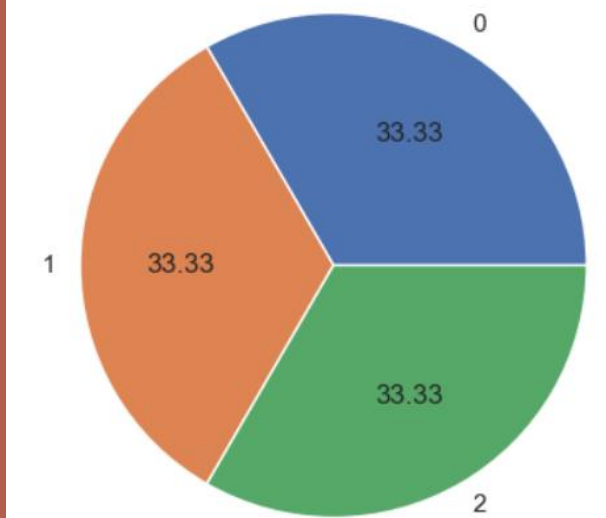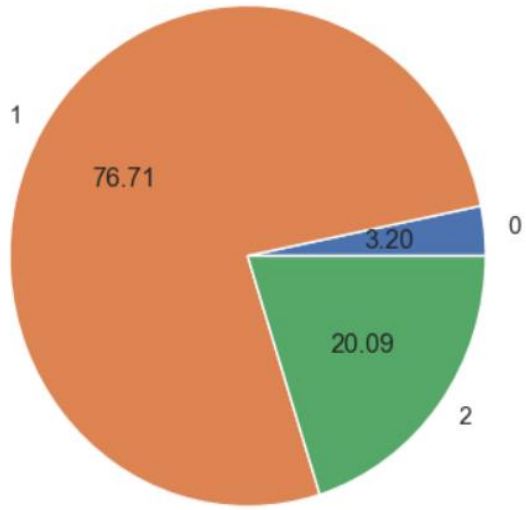# Data Preparation (Distribution)



- Different machine learning models have different assumptions about the distribution of the data.

- The distribution can provide insights into the normality or non-normality of the data.

# Data Preparation (Multicollinearity)
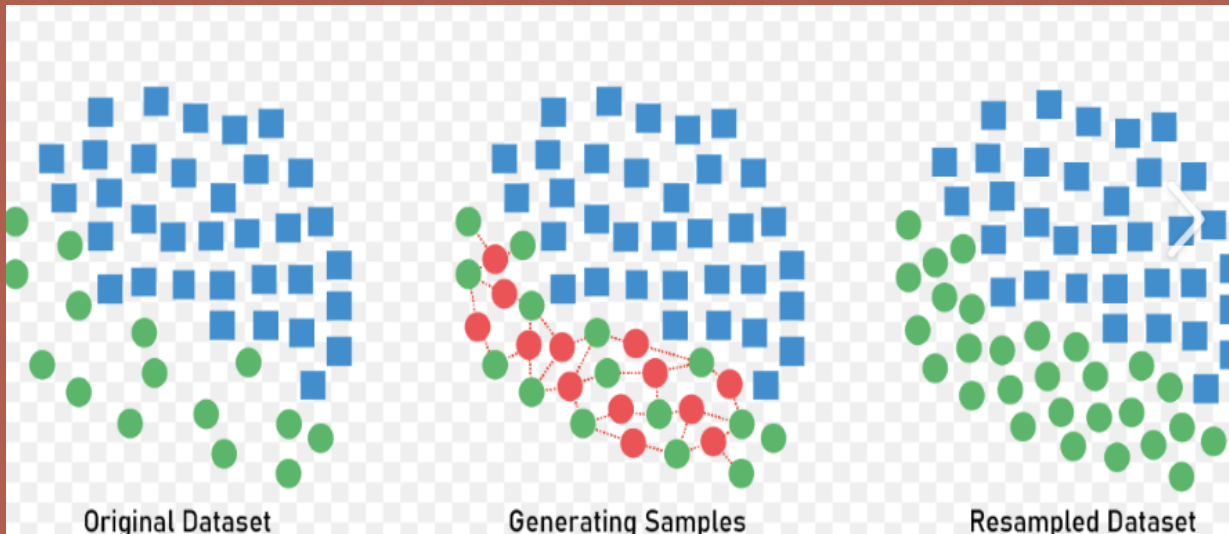

Correlation Heatmap

- Correlation is statistical technique which determines how one variables changes in relation with the other variable.

- Multicollinearity will severely affect performance outside of the original data sample.

# Data Preparation (Imbalance Dataset)







Original Dataset          Generating Samples          Resampled Dataset
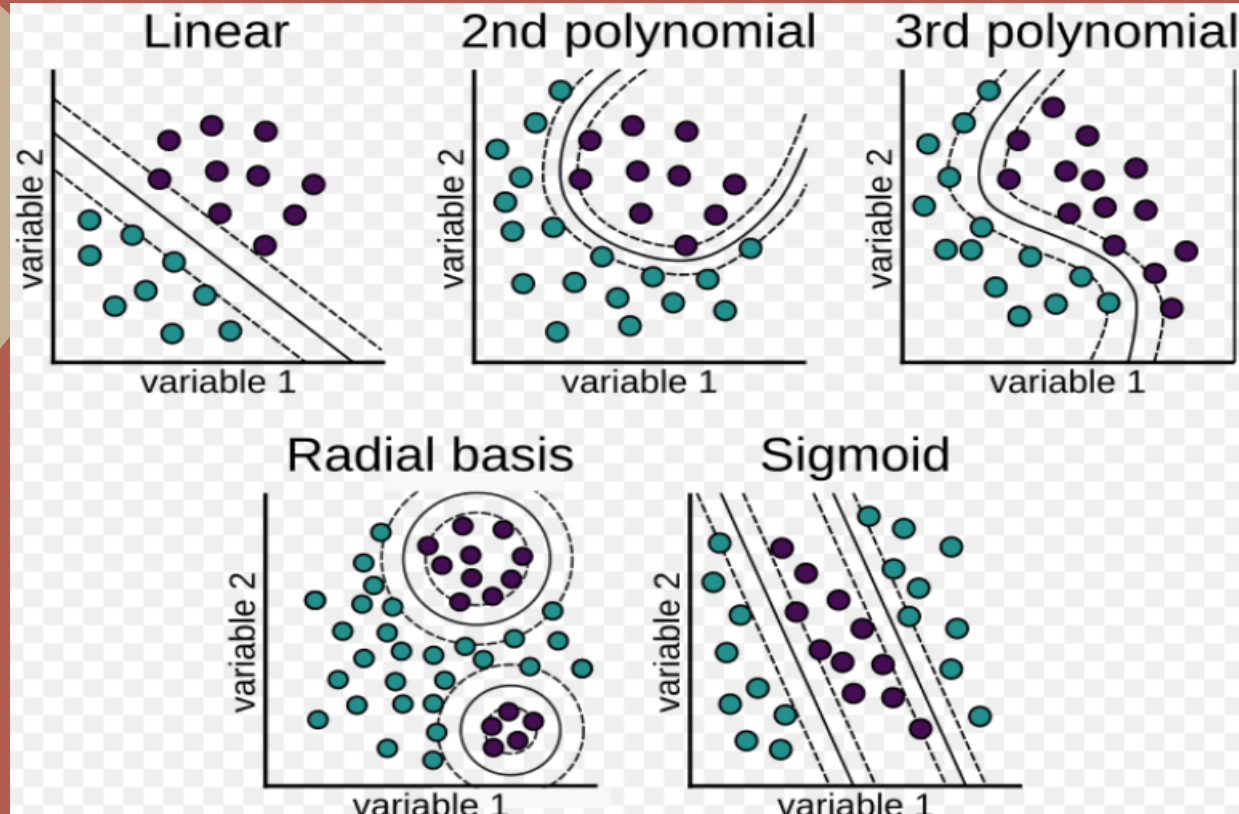
- Imbalance dataset arises when one set of classes dominate over another set of classes.

- it causes the machine learning model to be more biased towards majority class.

- SMOTE techniques generates the virtual training records by linear interpolation for the minority class.

Source: The Support Vector Machine: Basic Concept | by Aminah Mardiyyah Rufai | The Startup | Medium

- Support vector machine results from enlarging the feature space in a specific way using kernels

- The Radial Basis Function depends only on the distance between the input and some fixed point
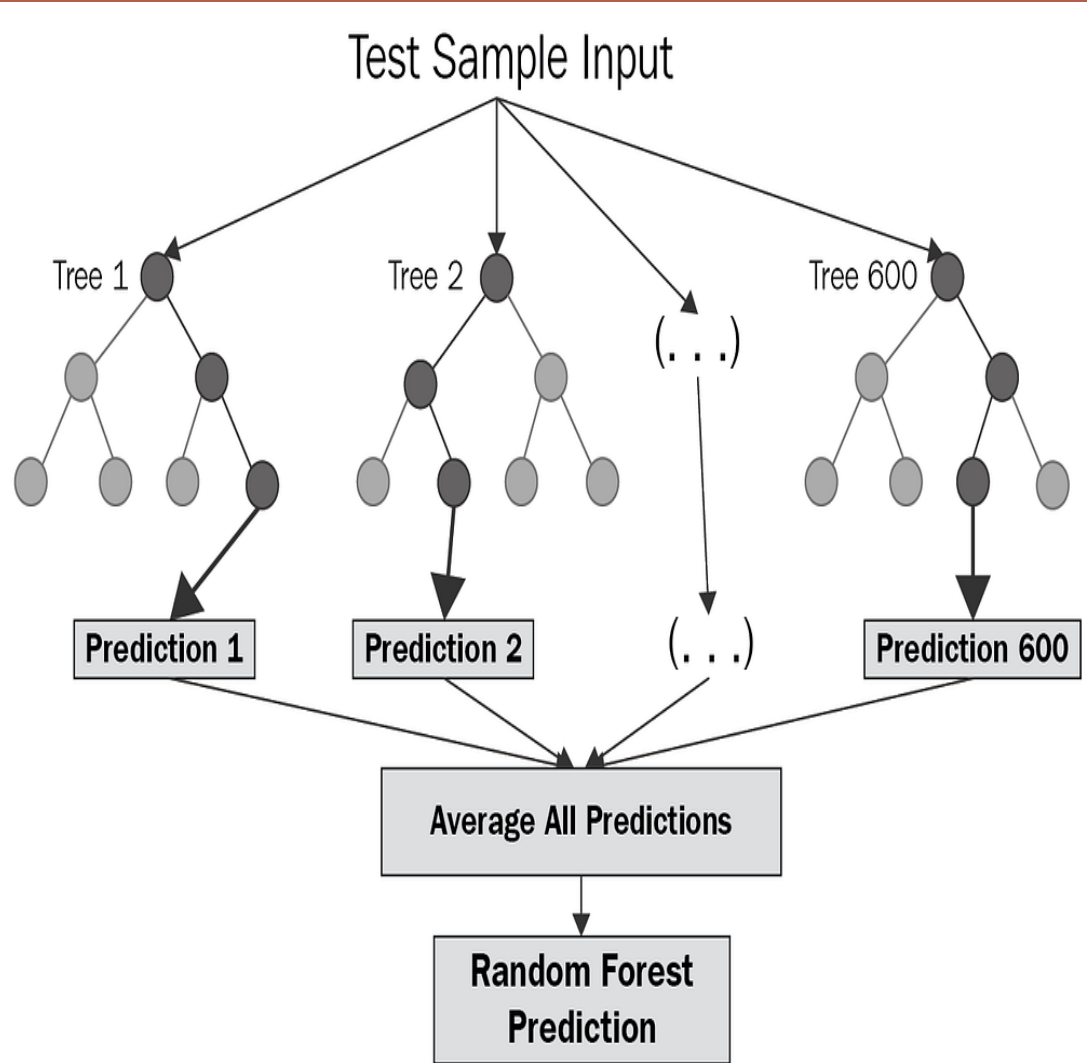
$$f(\mathbf{x}) = \sum_{i}^{N} \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b$$

weight (may be zero)

support vector

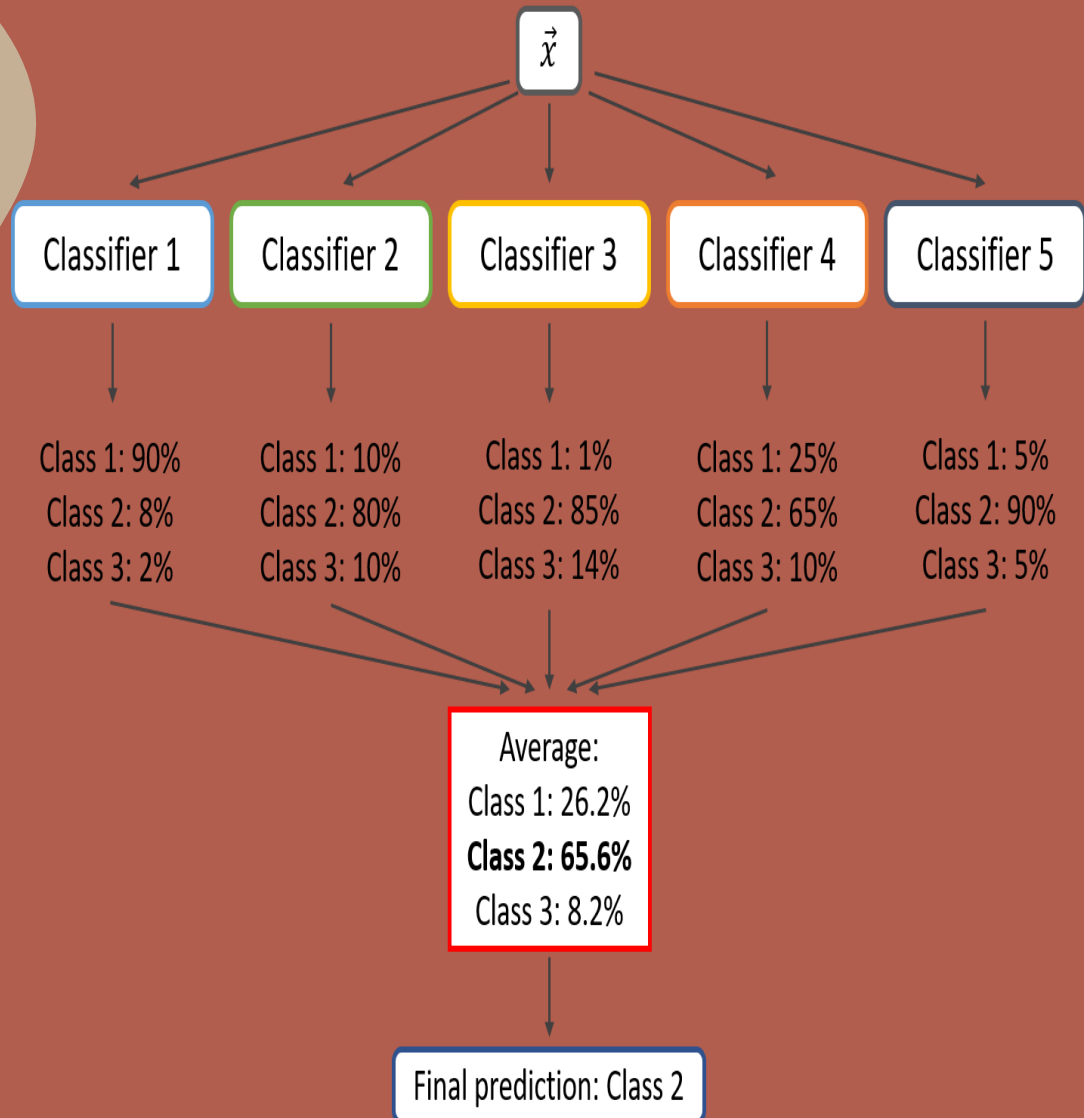Gaussian kernel $k(\mathbf{x}, \mathbf{x}') = \exp\left(-||\mathbf{x} - \mathbf{x}'||^2 / 2\sigma^2\right)$
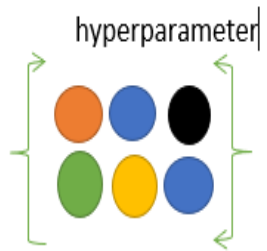
15

Source: Databricks

- Random forest ensemble is an ensemble of decision trees and a natural extension of bagging.

- Random forest handle large datasets with high dimensionality.

- It provides feature importance measures, allowing you to understand the relative importance of different features in the model.

# Model training (Voting Classifier)



- The Voting Classifier is an ensemble learning method that combines multiple individual classifiers to make predictions by voting.

- The predicted class probabilities from each classifier are averaged, and the class with the highest average probability is selected as the final prediction.

# Model training (Hyperparameter Tuning)
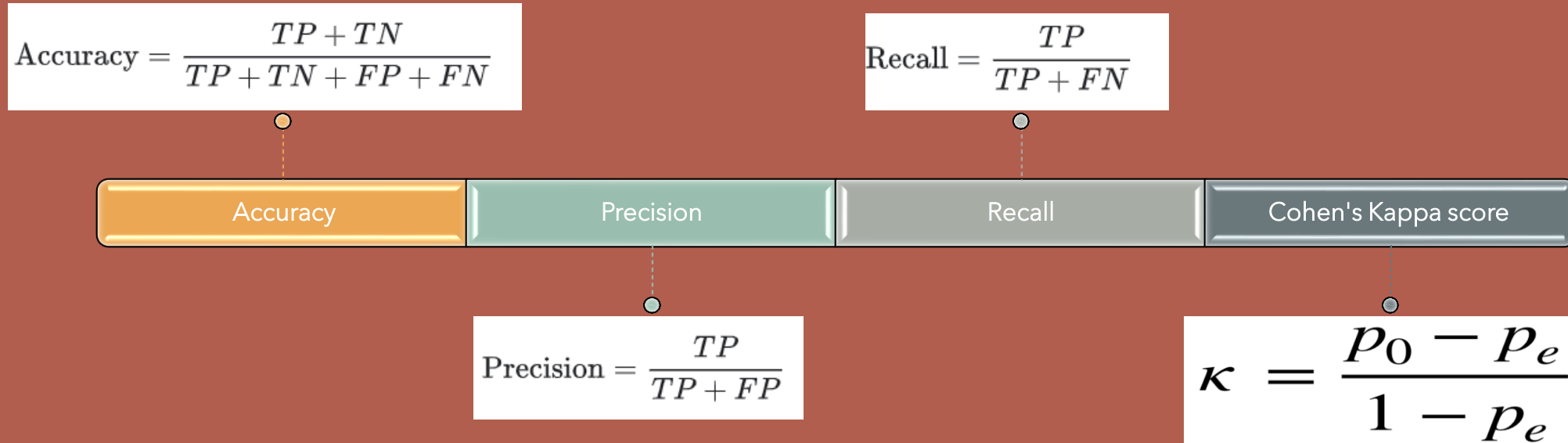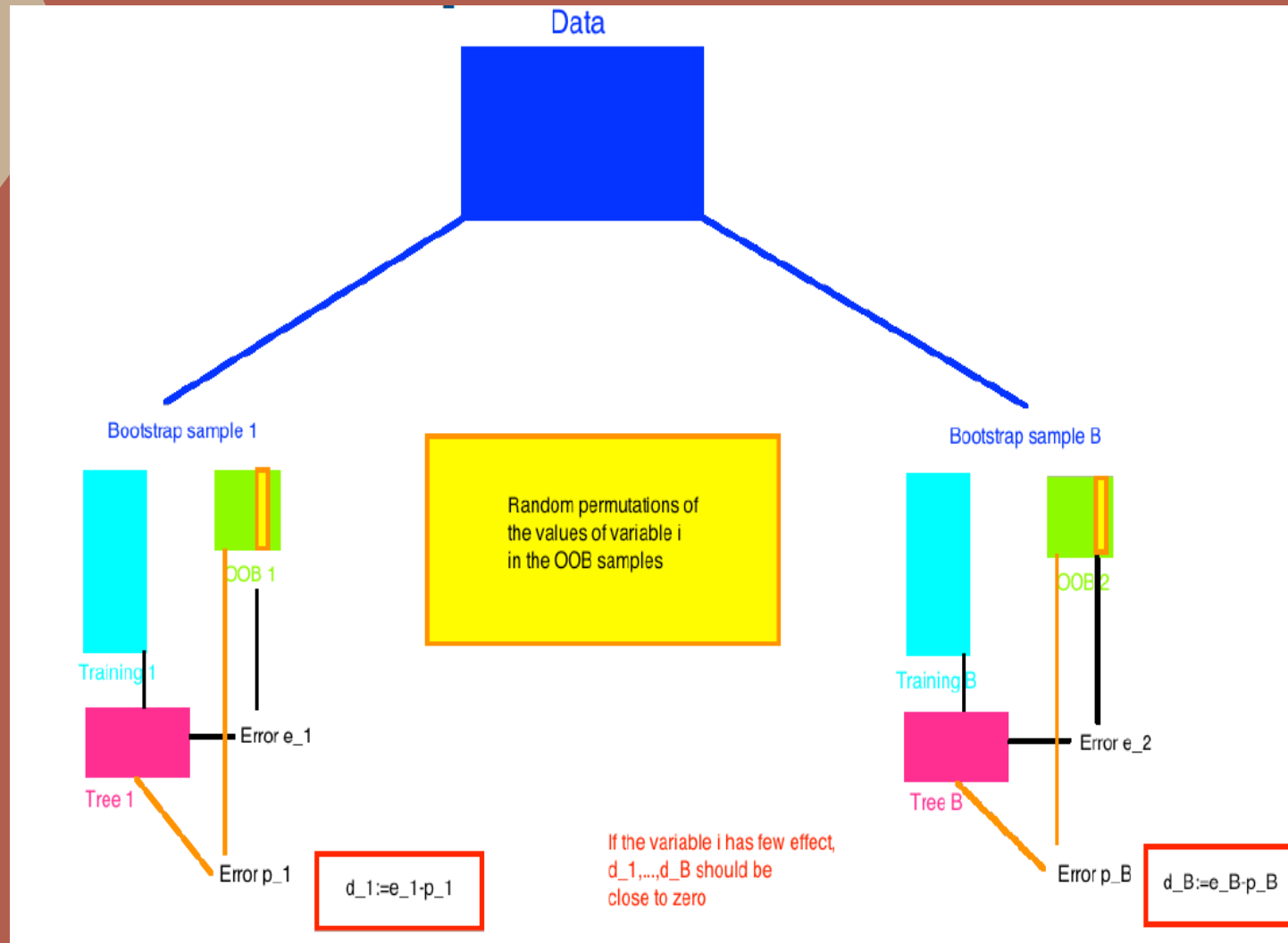


- Hyperparameter tuning is an essential step in machine learning to optimize the performance of your models.

- Gridsearch tries all possible combinations of hyperparameters and evaluates each combination using cross-validation or a validation set.

18

# Model Evaluation (Evaluation Metrics)

- Model evaluation is an essential step in machine learning to assess the performance and effectiveness of trained models.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

| Accuracy | Precision | Recall | Cohen's Kappa score |
|----------|-----------|--------|---------------------|

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\kappa = \frac{p_0 - p_e}{1 - p_e}$$

# Model Evaluation (Variable Importance)



**Data**

Bootstrap sample 1 — Bootstrap sample B

Random permutations of the values of variable i in the OOB samples

OOB 1 — OOB 2

Training 1 — Training B

Tree 1 — Tree B

Error e_1 — Error e_2

Error p_1 — Error p_B

$d\_1 := e\_1 - p\_1$ — $d\_B := e\_B - p\_B$

If the variable i has few effect, d_1,...,d_B should be close to zero
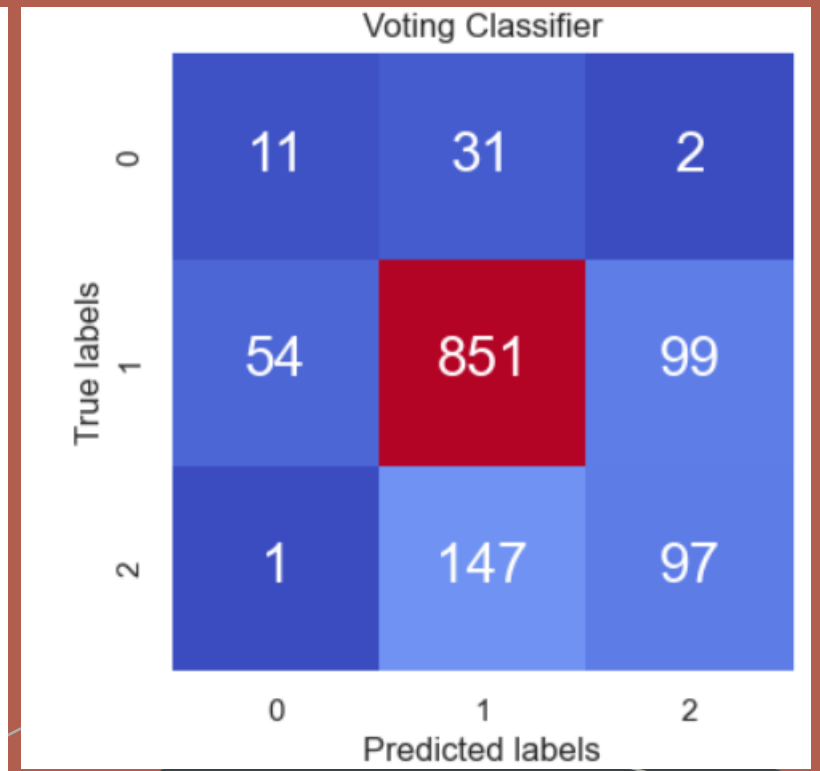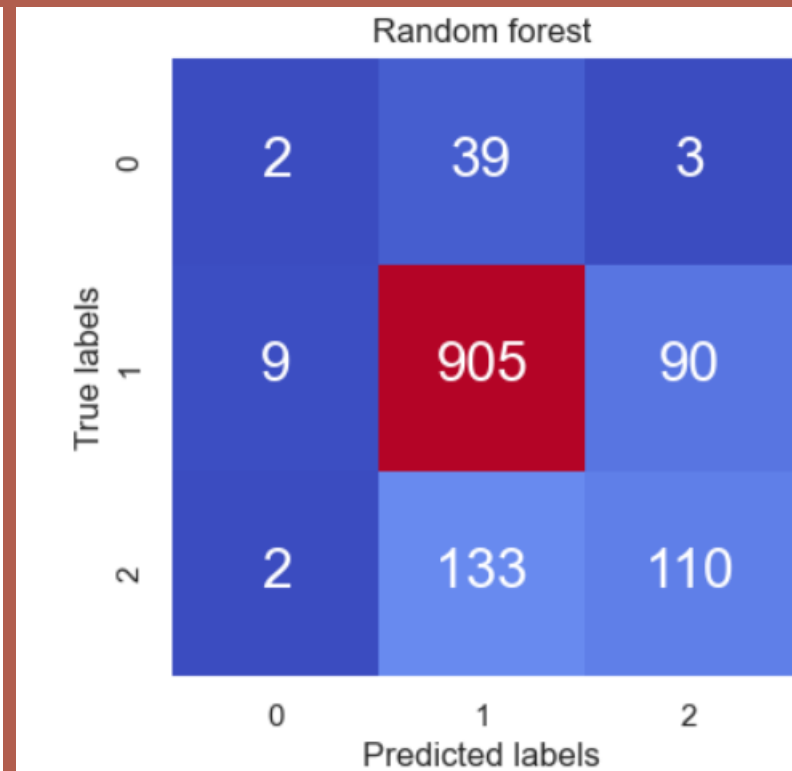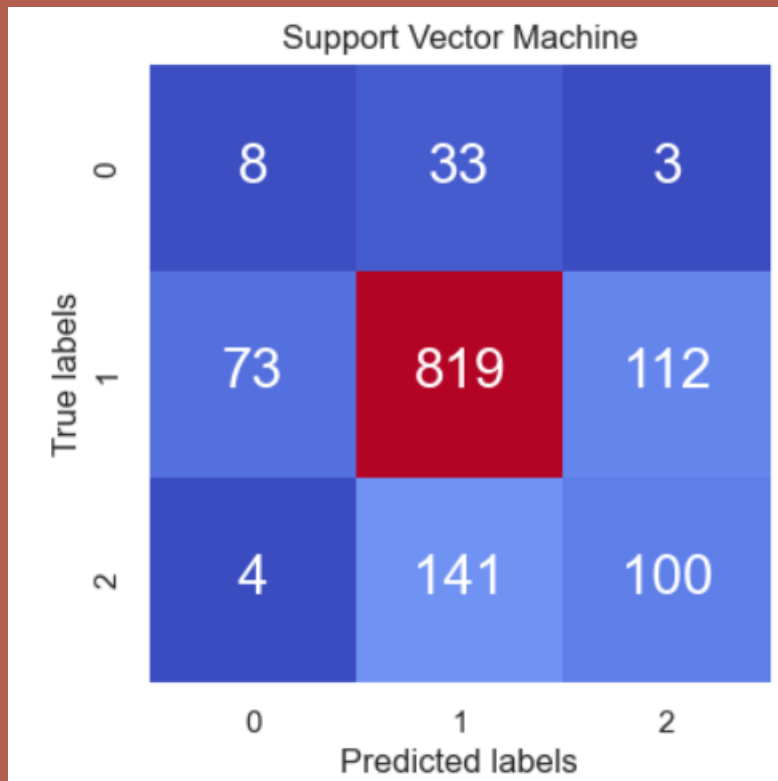
- Feature importance is a technique used in machine learning to determine the relevance or contribution of each feature in a predictive model.

# Result and Discussion (1)

- Voting classifiers can identity bad wines compared to the other models

- ML models are very good in identifying average wines

# Result and Discussion (2)

- Classification Report

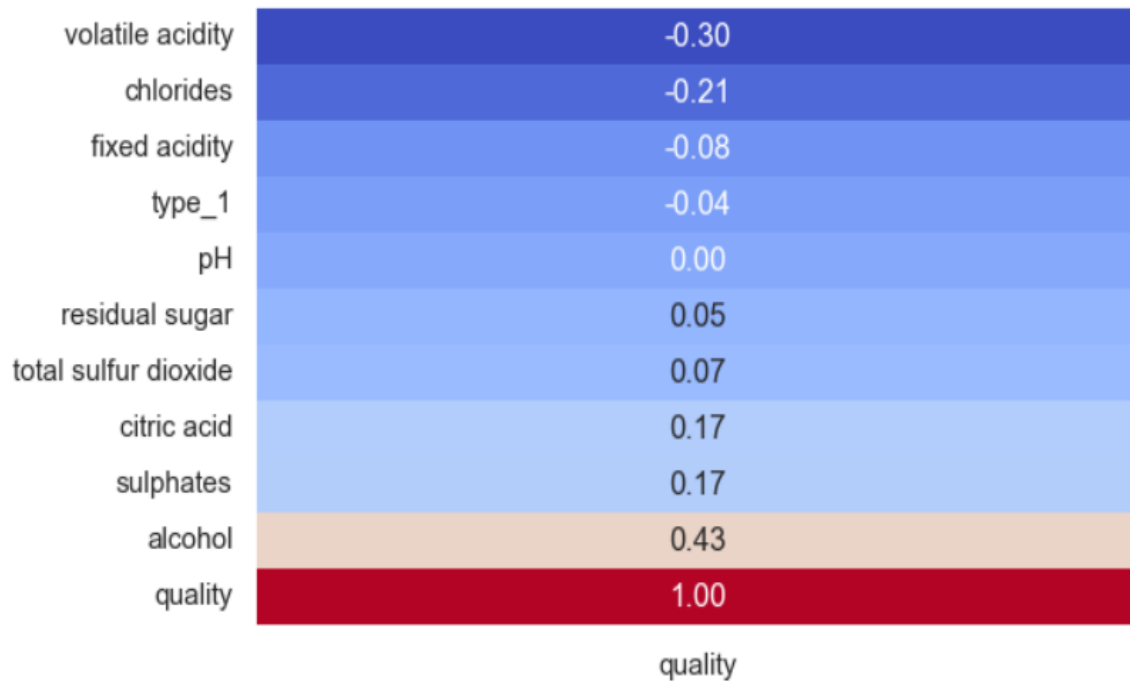| Models | Accuracy | Precision (weighted average) | Recall (weighted average) | Kohen's Kappa Score |
|---|---|---|---|---|
| Support Vector Machine | 0.72 | 0.73 | 0.72 | 0.23 |
| Random Forest | 0.78 | 0.75 | 0.78 | 0.31 |
| Voting Classifier | 0.74 | 0.74 | 0.74 | 0.27 |

- Feature Importance



Correlation between Target Variable and Predictors

| | |
|---|---|
| volatile acidity | -0.30 |
| chlorides | -0.21 |
| fixed acidity | -0.08 |
| type_1 | -0.04 |
| pH | 0.00 |
| residual sugar | 0.05 |
| total sulfur dioxide | 0.07 |
| citric acid | 0.17 |
| sulphates | 0.17 |
| alcohol | 0.43 |
| quality | 1.00 |

quality



Variable Importance

# Summary

- ML models are very good in identifying average wines, but it's terrible in finding poor wines.

- In Business perspective, the voting classify will be preferred over the Random forest

- More features are needed to improve the model performance

# References

- Olivier Lopez (2022), Machine Learning and Insurance : Workshop Actuarial Science

- Martin Gubri & Salah Ghamizi (2022), Applied Machine Learning Lecture Note

- Gareth James & Daniela Witten & Trevor Hastie & Robert Tisbshirani (2013), An Introduction to Statistical Learning

# Thank You

https://github.com/Franosei/Wine-Quality