

# Assignment 1 Report

## 1. Obtaining Dataset

In this assignment, I used the **Wine Quality** datasets from **UCI Dataset**.

The table below contains all the attribute for the dataset:

Variable Name	Role	Type
fixed_acidity	Feature	Continuous
volatile_acidity	Feature	Continuous
citric_acid	Feature	Continuous
residual_sugar	Feature	Continuous
chlorides	Feature	Continuous
free_sulfur_dioxide	Feature	Continuous
total_sulfur_dioxide	Feature	Continuous
density	Feature	Continuous
pH	Feature	Continuous
sulphates	Feature	Continuous
alcohol	Feature	Continuous
quality	Target	Integer

## 2. Preprocessing Dataset

The data in the dataset do have **missing value**. However, there exist some outlier for each features. In this assignment, I decided to keep the outlier since it does not affects our model.

Below contains the correlation between our target, **quality** and other features:

Feature	Correlation
fixed acidity	0.124052
volatile acidity	-0.390558
citric acid	0.226373
residual sugar	0.013732
chlorides	-0.128907
free sulfur dioxide	-0.050656
total sulfur dioxide	-0.185100
density	-0.174919
pH	-0.057731
sulphates	0.251397
alcohol	0.476166

We noticed that features such as **residual sugar**, **free sulfur dioxide** and **pH** do not play a huge role for our model so we can exclude the features from our model.

### 3. Regression Model

#### 3.1. SGDRegressor

We use **Standard Scaler** to normalize our dataset in order to have better performance. We divide the datasets into **80/20** split for training and testing.

After testing the SGDRegressor with fixed parameter of **max iterations: 100\_000**, **tolerance: 1e-3**, we have an **R2** score of around **0.35**.

I initially thought that the low score might be due to the outlier data in the dataset. I tried removing the outlier from the dataset. However, the **R2** score seems to be worse than the one without removed outliers.

After that I tried hyperparameter tuning using **Grid Search** from sklearn, with the parameter as below: - max-iterations: [1000, 100\_000, 1\_000,000] - tolerance: [1e-3, 1e-4, 1e-5, 1e-6]

After performing the **Grid Search** on **SGDRegressor**, we have the following best hyperparameter: **{'max\_iter': 1000000, 'tol': 0.001}** with an **R2** score of 0.35

In my conclusion for the low **R2** score, I believe that **Linear Regression** model is not suitable for this problem because the data that we have is not close to linear, causing our regression results to deviate from the ground truth.

#### 3.2. Ordinary Least Square(OLS)

I used the **OLS** model from **statsmodel** package. The OLS model gave us a similar **R2** score compared to the **SGDRegressor**. This confirms that our conclusion above is correct.