

```
[1]: # ignore the warnings
import warnings
warnings.filterwarnings('ignore')
warnings.filterwarnings('ignore')

# data visualisation and manipulation
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from matplotlib import style
import seaborn as sns
import missingno as msno
#configure
#sns.set(style='whitegrid',color_codes=True)
#plot the necessary modelling algo.

#classification
from sklearn.linear_model import LogisticRegression
from sklearn.svm import LinearSVC, SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier, AdaBoostClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.naive_bayes import GaussianNB

#regression
from sklearn.linear_model import LinearRegression, Ridge, Lasso, RidgeCV
from sklearn.ensemble import RandomForestRegressor, BaggingRegressor, GradientBoostingRegressor, AdaBoostRegressor
from sklearn.svm import SVR
from sklearn.neighbors import KNeighborsRegressor

#model selection
from sklearn.model_selection import train_test_split, cross_validate
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import GridSearchCV

#preprocessing
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import MinMaxScaler, StandardScaler, LabelEncoder

#evaluation metrics
from sklearn.metrics import mean_squared_log_error, mean_squared_error, r2_score, mean_absolute_error # for regression
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score # for classification
```

```
In [2]: df=pd.read_csv('winequality.csv')
```

```
In [3]: df.shape
```

```
Out[3]: (1599, 12)
```

```
In [4]: #checking columns
df.columns
```

```
Out[4]: Index(['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
        'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
        'pH', 'sulphates', 'alcohol', 'quality'],
        dtype='object')
```

```
In [5]: df.info()
```

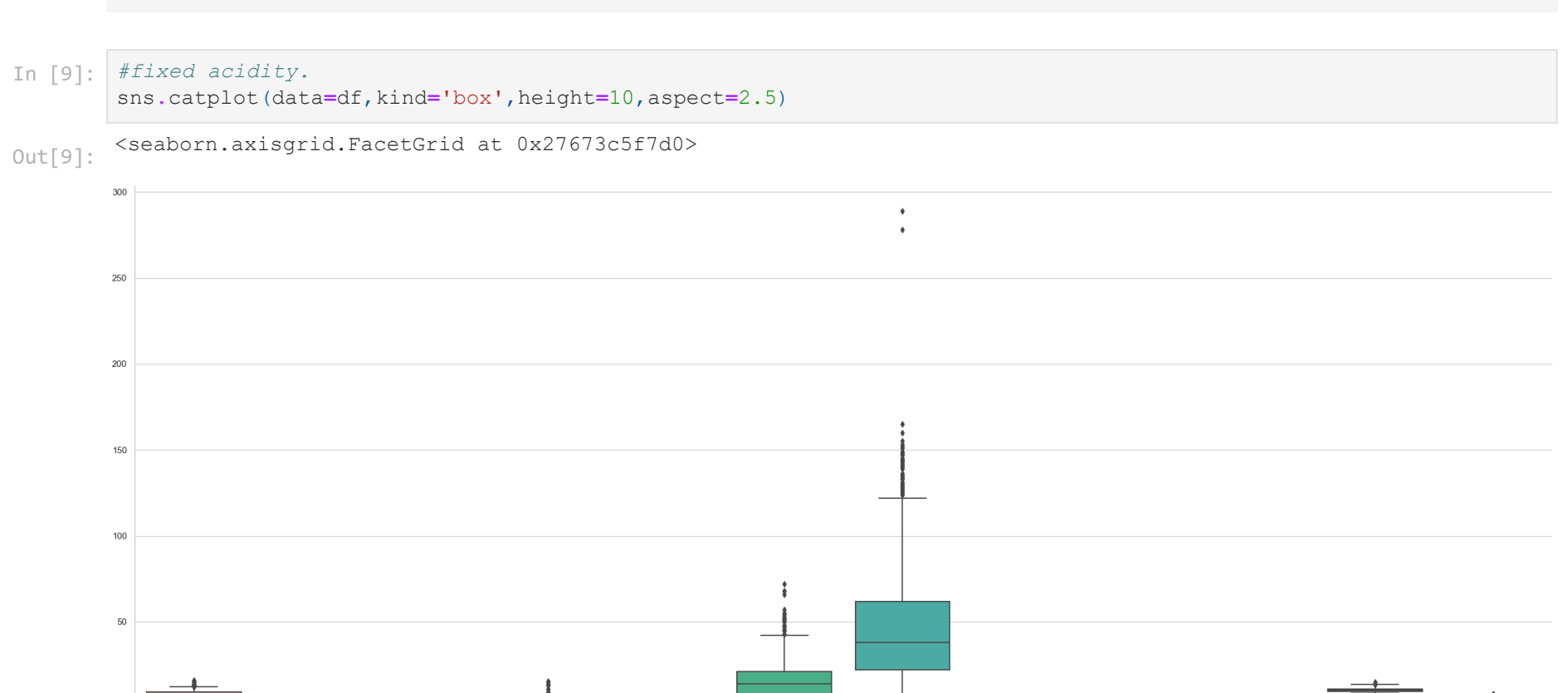
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1599 entries, 0 to 1598
Data columns (total 12 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   fixed acidity        1599 non-null    float64
 1   volatile acidity     1599 non-null    float64
 2   citric acid          1599 non-null    float64
 3   residual sugar       1599 non-null    float64
 4   chlorides            1599 non-null    float64
 5   free sulfur dioxide  1599 non-null    float64
 6   total sulfur dioxide 1599 non-null    float64
 7   density              1599 non-null    float64
 8   pH                  1599 non-null    float64
 9   sulphates           1599 non-null    float64
10   alcohol             1599 non-null    float64
11   quality              1599 non-null    int64
dtypes: float64(11), int64(1)
memory usage: 150.0 KB
```

```
In [6]: #checking null
df.isnull().sum()
```

```
Out[6]: fixed acidity      0
volatile acidity    0
citric acid        0
residual sugar     0
chlorides          0
free sulfur dioxide 0
total sulfur dioxide 0
density            0
pH                0
sulphates          0
alcohol            0
quality            0
dtype: int64
```

```
In [7]: mmo.matrix(df)
```

```
Out[7]: <Axes: >
```



UNIVARIATE ANALYSIS.

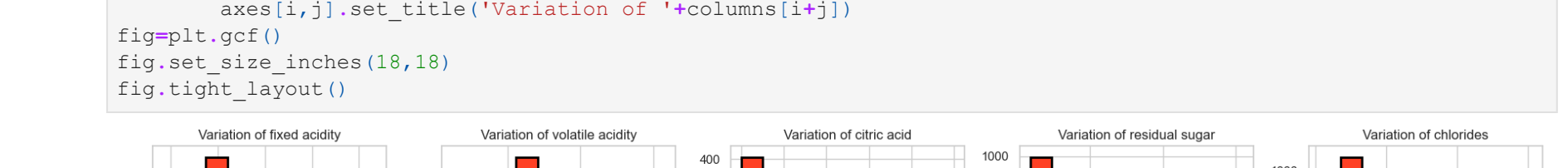
```
In [8]: df.describe(include='all')
```

```
Out[8]:
```

```
In [9]: #fixed acidity.
```

```
sns.catplot(data=df, kind='box', height=10, aspect=2.5)
```

```
Out[9]: <seaborn.axisgrid.FacetGrid at 0x27673c57d0>
```



```
In [10]: #using a histogram.
```

```
fig, axes=plt.subplots(5,5)
```

```
columns=['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
```

```
        'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
```

```
        'pH', 'sulphates', 'alcohol', 'quality']
```

```
for i in range(5):
```

```
    for j in range(5):
```

```
        axes[i,j].hist(x=columns[i+j], data=df, edgecolor='b', linewidth=2, color='ff4125')
```

```
        fig=plt.gcf()
```

```
        fig.set_size_inches(18,18)
```

```
        fig.tight_layout()
```



CORRELATION BETWEEN DIFFERENT FEATURES

```
In [11]: #correlation matrix.
```

```
cor_mat=corr()
```

```
mask = np.triu(corr_mat)
```

```
mask[np.triu_indices_from(mask)] = False
```

```
fig=plt.gcf()
```

```
fig.set_size_inches(30,12)
```

```
sns.heatmap(data=cor_mat, mask=mask, square=True, annot=True, cbar=True)
```

```
Out[11]: <Axes: >
```



1. The quality of wine is highly related to volatile acidity.
2. Also the quality of wine is highly correlated to alcohol.
3. pH and citric acid/ fixed acidity are highly inversely related as all of us know that acids have smaller pH values.
4. Self Relation ie of a feature to itself is 1 as expected.
5. some other similar inferences can be drawn.

HOW QUALITY VARIES WITH DIFFERENT NUMERIC FEATURES.

```
In [12]: def plot(feature_x, target='quality'):
```

```
sns.catplot(x=target, y=feature_x, data=df, kind='bar', height=5, aspect=1)
```

```
group = sns.factorplot(x=target, y=feature_x, data=df, kind='violin', height=5, aspect=1)
```

```
sns.catplot(x=target, y=feature_x, data=df, kind='swarm', height=5, aspect=1)
```

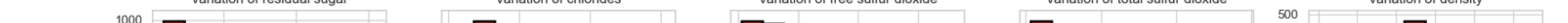
```
In [13]: # for fixed acidity.
```

```
plot('fixed acidity', 'quality')
```



```
In [14]: # for alcohol.
```

```
plot('alcohol', 'quality')
```



```
In [15]: # similarly we can check for other variables.
```

MODELLING

```
In [16]: bins = (2, 6.5, 8)
```

```
group_name = ['good', 'bad']
```

```
df['quality'] = pd.cut(df['quality'], bins = bins, labels = group_name)
```

```
In [17]: label_quality = LabelEncoder()
```

```
In [18]: #Bad becomes 0 and good becomes 1
```

```
df['quality'] = label_quality.fit_transform(df['quality'])
```

```
In [19]: x_train,x_test,y_train,y_test=train_test_split(df.drop('quality',axis=1),df['quality'], test_size=0.25,random_state=42)
```

```
In [20]: models=[LogisticRegression(), LinearSVC(), SVC(kernel='rbf'), KNeighborsClassifier(), RandomForestClassifier(),
```

```
            DecisionTreeClassifier(), GradientBoostingClassifier(), GaussianNB()]
```

```
model_names=['LogisticRegression', 'LinearSVC', 'rbfSVM', 'KNearestNeighbors', 'RandomForestClassifier', 'DecisionTree',
```

```
            'GradientBoostingClassifier', 'GaussianNB']
```

```
acc=[]
```

```
for model in range(len(models)):
```

```
    clf=models[model]
```

```
    clf.fit(x_train,y_train)
```

```
    pred=clf.predict(x_test)
```

```
    acc.append(accuracy_score(pred,y_test))
```

```
d={'Modelling Algo':model_names, 'Accuracy':acc}
```

```
d
```

```
Out[20]: {'Modelling Algo': ['LogisticRegression',
```

```
                        'LinearSVM',
```

```
                        'rbfSVM',
```

```
                        'KNearestNeighbors',
```

```
                        'RandomForestClassifier',
```

```
                        'DecisionTree',
```

```
                        'GradientBoostingClassifier',
```

```
                        'GaussianNB'],
```

```
        'Accuracy': [0.875, 0.875, 0.87, 0.8625, 0.91, 0.895, 0.8775, 0.8525]}
```

```
In [21]: acc_frame=pd.DataFrame(d)
```

```
acc_frame
```

```
Out[21]:
```

	Modelling Algo	Accuracy
0	LogisticRegression	0.8750
1	LinearSVM	0.8750
2	rbfSVM	0.8700
3	KNearestNeighbors	0.8625
4	RandomForestClassifier	0.9100
5	DecisionTree	0.8950
6	GradientBoostingClassifier	0.8775
7	GaussianNB	0.8525

```
In [22]: sns.barplot(y='Modelling Algo',x='Accuracy',data=acc_frame)
```

```
Out[22]: <Axes: xlabel='Accuracy', ylabel='Modelling Algo'>
```



```
In [23]: sns.catplot(x='Modelling Algo',y='Accuracy',data=acc_frame, kind='point', height=4, aspect=3.5)
```

```
Out[23]: <seaborn.axisgrid.FacetGrid at 0x27676912f90>
```



THIS IS WITHOUT FEATURE SCALING. NOW SINCE FEATURES HAVE DIFFERENT SCALES LET US TRY TO DO FEATURE SCALING AND SEE THE IMPACT.

```
In [24]: def func(x_train,x_test,y_train,y_test,name_scaler):
```

```
    models=[LogisticRegression(), LinearSVC(), SVC(kernel='rbf'), KNeighborsClassifier(), RandomForestClassifier(),
```

```
            DecisionTreeClassifier(), GradientBoostingClassifier(), GaussianNB()]
```

```
    acc=[]
```

```
    for model in range(len(models)):
```

```
        clf=models[model]
```

```
        clf.fit(x_train,y_train)
```

```
        pred=clf.predict(x_test)
```

```
        acc.append(accuracy_score(pred,y_test))
```

```
    acc_frame[name_scaler]=np.array(acc)
```

```
acc_frame
```

```
In [25]: scaler=[MinMaxScaler(), StandardScaler()]
```

```
names=['Acc_Min_Max_Scaler', 'Acc_Standard_Scaler']
```

```
for scaler in range(len(scalers)):
```

```
    scaler=scaler[scaler]
```

```
    scaled_df=scaler.transform(df)
```

```
    X=scaled_df[:,0:11]
```

```
    y=df['quality'].to_numpy()
```

```
    x_train,x_test,y_train,y_test=train_test_split(X,y, test_size=0.25, random_state=42)
```

```
    func(x_train,x_test,y_train,y_test,name_scaler)
```

```
In [26]: acc_frame
```

```
Out[26]:
```

	Modelling Algo	Accuracy	Acc_Min_Max_Scaler	Acc_Standard_Scaler
0	LogisticRegression	0.8750	0.8800	0.8775
1	LinearSVM	0.8750	0.8825	0.8825
2	rbfSVM	0.8700	0.8850	0.8900
3	KNearestNeighbors	0.8625	0.8850	0.8950
4	RandomForestClassifier	0.9100	0.9100	0.9025
5	DecisionTree	0.8950	0.8925	0.8925
6	GradientBoostingClassifier	0.8775	0.8750	0.8775
7	GaussianNB	0.8525	0.8425	0.8425

1. Note that here the accuracies increase marginally on scaling.
2. Also for this data, StandardScaling seems to give slightly better results than the MinMaxScaling.
3. For some modelling algo there is a considerable increase in accuracies upon scaling the features like SVM, KNN whereas for others there isn't a considerable increase in accuracies upon scaling.

```
In [27]: # Just to visualize the accuracies.
```

```
sns.barplot(y='Modelling Algo',x='Accuracy',data=acc_frame)
```

```
Out[27]: <Axes: xlabel='Accuracy', ylabel='Modelling Algo'>
```



```
In [28]: sns.barplot(y='Modelling Algo',x='Acc_Min_Max_Scaler',data=acc_frame)
```

```
Out[28]: <Axes: xlabel='Acc_Min_Max_Scaler', ylabel='Modelling Algo'>
```



```
In [29]: sns.barplot(y='Modelling Algo',x='Acc_Standard_Scaler',data=acc_frame)
```

```
Out[29]: <Axes: xlabel='Acc_Standard_Scaler', ylabel='Modelling Algo'>
```



```
In [30]: # Preparing the features by using a StandardScaler as it gave better results.
```

```
scaler=StandardScaler()
```

```
scaled_df=scaler.fit_transform(df)
```

```
rescaled_df=scaled_df[:,0:11]
```

```
Y=df['quality'].to_numpy()
```

```
x_train,x_test,y_train,y_test=train_test_split(X,Y, test_size=0.25, random_state=42)
```

PARAMETER TUNING

1. LOGISTIC REGRESSION.

```
In [31]: param_dict={'C':[0.001, 0.01, 0.1, 1, 10, 100, 1000], 'penalty':['l1', 'l2']}
```

```
clf_lr=GridSearchCV(estimator=LogisticRegression(), param_grid=param_dict, scoring='accuracy', cv=10)
```

```
Out[31]:
```



```
In [32]: clf_lr.best_params_
```

```
Out[32]: {'C': 0.1, 'penalty': 'l2'}
```

```
In [33]: clf_lr.best_score_ # the best accuracy obtained by Grid search on the train set.
```

```
Out[33]: 0.881592630532213
```

```
In [34]: clf_lr.cv_results_
```