

# Machine Learning Report

Siddharth Bahekar  
College of Engineering  
Northeastern University  
Toronto, ON  
[bahekar.si@northeastern.edu](mailto:bahekar.si@northeastern.edu)

## Abstract

In this report, I have applied and explored the performance of supervised learning algorithm on two datasets. This report outlines the exploratory data analysis (EDA) and modeling process undertaken to train the Artificial Neural Network (ANN) algorithm. The two datasets included are **Heart Attack Analysis & Prediction** dataset and **Wine Quality** dataset. The learning curve, model complexity and model training has been performed and analyzed on both the datasets.

### 1. Dataset

The **heart attack analysis** dataset provides a comprehensive set of attributes for heart disease prediction, allowing for in-depth analysis of various factors contributing to heart health. The heart disease dataset contains 303 rows and 14 columns. The problem at hand is to classify whether a patient has heart disease or not based on various medical attributes. This is a binary classification problem using a dataset containing features and other relevant medical information.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   age         303 non-null   int64
1   sex         303 non-null   int64
2   cp          303 non-null   int64
3   trtbps      303 non-null   int64
4   chol        303 non-null   int64
5   fbs         303 non-null   int64
6   restecg     303 non-null   int64
7   thalachh    303 non-null   int64
8   exng        303 non-null   int64
9   oldpeak     303 non-null   float64
10  slp         303 non-null   int64
11  caa         303 non-null   int64
12  thall       303 non-null   int64
13  output      303 non-null   int64
dtypes: float64(1), int64(13)
memory usage: 33.3 KB
```

The **wine quality** dataset contains observations related to the physicochemical properties of red wine, including features and the target variable: wine quality. The dataset used in this analysis comprises 1,599 rows and 12 columns. Each row represents a unique wine sample with its respective measurements. The target variable, "quality," is an ordinal feature ranging from 3 to 8, indicating the wine's quality rating. This dataset provides a comprehensive view of various chemical attributes that potentially influence wine quality, allowing for in-depth exploratory data analysis and predictive modeling.

```
<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
```

2. Null Values

```
: # Check for missing values
print(data.isnull().sum())
```

```
# Check for missing values
df.isnull().sum()
```

```
age          0
sex          0
cp           0
trtbps       0
chol         0
fbs          0
restecg      0
thalachh     0
exng         0
oldpeak      0
slp          0
caa          0
thall        0
output       0
dtype: int64
```

```
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
```

Both the datasets do not contain any null or missing values, so there is no need to clean the data.

3. Dataset Description

Heart Attack dataset

```
# Display summary statistics
data.describe()
```

	age	sex	cp	trtbps	chol	fbs	restecg	thalachh	exng	oldpeak	slp	caa
count	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000
mean	54.366337	0.683168	0.966997	131.623762	246.264026	0.148515	0.528053	149.646865	0.326733	1.039604	1.399340	0.729373
std	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860	22.905161	0.469794	1.161075	0.616226	1.022606
min	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000	71.000000	0.000000	0.000000	0.000000	0.000000
25%	47.500000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000	133.500000	0.000000	0.000000	1.000000	0.000000
50%	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000	153.000000	0.000000	0.800000	1.000000	0.000000
75%	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.000000	166.000000	1.000000	1.600000	2.000000	1.000000
max	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000	202.000000	1.000000	6.200000	2.000000	4.000000
sex	cp	trtbps	chol	fbs	restecg	thalachh	exng	oldpeak	slp	caa	thall	output
000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000
683168	0.966997	131.623762	246.264026	0.148515	0.528053	149.646865	0.326733	1.039604	1.399340	0.729373	2.313531	0.544554
.466011	1.032052	17.538143	51.830751	0.356198	0.525860	22.905161	0.469794	1.161075	0.616226	1.022606	0.612277	0.498835
000000	0.000000	94.000000	126.000000	0.000000	0.000000	71.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
000000	0.000000	120.000000	211.000000	0.000000	0.000000	133.500000	0.000000	0.000000	1.000000	0.000000	2.000000	0.000000
000000	1.000000	130.000000	240.000000	0.000000	1.000000	153.000000	0.000000	0.800000	1.000000	0.000000	2.000000	1.000000
000000	2.000000	140.000000	274.500000	0.000000	1.000000	166.000000	1.000000	1.600000	2.000000	1.000000	3.000000	1.000000
000000	3.000000	200.000000	564.000000	1.000000	2.000000	202.000000	1.000000	6.200000	2.000000	4.000000	3.000000	1.000000

Wine Quality dataset

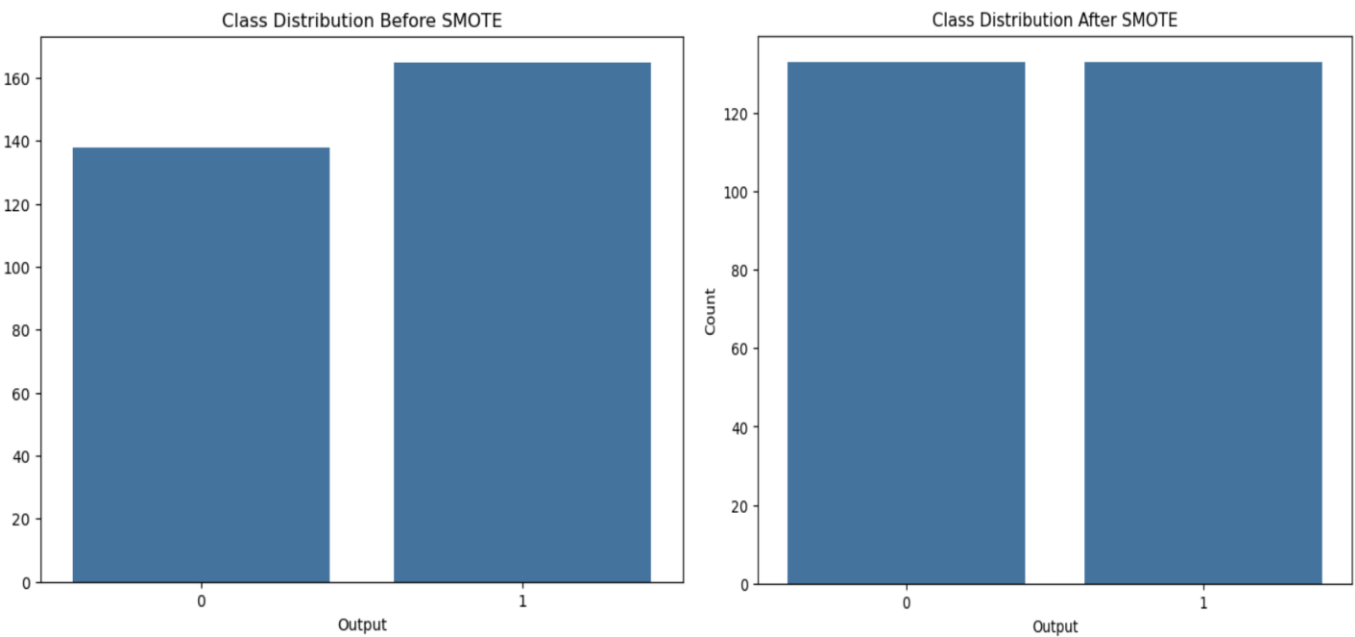
```
# Display summary statistics
df.describe()
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	
count	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000
mean	8.319637	0.527821	0.270976	2.538806	0.087467	15.874922	46.467792	0.996747	3.311113	0.658149	10.422983	5.636023
std	1.741096	0.179060	0.194801	1.409928	0.047065	10.460157	32.895324	0.001887	0.154386	0.169507	1.065668	0.807569
min	4.600000	0.120000	0.000000	0.900000	0.012000	1.000000	6.000000	0.990070	2.740000	0.330000	8.400000	3.000000
25%	7.100000	0.390000	0.090000	1.900000	0.070000	7.000000	22.000000	0.995600	3.210000	0.550000	9.500000	5.000000
50%	7.900000	0.520000	0.260000	2.200000	0.079000	14.000000	38.000000	0.996750	3.310000	0.620000	10.200000	6.000000
75%	9.200000	0.640000	0.420000	2.600000	0.090000	21.000000	62.000000	0.997835	3.400000	0.730000	11.100000	6.000000
max	15.900000	1.580000	1.000000	15.500000	0.611000	72.000000	289.000000	1.003690	4.010000	2.000000	14.900000	8.000000

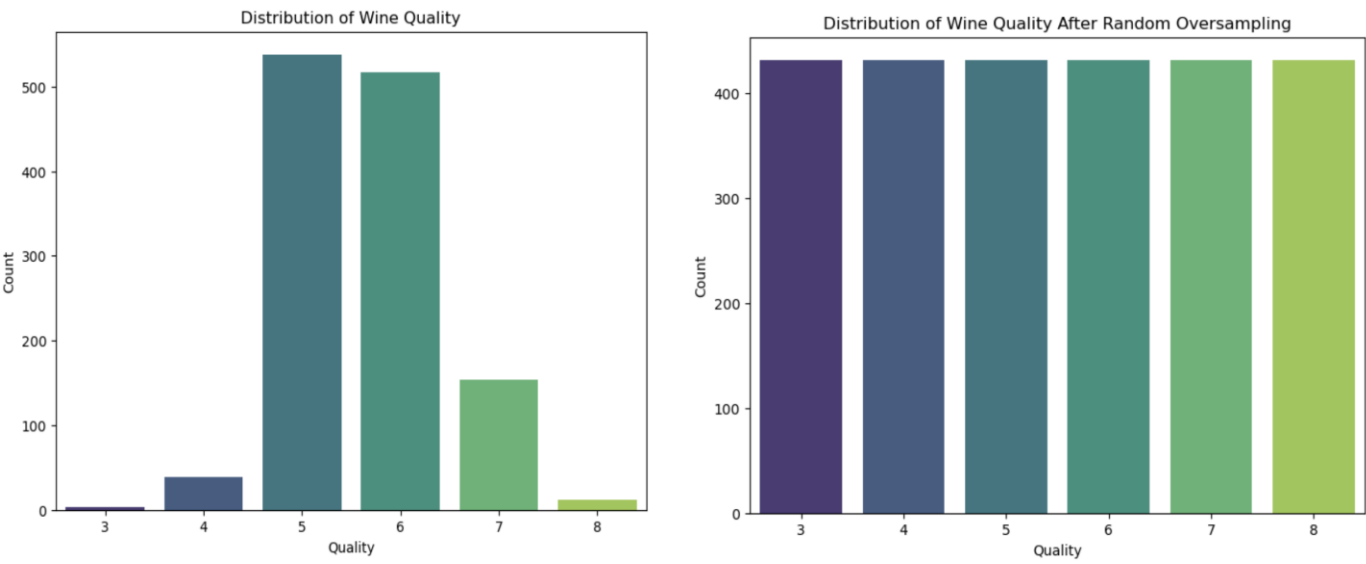
	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
count	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000
mean	8.319637	0.527821	0.270976	2.538806	0.087467	15.874922	46.467792	0.996747	3.311113	0.658149	10.422983	5.636023
std	1.741096	0.179060	0.194801	1.409928	0.047065	10.460157	32.895324	0.001887	0.154386	0.169507	1.065668	0.807569
min	4.600000	0.120000	0.000000	0.900000	0.012000	1.000000	6.000000	0.990070	2.740000	0.330000	8.400000	3.000000
25%	7.100000	0.390000	0.090000	1.900000	0.070000	7.000000	22.000000	0.995600	3.210000	0.550000	9.500000	5.000000
50%	7.900000	0.520000	0.260000	2.200000	0.079000	14.000000	38.000000	0.996750	3.310000	0.620000	10.200000	6.000000
75%	9.200000	0.640000	0.420000	2.600000	0.090000	21.000000	62.000000	0.997835	3.400000	0.730000	11.100000	6.000000
max	15.900000	1.580000	1.000000	15.500000	0.611000	72.000000	289.000000	1.003690	4.010000	2.000000	14.900000	8.000000

4. Dataset Imbalance

In this analysis of the heart disease dataset, I encountered a significant class imbalance issue. Upon examination of the target variable distribution, I found that the number of patients without heart disease (class 0) was considerably higher than those with heart disease (class 1). To address this issue, I implemented the **Synthetic Minority Over-sampling Technique (SMOTE)**. SMOTE is an oversampling method that creates synthetic examples of the minority class to balance the dataset. It effectively addressed the class imbalance, I was cautious to avoid potential overfitting and ensured a more equitable treatment of both classes, which is particularly important in medical diagnostics where false negatives can have serious consequences.

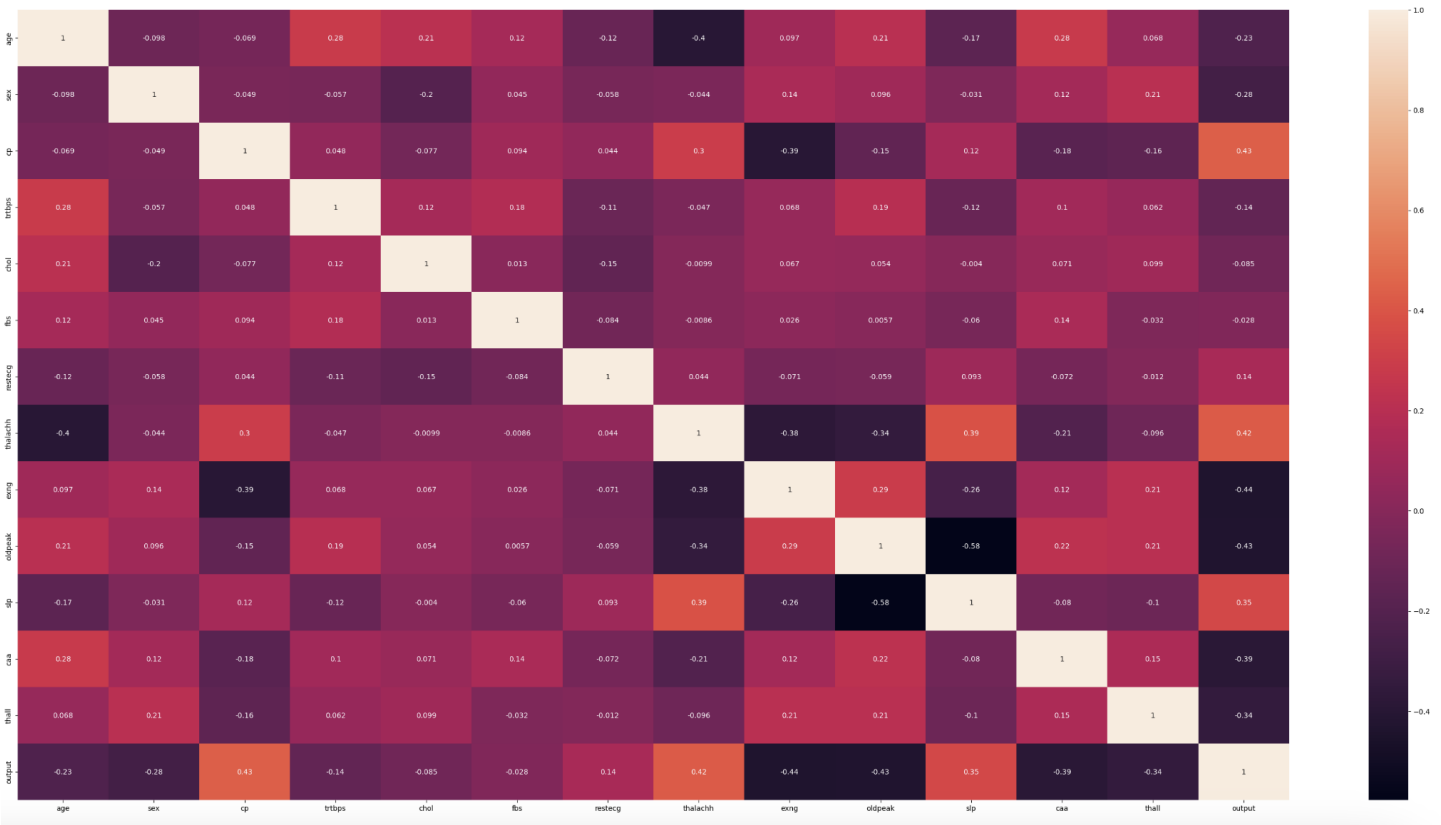


The wine quality dataset exhibited a significant class imbalance, as illustrated in the distribution chart of wine quality. Most samples were concentrated around quality scores of 5 and 6, with fewer instances of other ratings like 3, 4, 7, and 8. This imbalance can skew model training, causing it to favor the majority classes and potentially overlook the minority ones. To address this issue, I employed **Random Oversampling**. This technique duplicates samples from the minority classes to create a more balanced dataset. By doing so, I ensured that each class was represented more equally during model training, which helps improve the model's ability to generalize across all quality ratings.



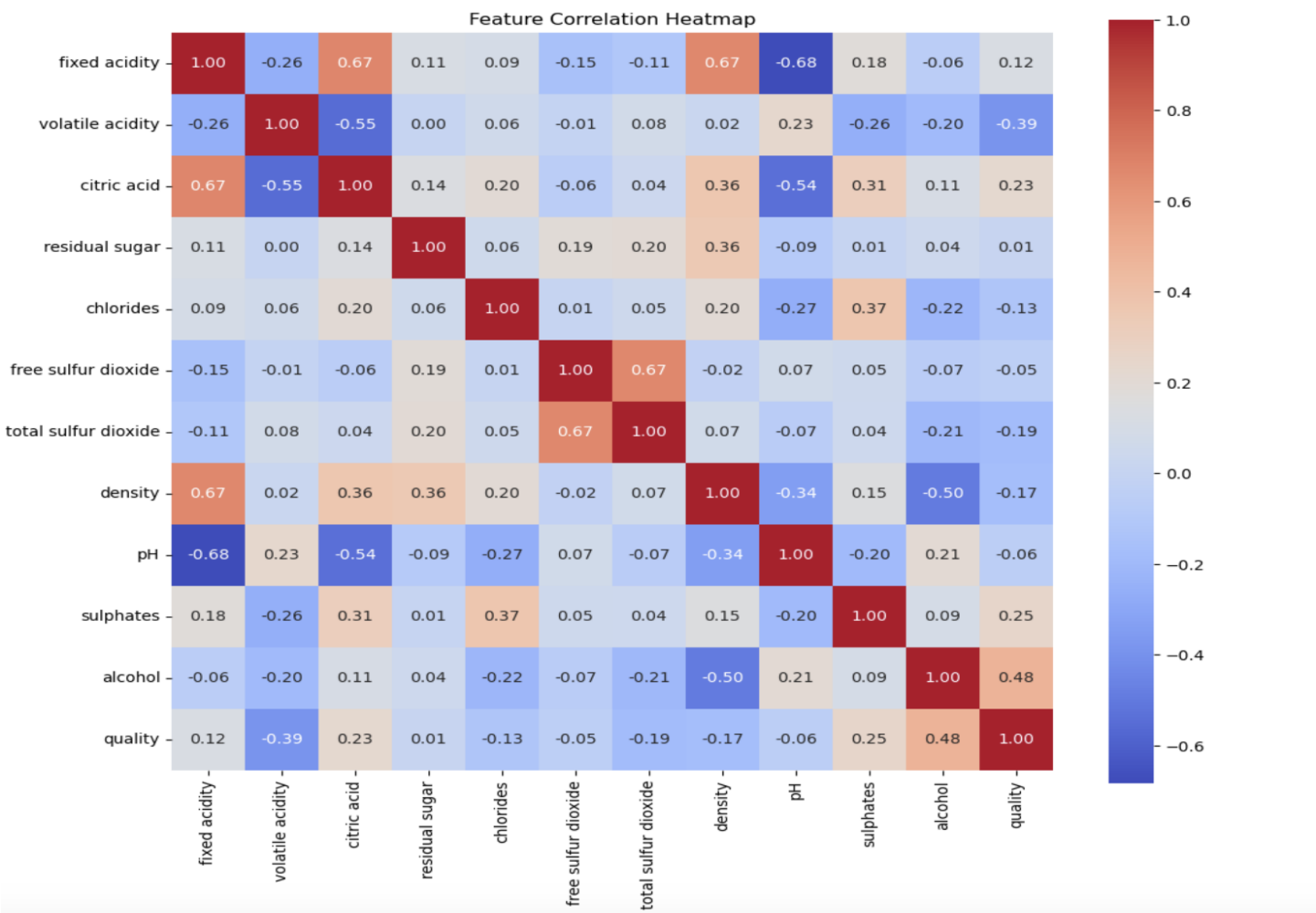
5. Correlation Matrices

I generated a correlation heatmap to visualize the relationships between different features in our heart disease dataset. A strong positive correlation was observed between 'chest pain' and the target variable, indicating that certain types of chest pain are closely associated with the presence of heart disease. 'Max heart rate' showed a moderate negative correlation with the target variable, suggesting that lower maximum heart rates might be indicative of heart disease. 'ST depression' exhibited a positive correlation with the target, implying that higher ST depression values are associated with an increased likelihood of heart disease. These correlations provided valuable insights into which features might be most predictive of heart disease, guiding our feature selection process and helping us understand the underlying patterns in the data.



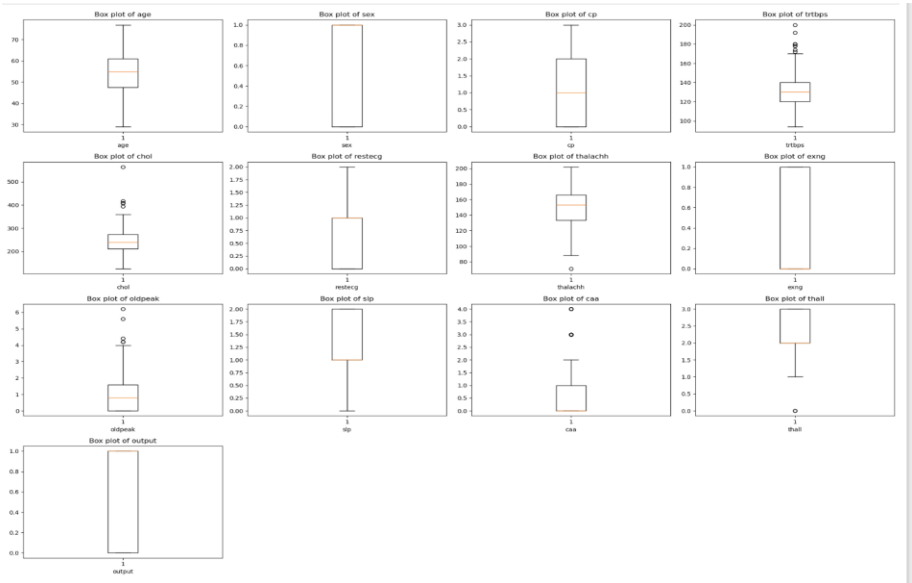
The correlation heatmap provides a visual representation of the relationships between different features in the dataset. It highlights how each feature correlates with wine quality and with each other. Notably, alcohol content

shows a strong positive correlation with quality, suggesting its importance in predicting higher quality wines. Volatile acidity has a negative correlation, indicating that higher acidity might detract from perceived quality. Features like residual sugar and chlorides exhibit low correlation with quality, suggesting they may have less impact on the model's predictive power.

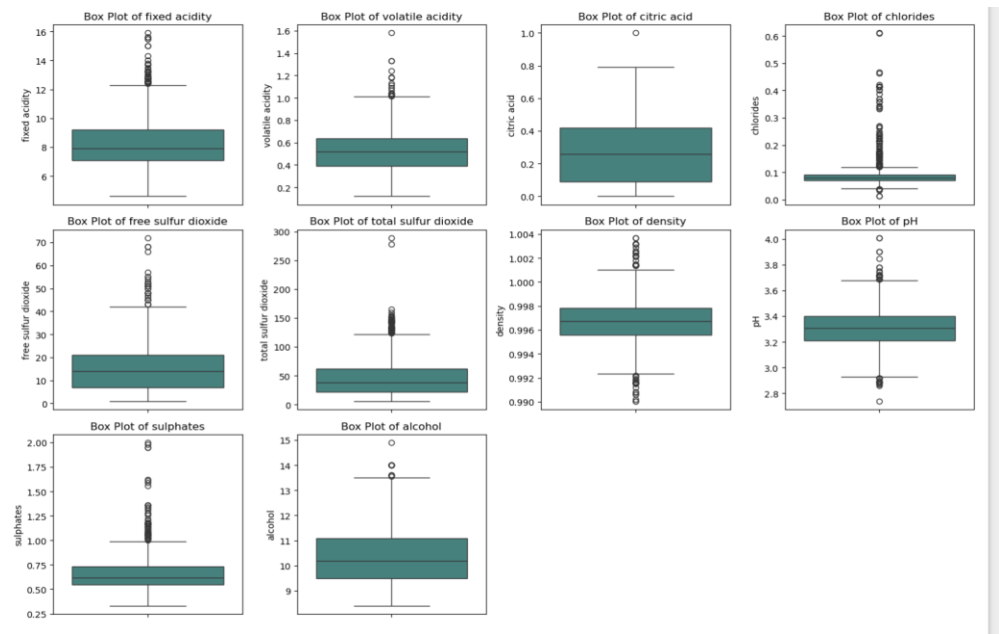


6. Outliers

To identify potential outliers in heart disease dataset, I employed box plots for each numerical feature. This analysis revealed the 'cholesterol' feature contained several high outliers, with some values significantly above the upper quartile. 'Age' and 'resting blood pressure' also showed a few outliers, but these were less extreme compared to cholesterol. Given the medical nature of our data, I decided to retain these outliers as they could represent genuine, rare, cases that are important for our model to learn from. I made note of these outliers to monitor their potential impact on our model's performance.



Outlier detection was performed using box plots for wine quality dataset, revealing extreme values in features such as residual sugar and sulphates. These outliers can skew analysis and model performance if not addressed. By identifying and potentially removing these outliers, I aimed to improve model robustness. I have used Inter Quartile Range to detect the outliers. I will keep the outliers within  $1.5 \times (\text{lower bound of IQR}) - 1 \times 5 (\text{upper bound of IQR})$  range test.

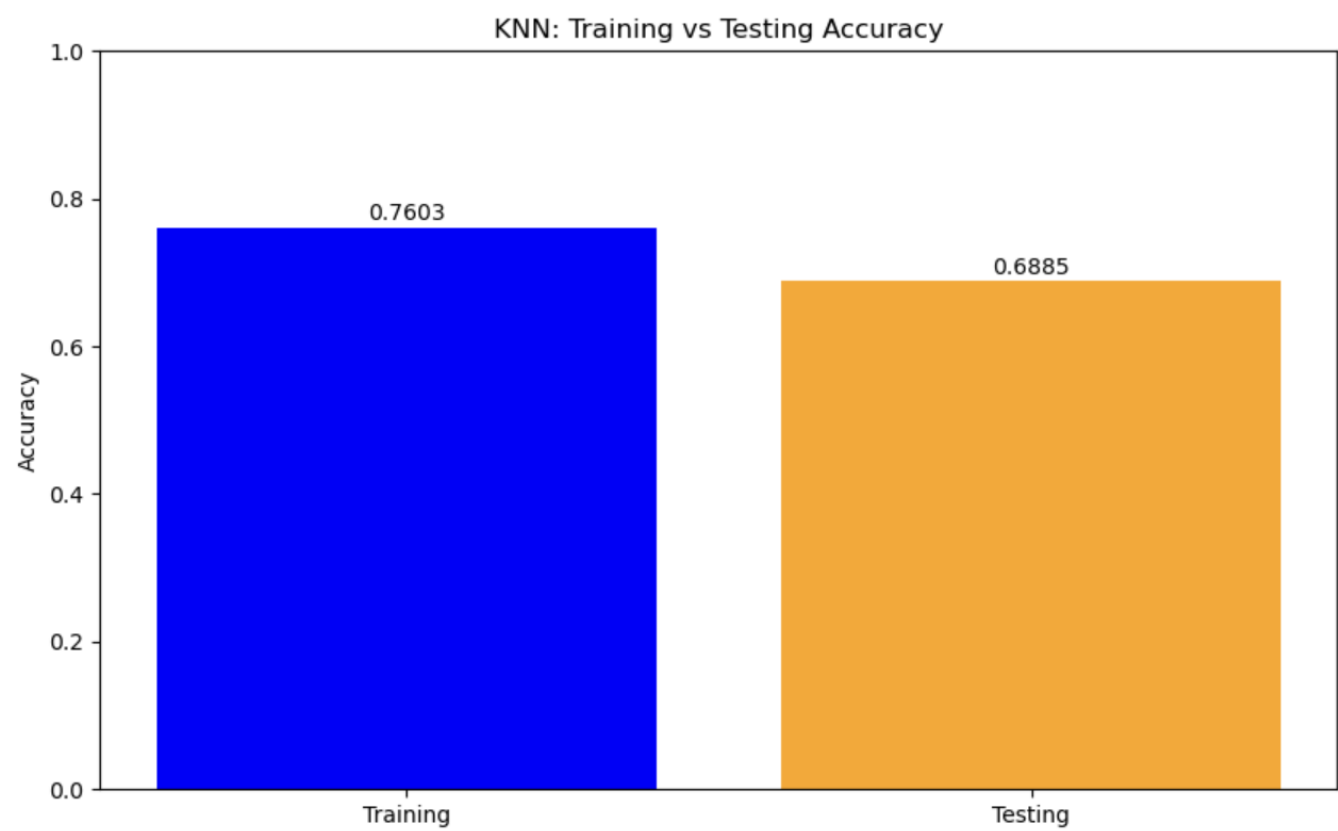


7. Why are these datasets interesting?

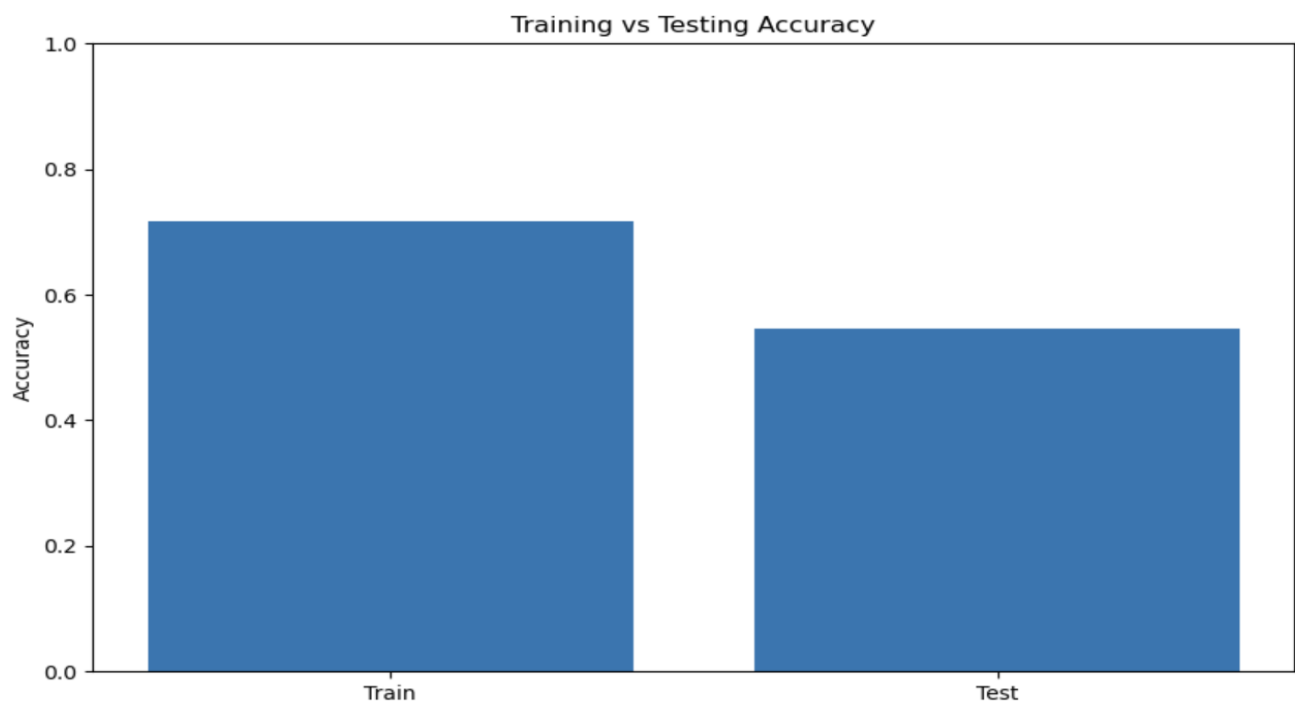
The dataset I have analyzed presents a compelling case study that bridges the gap between medical science and data analytics. Its real-world implications in potentially aiding early diagnosis of a leading cause of mortality worldwide make it particularly significant. The wine quality dataset is particularly interesting due to its practical application in the wine industry and its potential to reveal insights into the factors influencing wine quality. By analyzing physicochemical properties such as acidity, sugar content, and alcohol levels, I can identify key attributes that contribute to higher quality ratings. The challenge of addressing class imbalance and outliers provides a rich opportunity for applying advanced data preprocessing and machine learning techniques. This makes the dataset a valuable resource for both academic research and industry applications.

8. KNN Algorithm

In heart dataset, the K-Nearest Neighbors (KNN) algorithm was implemented using the KNeighborsClassifier from the sklearn library, with the number of neighbors set to 5. The performance of the model was evaluated using several metrics: accuracy, precision, recall, and F1 score, all calculated using weighted averages to account for class imbalances. The results showed a training accuracy of 76.03% and a testing accuracy of 68.85%, indicating that the model performs better on the training data than on unseen data. Precision, recall, and F1 scores followed a similar pattern, with values of 76.05%, 76.03%, and 75.85% for training, and 68.88%, 68.85%, and 68.70% for testing, respectively. These metrics suggest that while the model has learned patterns from the training data reasonably well, it may not generalize as effectively to new data, potentially due to overfitting or insufficient feature representation.

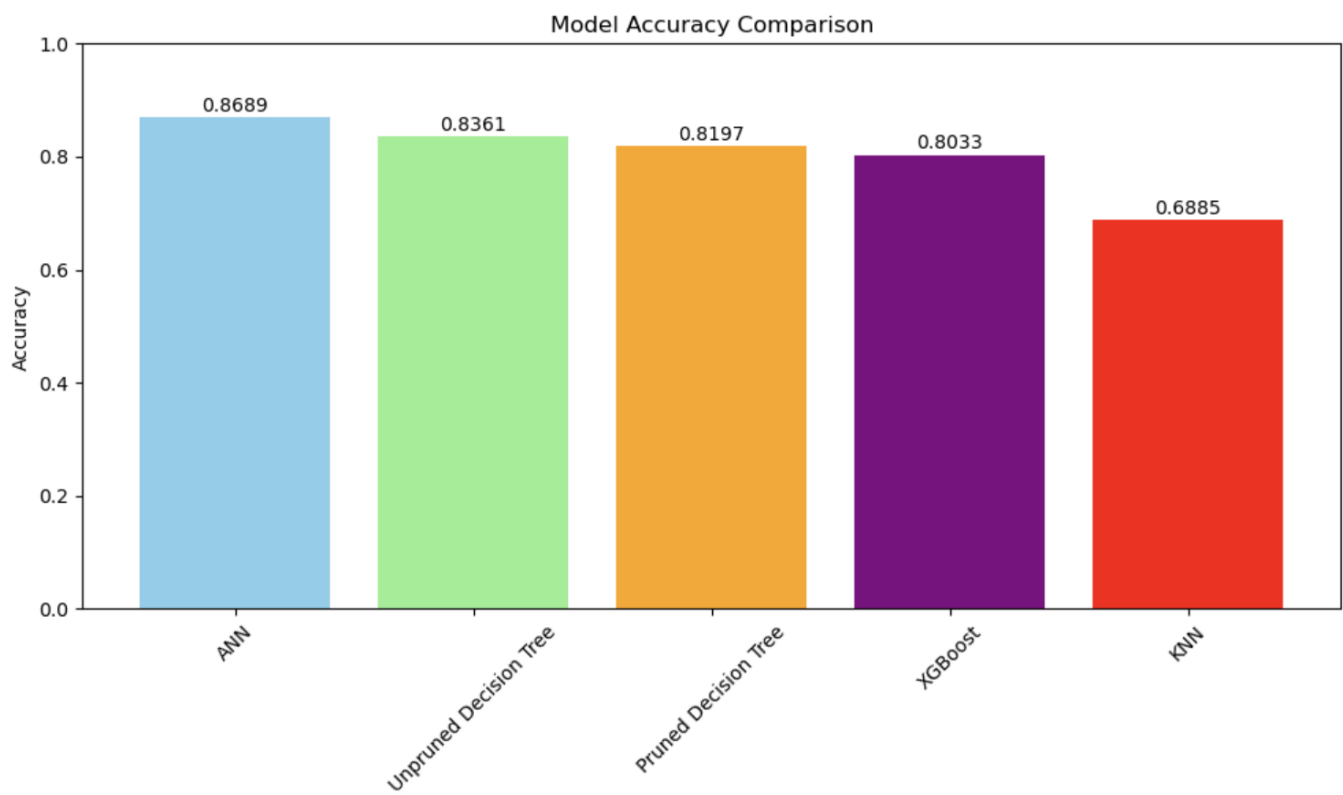


In wine dataset, I implemented a K-Nearest Neighbors (KNN) algorithm to classify wine quality based on various features. The dataset was preprocessed by separating the features from the target variable, 'quality'. The KNN model was initialized with 5 neighbors, though this parameter can be adjusted for optimization. After training the model on the scaled training data, predictions were made for both training and testing datasets. The model's performance was evaluated using several metrics: train accuracy was 71.77%, while test accuracy was 54.69%, indicating potential overfitting. Precision, recall, and F1 score were calculated as 52.24%, 54.69%, and 53.09% respectively, using a weighted average to account for class imbalances. These results suggest that while the model performs moderately well on the training data, its generalization to unseen data is limited, highlighting areas for potential improvement such as feature selection, hyperparameter tuning, or exploring different algorithms.



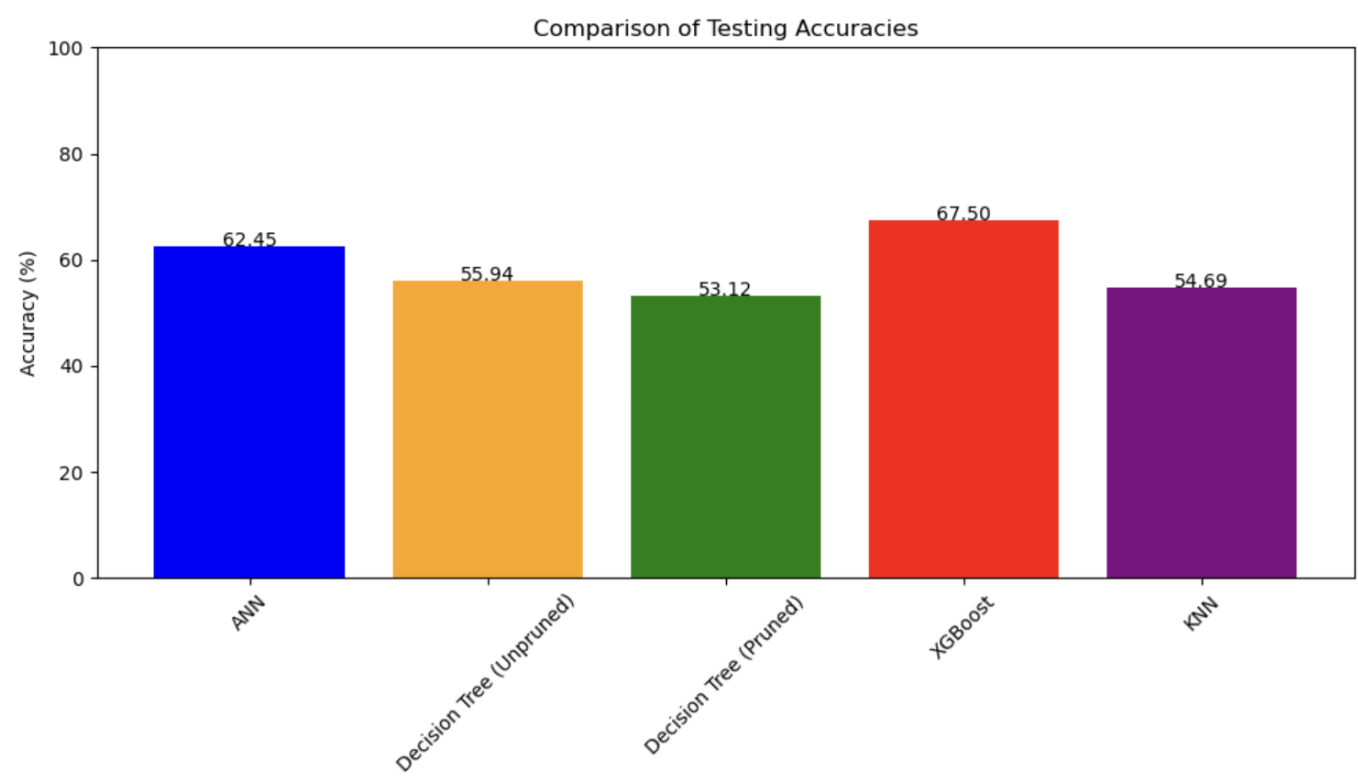
9. Conclusion

For heart dataset, the Artificial Neural Network (ANN) achieved the highest training accuracy at 86.89%, indicating its strong capability to capture complex patterns in the data, making it ideal for tasks requiring high model complexity and non-linear relationships. The Unpruned Decision Tree demonstrated a test accuracy of 83.61%, suggesting its robustness in handling datasets with clear decision boundaries without overfitting. The Pruned Decision Tree, with a slightly lower test accuracy of 81.97%, offers a more generalized model by reducing complexity, which is beneficial when interpretability and simplicity are prioritized. XGBoost, a boosting algorithm, achieved a training accuracy of 80.33%, showcasing its effectiveness in improving model performance through ensemble learning and is particularly useful in competitions or scenarios where incremental performance gains are crucial. Lastly, the K-Nearest Neighbors (KNN) algorithm showed a test accuracy of 68.85%, indicating its suitability for smaller datasets or when computational simplicity is desired, despite its lower performance compared to other models in this context.





For wine quality dataset, XGBoost achieved the highest accuracy at 67.50%, indicating its strong capability in handling complex datasets and capturing intricate patterns. ANN followed with an accuracy of 62.45%, showcasing its effectiveness in modeling non-linear relationships. The unpruned Decision Tree had an accuracy of 55.94%, slightly outperforming the pruned version at 53.12%, suggesting that while pruning can prevent overfitting, it might also reduce model complexity too much for this dataset. KNN, with an accuracy of 54.69%, demonstrated moderate performance, which may be attributed to its sensitivity to feature scaling and choice of neighbors. XGBoost is recommended for tasks requiring high accuracy and robustness, while ANN is suitable for problems involving complex feature interactions. Decision Trees can be useful for interpretability and simplicity, whereas KNN might be preferred for smaller datasets or when model simplicity is crucial.



10. References

- <https://www.kaggle.com/datasets/rashikrahmanpritom/heart-attack-analysis-prediction-dataset>
- <https://labelyourdata.com/articles/machine-learning-for-wine-quality-prediction#>
- <https://archive.ics.uci.edu/>
- [ChatGPT](#)