Datesets :

# 1-Pima Indians Diabetes Database

- **Number of Classes**: There are two classes in this dataset:
- 0: Represents individuals without diabetes
- 1: Represents individuals with diabetes
- **Labels**: The dataset includes the following attributes (features) along with the class label:

- 
  1. Number of times pregnant

2. Plasma glucose concentration a 2 hours in an oral glucose tolerance test
3. Diastolic blood pressure (mm Hg)
4. Triceps skinfold thickness (mm)
5. 2-Hour serum insulin (mu U/ml)
6. Body mass index (weight in kg/(height in m)^2)
7. Diabetes pedigree function
8. Age (years)
9. Class variable (0 or 1)

**Total Number of Samples**: The dataset contains 768 samples/instances.

# 2-Boston Housing

- **Number of Classes**: This dataset is typically used for regression tasks, so it doesn't have distinct classes like classification datasets. Instead, the target variable is usually the median value of owner-occupied homes in $1000s (MEDV), which is a continuous variable representing house prices.

- **Labels**: The dataset includes the following features:

1. CRIM: Per capita crime rate by town

2. ZN: Proportion of residential land zoned for lots over 25,000 sq. ft.

.3 INDUS: Proportion of non-retail business acres per town

.4 CHAS: Charles River dummy variable (1 if tract bounds river; 0 otherwise)

.5 NOX: Nitric oxides concentration (parts per 10 million)

.6 RM: Average number of rooms per dwelling

.7 AGE: Proportion of owner-occupied units built prior to 1940

.8 DIS: Weighted distances to five Boston employment centers

.9 RAD: Index of accessibility to radial highways

.10 TAX: Full-value property tax rate per $10,000

11. PTRATIO: Pupil-teacher ratio by town

12. B: 1000(Bk - 0.63)^2 where Bk is the proportion of Black residents by town

13. LSTAT: Percentage of lower status of the population

14. MEDV: Median value of owner-occupied homes in $1000s (target variable)

- **Total Number of Samples**: The dataset contains 506 samples/instances.

- **Size of Each Sample**: Each sample in this dataset consists of 13 input features (attributes) and the target

variable (median value of owner-occupied homes).

## Raisin binary classification-3

- **Number of Classes**: There are two classes in this dataset:

- 0: Indicates non-defective raisins

- 1: Indicates defective raisins

- **Labels**: The dataset includes the following attributes along with the class label:

1. **Area** - The area of the raisin

2. **Perimeter** - The perimeter of the raisin

3. **MajorAxisLength** - The major axis length of the raisin

4. **MinorAxisLength** - The minor axis length of the raisin

5. **AspectRation** - The aspect ratio of the raisin

6. **Eccentricity** - The eccentricity of the raisin

7. **ConvexArea** - The convex area of the raisin

8. **EquivDiameter** - The equivalent diameter of the raisin

9. **Extent** - The extent of the raisin

10. **Solidity** - The solidity of the raisin.

11. **Roundness** - The roundness of the raisin

12. **Compactness** - The compactness of the raisin

**ShapeFactor1** - Shape factor 1 of the .13 raisin

**ShapeFactor2** - Shape factor 2 of the .14 raisin

**ShapeFactor3** - Shape factor 3 of the .15 raisin

**ShapeFactor4** - Shape factor 4 of the .16 raisin

**Class** - Class variable (0 or 1).17

**Total Number of Samples**: The . dataset contains a total of 106 samples/instances.

**Size of Each Sample**: Each sample in . this dataset contains 17 attributes including the class label.

# Decision Tree Model:

## Feature Extraction Details:

- **Number of Features Extracted:** 7 features were extracted.

- **Feature Names:**

  - 'pregnant'

  - 'insulin'

  - 'bmi'

  - 'age'

  - 'glucose'

  - 'bp'

  - 'pedigree'

- **Dimension of Resulted Features:** Each instance has 7 features.

# Results Details (on Testing Data):

## Accuracy with Hyperparameter .
## Tuning: 77%

**Cross-Validation:**

- **Implemented Model:** Decision Tree Classifier

- **Cross-Validation Usage:** Yes, 10-fold cross-validation was used during hyperparameter tuning.

- **Ratio of Training/Validation:** 9:1 (90% training, 10% validation)

# Confusion Matrix:

# Feature Importance:



Feature Importances

# Decision tree:

# SVR Model :

**Feature Extraction Details:**

1. **How many features were extracted?**
   - The feature extraction process included selecting the top 5 features based on the Recursive Feature Elimination (RFE) technique. These features are:
     - CRIM (per capita crime rate by town)
     - RM (average number of rooms per dwelling)
     - DIS (weighted distances to five Boston employment centers)
     - PTRATIO (pupil-teacher ratio by town)
     - LSTAT (% lower status of the population)

2. **Cross-Validation Usage:**
   - Yes, cross-validation was used in the implemented Support Vector Regression (SVR) model.
   - Number of folds: 5 (specified in the KFold method)
   - Ratio of training/validation: Each fold used 80% for training and 20% for validation.

3. **Hyperparameters Used:**
   - Initial Learning Rate: Not explicitly mentioned in the provided code.
   - Optimizer: Support Vector Regression (SVR) with a Radial Basis Function (RBF) kernel.
   - Regularization: Controlled by the hyperparameters `C` and `epsilon` in SVR.
   - Batch Size: Not applicable to SVR as it's not trained in batches like neural networks.
   - Number of Epochs: 8 epochs for each fold (as per the code), totaling 40 iterations (5 folds * 8 epochs).
   - Other Hyperparameters: Gamma value (0.1), kernel type (RBF), and epsilon (1.0) were selected through GridSearchCV.
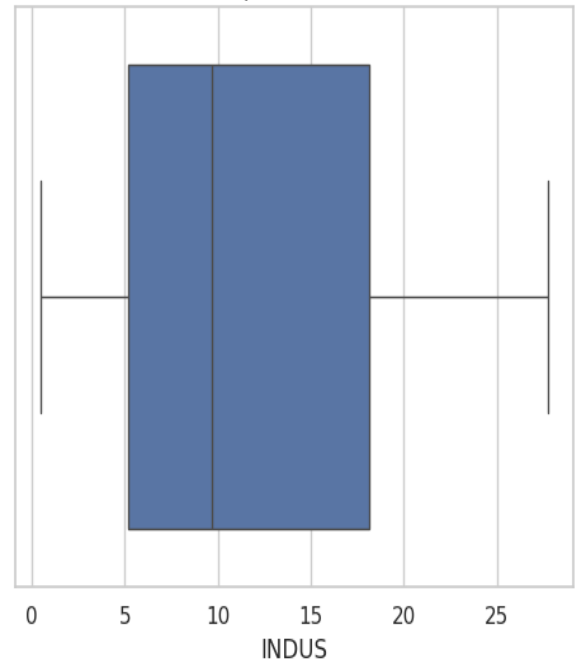
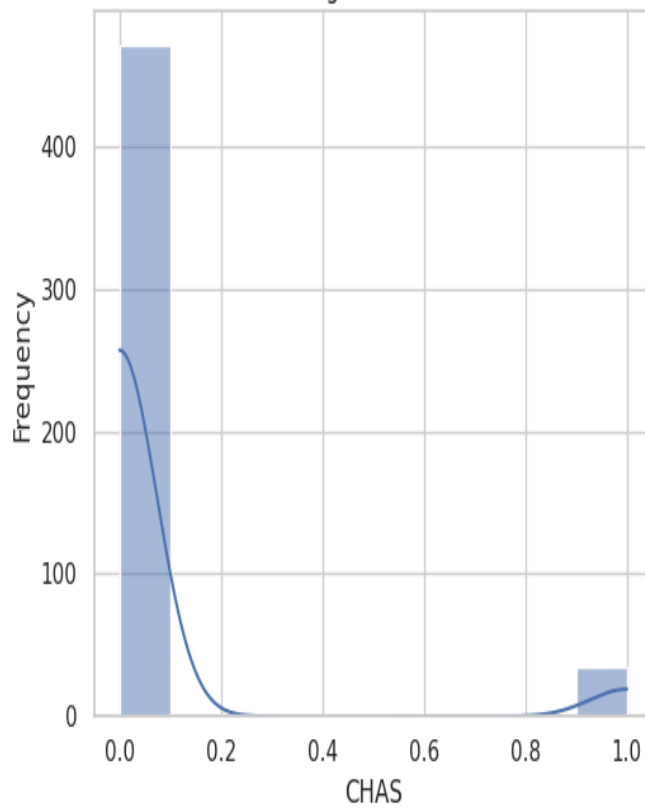# histogram & boxplot to check for outliers visually:
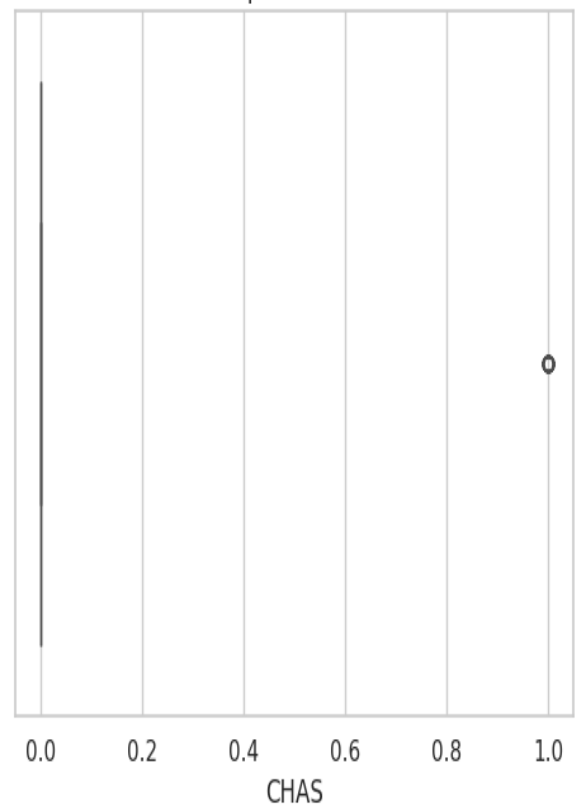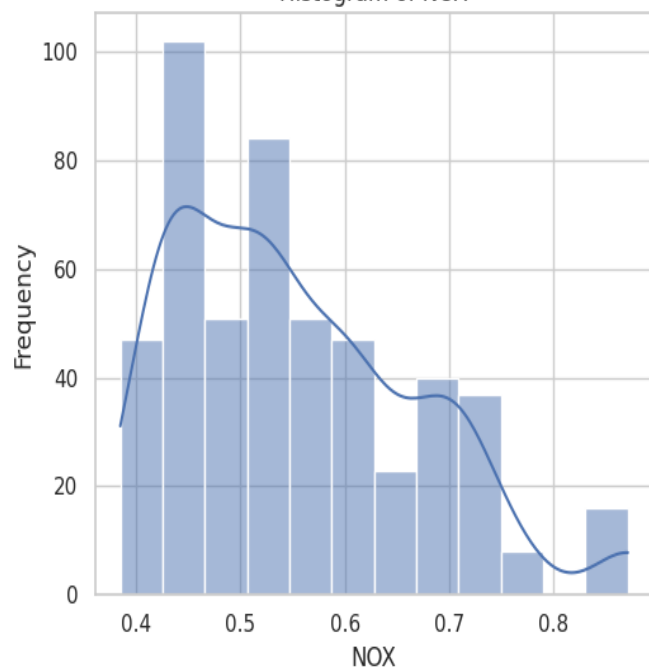
Histogram of INDUS

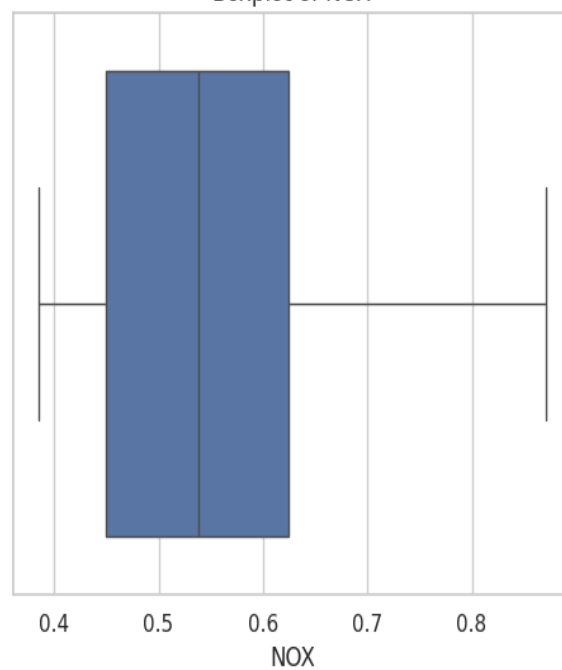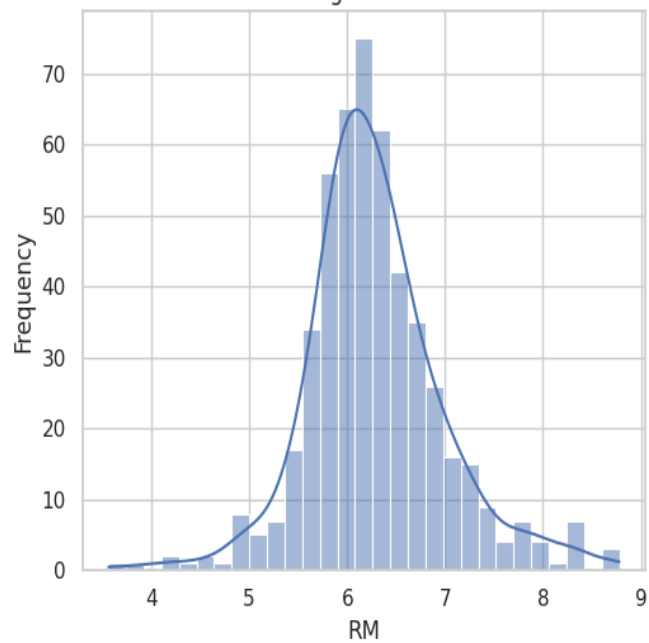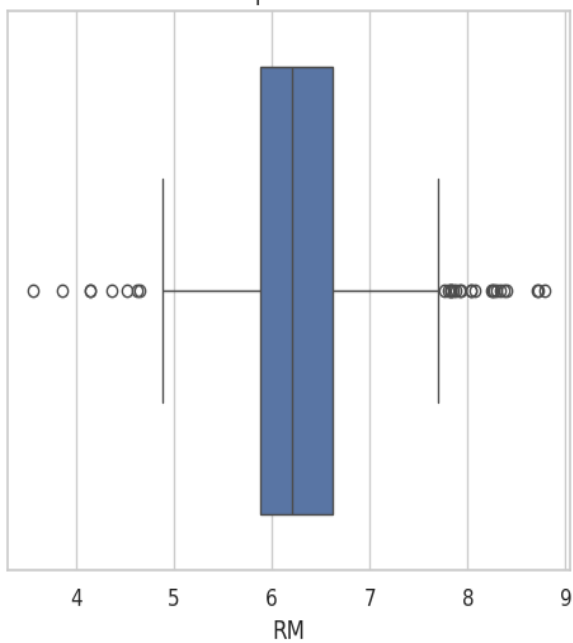Boxplot of INDUS

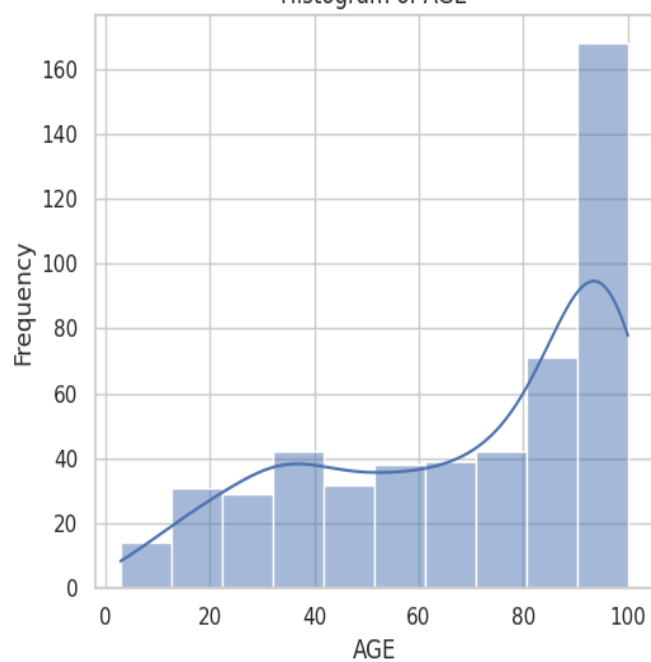Histogram of CHAS

Boxplot of CHAS

Histogram of NOX · Boxplot of NOX · Histogram of RM · Boxplot of RM
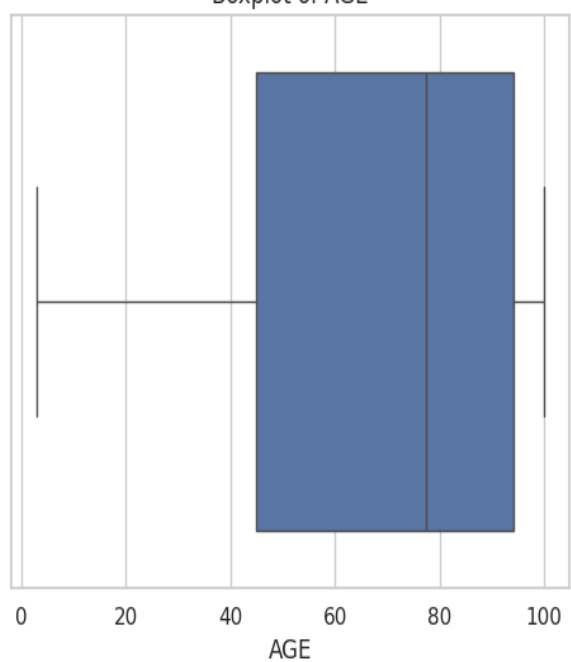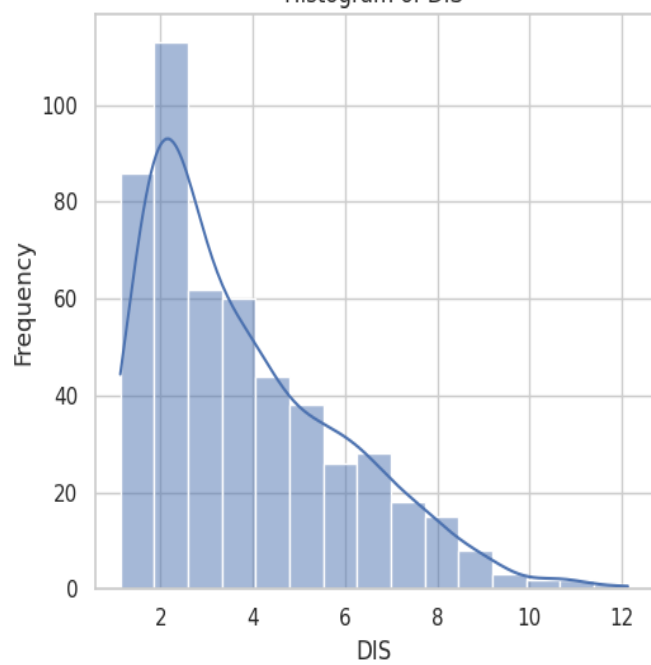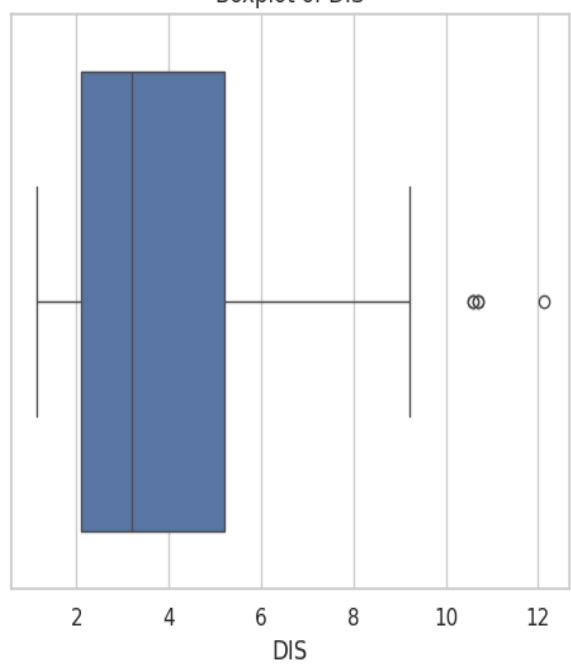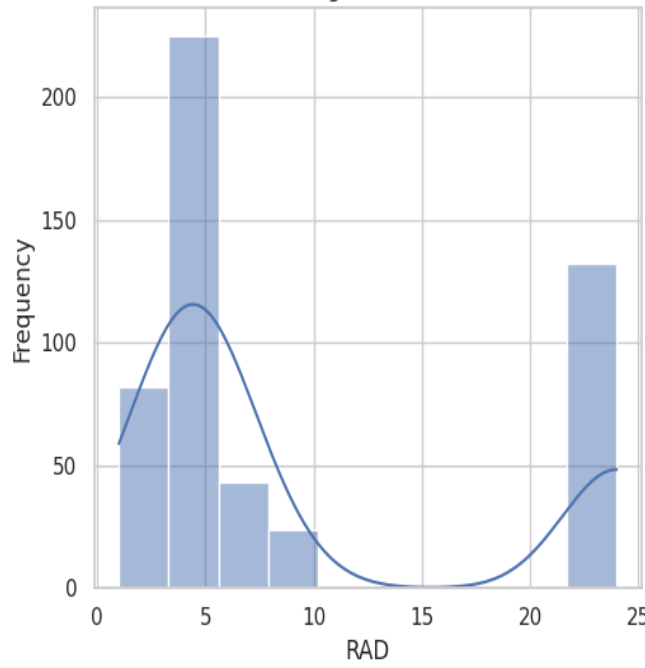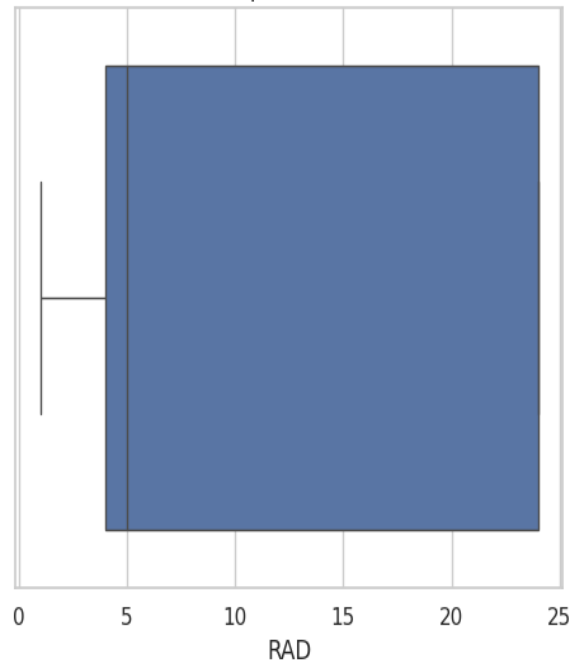
Histogram of AGE

Boxplot of AGE

Histogram of DIS

Boxplot of DIS

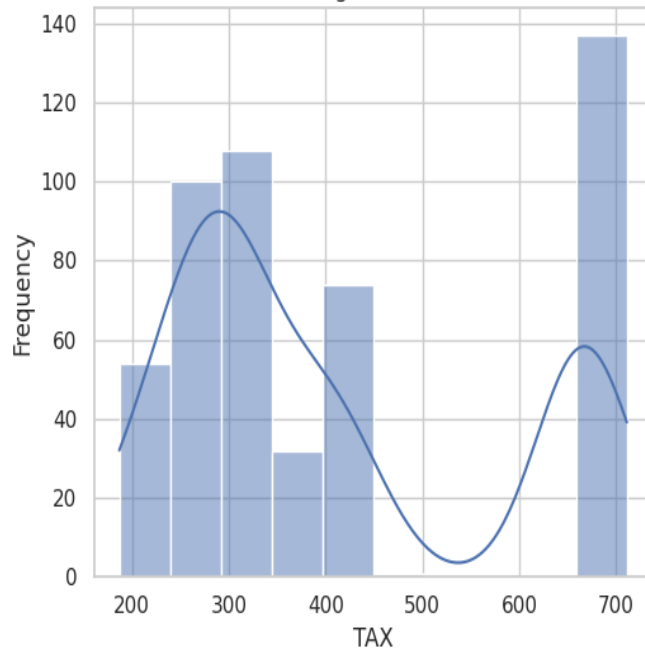Histogram of RAD
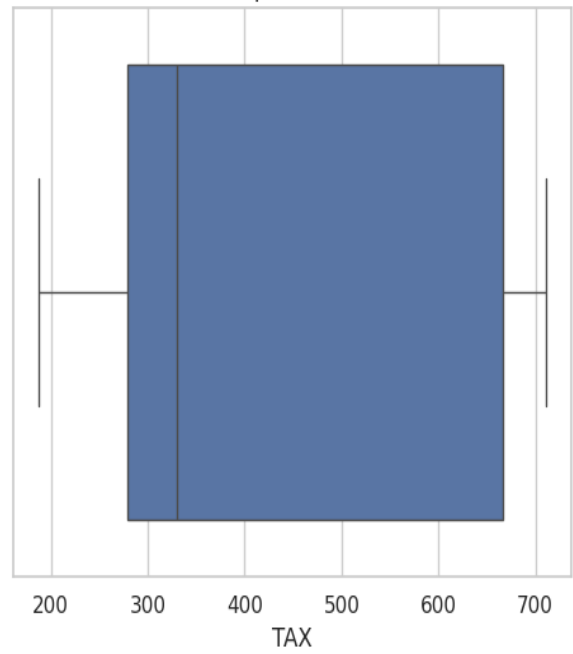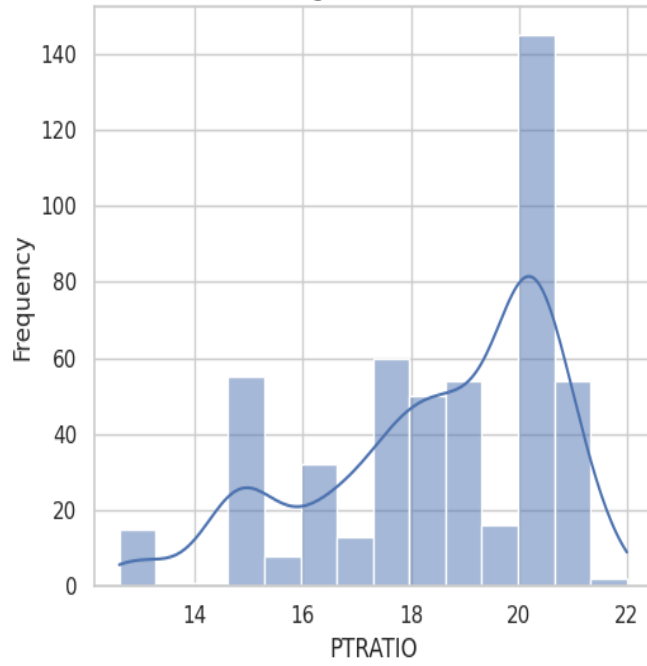
Boxplot of RAD

Histogram of TAX
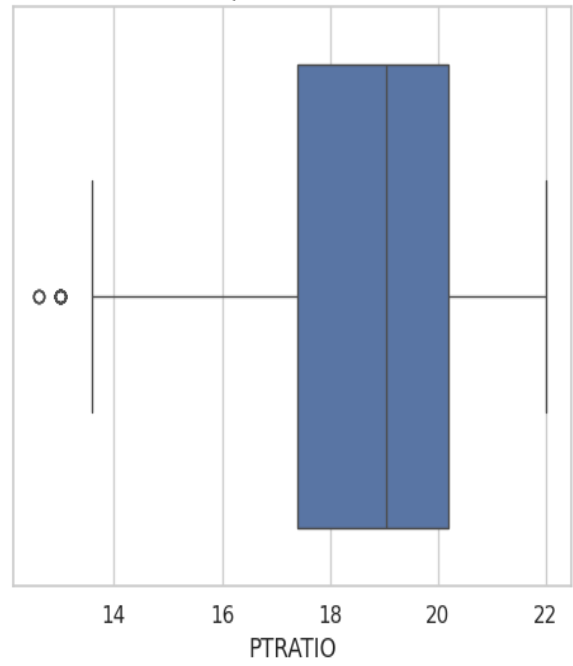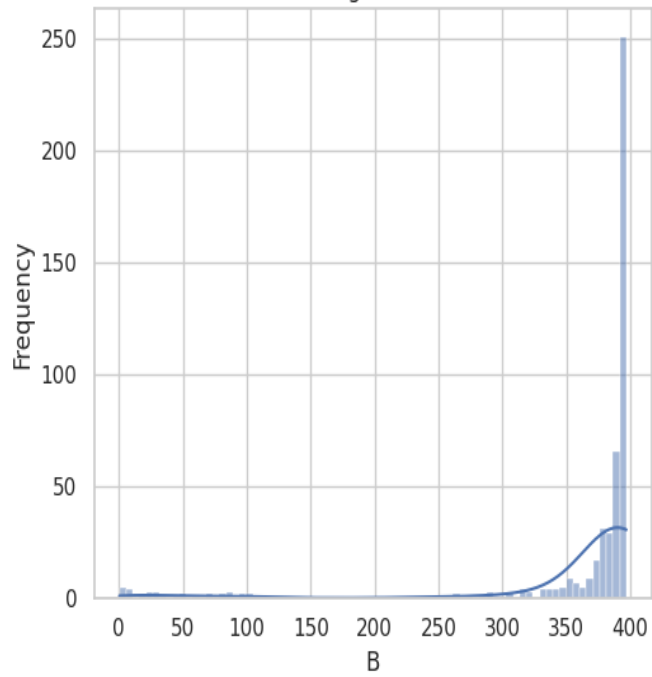
Boxplot of TAX

Histogram of PTRATIO
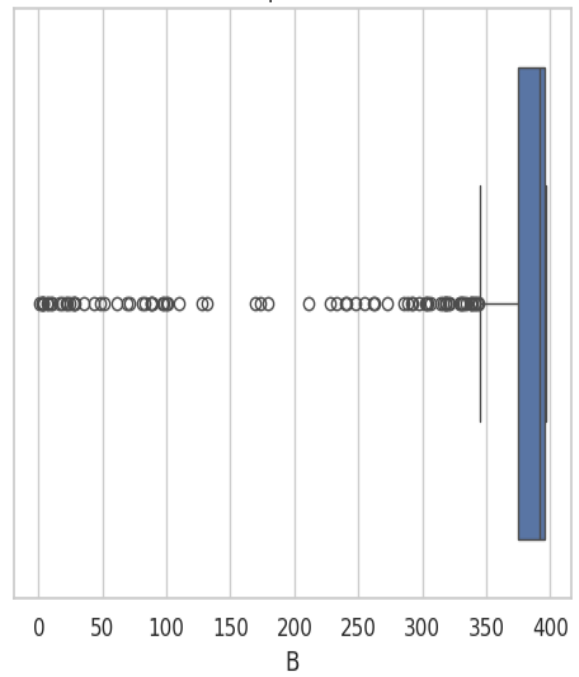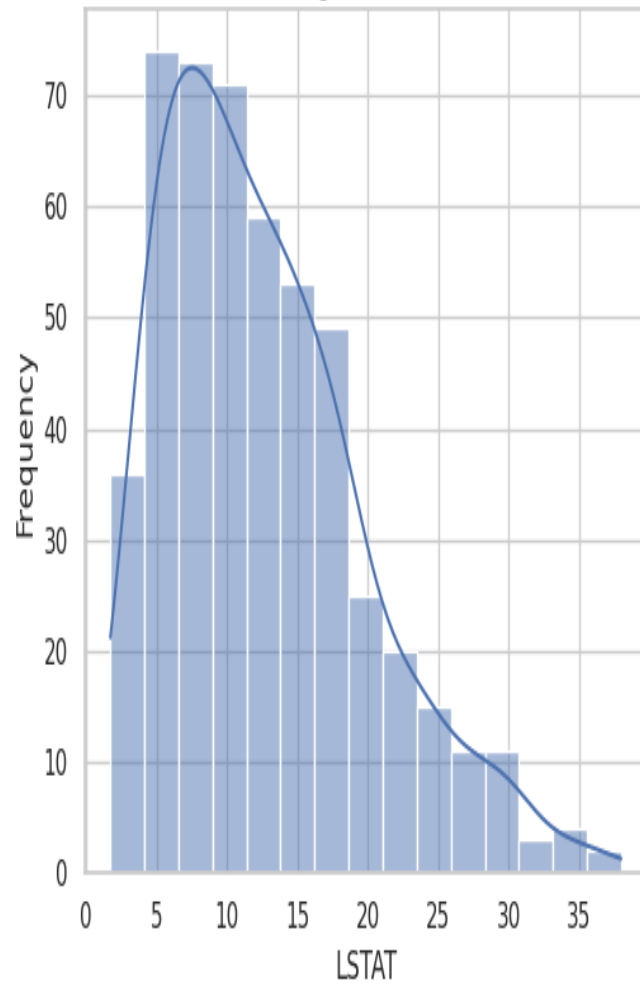
Boxplot of PTRATIO

Histogram of B

Boxplot of B

Histogram of LSTAT — Boxplot of LSTAT

# NN Model:

## Features extracted details:

**1. Area (not used)**

Gives the number of pixels within the boundaries of the raisin.

**2. MajorAxisLength**

Gives the length of the main axis, which is the longest line that

can be drawn on the raisin.

**3. MinorAxisLength**

Gives the length of the small axis, which is the shortest line that

can be drawn on the raisin.

**4. Eccentricity**

It gives a measure of the eccentricity of the ellipse, which has the

same moments as raisins.

**5. ConvexArea (not used)**

Gives the number of pixels of the smallest convex shell of the
region formed by the raisin.

**6. Extent**

Gives the ratio of the region formed by the raisin to the total

pixels in the bounding box.

**7. Perimeter**

It measures the environment by calculating the distance between

the boundaries of the raisin and the pixels

# Results :

## Accuracy:-

```
6/6 ──────────────────────── 0s 1ms/step - accuracy: 0.8713 - loss: 0.3517
Test accuracy: 0.88
6/6 ──────────────────────── 0s 8ms/step
Precision: 0.90
Recall: 0.89
F1 score: 0.90
```
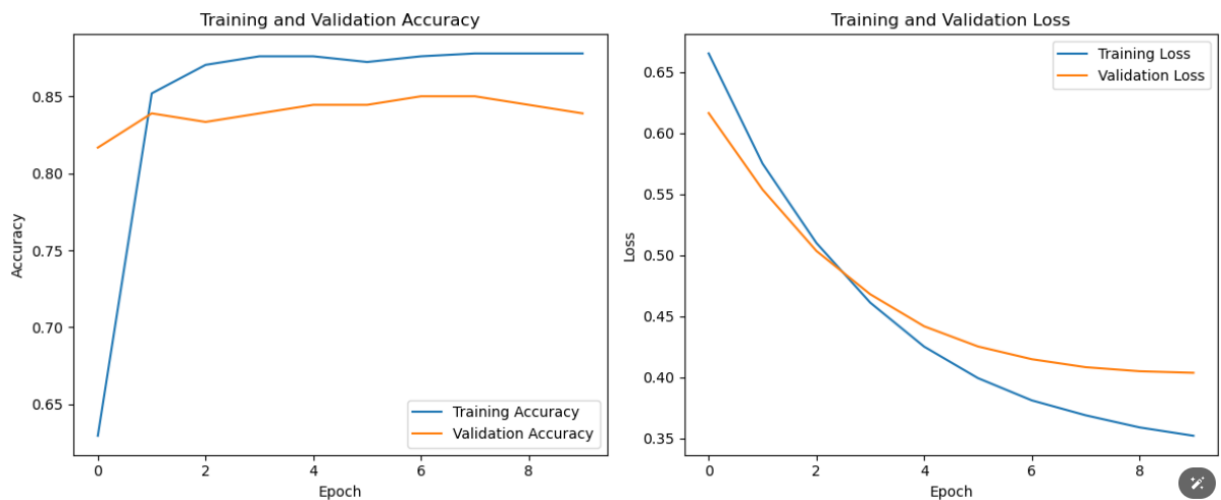
## Training vs validation:-



## Precision, Recall and F1 Score:-

## Metrics



|             | precision | recall | f1-score | support |
|-------------|-----------|--------|----------|---------|
| 0           | 0.86      | 0.87   | 0.86     | 76      |
| 1           | 0.90      | 0.89   | 0.90     | 104     |
|             |           |        |          |         |
| accuracy    |           |        | 0.88     | 180     |
| macro avg   | 0.88      | 0.88   | 0.88     | 180     |
| weighted avg | 0.88     | 0.88   | 0.88     | 180     |

Validation Size :-

yes validation was used.

```
TEST_SIZE = 0.2
VALIDATION_SIZE = 0.2
RANDOM_STATE = 42
```

# Hyperparameters:-

several hyperparameters are used:

1. TEST_SIZE: This hyperparameter determines the proportion of the dataset to include in the test split when splitting the data into training, validation, and test sets. It is set to 0.2, meaning 20% of the data will be used for testing.

2. VALIDATION_SIZE: This hyperparameter determines the proportion of the dataset to include in the validation split when splitting the data into training, validation, and test sets. It is also set to 0.2, indicating that 20% of the data will be used for validation.

3. RANDOM_STATE: This hyperparameter is used for randomizing the dataset splitting. It ensures reproducibility of results by fixing the random state. It is set to 42.

4. Number of Epochs: In the model training part, the number of epochs determines how many times the

learning algorithm will work through the entire training dataset. Here it's set to 10.
Batch Size: The batch size specifies the number of .5 samples that will be propagated through the network. Here it's set to 32.