Q1) Ans)A) A.T.Q

Let the hypothetical equation for linear regression model $y_i = mx_i + c$

$$\left[\begin{array}{l} Y_i : \text{Actual output for } i^{th} \text{ sample} \\ x_i : \text{Input feature for } i^{th} \text{ data sample} \end{array}\right]$$

Now let equation of the least square fit line

$$h_0(x) = \Theta_0 x_0 + \Theta_1 x_1 + \Theta_2 x_2 \cdots \cdots \Theta_n x_n$$

$$\left[\begin{array}{l} h_0(x) \text{ is hypothetical equation} \\ \Theta_n : \text{weights} \\ X_n : \text{Independent variable for 0 to } n \end{array}\right]$$

To prove: $\bar{y} = \Theta_0 \bar{x}_0 + \Theta_1 \bar{x}_1 + \cdots + \Theta_n \bar{x}_n$

→ if we can prove it for $n=2$ then it will be true for $n = n$ also

We know $L_2$ loss function as

$$J(\Theta_j) = \frac{1}{m} \sum_{i=1}^{m} (h_0(x^{(i)}) - y^{(i)})^2$$

To minimize the loss function we will differentiate it

For $j = \Theta_0$:

$$\frac{S(\Theta_0)}{S\Theta_0} = \frac{1}{m} \sum_{i=i}^{m} (h_0(x^{(i)}) - y^{(i)}) = 0 \quad \text{①}$$

For $j = \Theta_1$

$$\frac{S(\Theta_1)}{S\Theta_1} = \frac{1}{m} \sum_{i=1}^{m} (h_0(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} = 0 \quad \text{②}$$

Solving ①

→ $\frac{1}{m} \sum_{i=i}^{m} (h_0(x^{(i)}) - y^{(i)}) = 0$

→ $\frac{1}{m} \sum_{i=i}^{2} (\Theta_0 x_0^{(i)} + \Theta_1 x_1^{(i)} + \cdots \cdots \Theta_1 x_1^{(i)} - y^{(i)}) = 0$

Taking $\Sigma$ inside

$$\rightarrow \frac{1}{m}\left[\sum_{i=1}^{m} \theta_0 x_0^{(i)} + \sum_{i=1}^{m} \theta_1 x_1^{(i)} + \dots + \sum_{i=1}^{m} \theta_n x_n^{(i)} - \sum_{i=1}^{m} y^{(i)}\right] = 0$$

$$\rightarrow \frac{1}{m}\left[\theta_0 x_0' + \theta_1 x_1' + \theta_2 x_2' + \dots \theta_n x_n' - y'\right] = 0$$

$$\left[\text{where } x_n' = \text{sum of all } x_n \atop y_x' = \text{sum of all } y\right]$$

$$\rightarrow \frac{\theta_0 x_0'}{m} + \frac{\theta_1 x_1'}{m} + \frac{\theta_2 x_2'}{m} + \dots \frac{\theta_n x_n'}{m} - \frac{y'}{m} = 0$$

$$\rightarrow \theta_0 \bar{x}_0 + \theta_1 \bar{x}_1 + \theta_2 \bar{x}_2 + \dots \theta_n \bar{x}_n - \bar{y} = 0$$

$$\left[\text{Let know mean} = \frac{\text{Sum of all } x}{\text{Total No. of } x} \Rightarrow \bar{x}\right]$$

Here proved linear regression, the least square fit line always passes through the point $\bar{x}, \bar{y}$

                    $\underset{\text{mean}}{\vee}$

Q C) b) $\underline{A.T.Q}$

Weak law of large number state that as exp data size is increase the expected value of error of prediction come close to the Actual true error

Let $\lambda_i$ be iid random variable with mean $= \mu$, $sd = \sigma$

The law state:

$$\lim_{n \to \infty} P \left\{ \right.$$

$$\underline{\text{In Mathematical terms law state that}}$$

$$\lim_{n \to \infty} P\left\{ |A_n - \mu| \geq \epsilon \right\} = 0$$

where $\left[ A_n \rightarrow \text{mean of first } n \text{ trials} \right]$

$$A_n = \frac{x_1 + x_2 \cdots x_n}{n}$$

$$E[A_n] = E\left[\frac{x_1 + x_2 \cdots + x_n}{n}\right]$$

$$E[A_n] = E\left[\frac{x_1 + x_2 \cdots x_n}{n}\right]$$

$$E[A_n] = \frac{n \times E[x_1]}{n} = \mu \quad \text{①} \quad \left[\begin{array}{l} \text{formula used :} \\ \text{let know that :} \\ E(x_1) = E(x_2) = \cdots E(x_n) \end{array}\right]$$

$$\text{Var}(A_n) = \text{Var}\left(\frac{x_1 + x_2 \cdots x_n}{n}\right) = \frac{1}{n^2}\left[\text{Var}(x_1) + \cdots \text{Var}(x_n)\right]$$

$$= \frac{1}{n^2} \cancel{\phantom{X}}$$

$$= \frac{1}{n^2} \times n \times \sigma^2 = \frac{\sigma^2}{n}$$

By from ~~Chaly the~~ Chebyshev's ~~eq~~ Inequality

$$\rightarrow \quad P\{|x-\mu| \geq \varepsilon\} \leq \frac{\sigma^2}{\varepsilon^2}$$

$$\rightarrow \quad P\{|A_n - \mu| \geq \varepsilon\} \leq \frac{\sigma^2}{n\varepsilon^2}$$

$$\rightarrow \quad \lim_{x \to \infty} P\{|A_n - \mu| \geq \varepsilon\} \leq \lim_{n \to \infty} \frac{\sigma^2}{n\varepsilon^2}$$

$$\rightarrow \quad n \longrightarrow \infty \quad \text{so value of } \frac{\sigma^2}{n\varepsilon^2} \text{ will tends to zero}$$

$$\rightarrow \quad \text{so we will get } \lim_{n \to \infty} P\{|A_n - \mu| \geq \varepsilon\} = 0$$

Hence proved

(C) Pseudo-Code

```
for(t=1; t≤100000000; t=tx10)
    sum=0

for(

for (j=1; j≤100000000; j=jx200) {
    sum=0;
    for(i=0; i<j; i++){
        int arr = random[]  → arr size j
        sum = sum.(arr);
    }
    mean = sum/ arr.size();
    print(mean)
}
```

Random fn() is giving Number between 1 to 200

as the j is increasing our mean will tends to move toward

100

(d) A - T. O

The Map algorithm maximizes the posterior probability

→ Bays law

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

→ Now for Map    A = weights, B = data

$$P(\omega|D) = \frac{1}{P(D)} \cdot P(D|\omega) \cdot P(\omega)$$

$$\left[ \text{where } \frac{1}{P(D)} \to \text{Normalization}, \quad P(D|\omega) \to \text{likelihood} \right.$$
$$\left. \text{and } P(\omega) \text{ is prior knowledge} \right]$$

→ we have a gaussian distribution of weights

$$P(\omega) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x-\mu}{2\sigma^2}}$$

Now $\omega \sim N(0, \sigma^2)$

→ $P(\omega) = \frac{1}{\sqrt{2\pi\sigma^2}} \times e^{-\left(\frac{\omega^T\omega}{2\sigma^2}\right)}$

→ Now for our $\omega \sim N(0,\sigma^2)$ let can say that
$$D|\omega \sim N(\upsilon^T\mu, \sigma^2)$$

Now, $P(D|\omega) = \prod_{k=1}^{n} U(y_k|\omega) = \prod_{k=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\left(\frac{1}{2\sigma}(\omega x - y_k)^2\right)}$

Taking log in

$$\log_e (P(\omega|D)) = \log_e (P(D|\omega)) + \log_e (P(\omega)) - \log_e (P(D))$$

As

According to MAP

Now let with maximise $P(w \mid D)$

~~if let most~~ maximising $\log_e (P(w \mid D))$

$w = \arg\max \log_e (P(w \mid D)) = \arg\max [\log_e (P(D \mid w)) + \log(P(w)) - \log_e(P(D))]$

$\rightarrow w = \arg\max_w [\log_e(P(D \mid w)) + \log_e (P(w))]$

$\rightarrow \log(P(D \mid w)) = \sum_{k=1}^{D} \log_e\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \frac{1}{2\sigma^2} \sum_{k=1}^{2} (w^T x - y_2)^2$

$\rightarrow \log(P \mid w)) = \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \left(\frac{w^T w}{2\sigma^2}\right)$

Now

$w = \arg\max_w \left[ A\log\left(\frac{1}{2\pi\sigma^2}\right) - \frac{1}{2\sigma^2} \sum_{k=1}^{n} (w^T x_k - y_k)^2 + \log\left(\frac{1}{2\pi\sigma^2}\right) - \frac{w^T w}{2\sigma^2} \right]$

$w = \arg\max_w \left[ \frac{-1}{2\sigma^2} \sum_{k=1}^{n} (w^T x_k - y_k)^2 - \frac{1}{2\sigma^2}(w^T w) \right]$

$w = \arg\min_w \left( \frac{1}{2\sigma^2} \sum_{k=1}^{n} (w^T x_k - y_k)^2 + \frac{w^T w}{2\sigma^2} \right)$

$$\left[ \sigma_o^2 = 1 \text{ \& divide equation by } x \right]$$

Hence $J(w) = \frac{1}{x} \arg\min_w \left( \sum_{k=1}^{n} (w^T x_k - y_k)^2 + \frac{w^T w}{2} \right)$

Later

**Answer B)**

Lets take A = x , B = Y where x and y have range of real numbers

Z = A*B   here we can see that Z have a strong correlation with A and B, as any one of them changes value of Z will change.

If A changes value of Z will change and similarly if B changes value of Z will change

But change in value of A doesn't change value of B so A and B have no correlation with each other same.

# Section- B

2. (15 points) **Section B (Scratch Implementation)**

Linear Regression

Implement Linear Regression on the given Dataset. You need to implement gradient descent from scratch i.e. you cannot use any libraries for training the model (You may use numpy, but libraries like sklearn are not allowed).

Dataset: **Housing Price Prediction Dataset**

(a) (6 marks) You will need to perform K-Fold cross-validation (K=2-5) in this exercise (implement from scratch). What is the optimal value of K? Justify it in your report along with the table for the mean RMSE of K-values and K-value.

(b) (3 marks) Plot the RMSE V/s iteration graph for all models trained with optimal value of K for K-Fold cross-validation. RMSE should be reported on the train and val set.

(c) (4 marks) Modify your Regression implementation by including L1 (LASSO) and L2 (Ridge Regression) regularization. Implement both regularization functions from scratch and train the model again. Try different values of the regularization parameter and report the best one. Plot similar RMSE V/s iteration graph as before (train and val loss).

(d) (2 marks) Implement the normal equation (closed form) for linear regression and get the optimal parameters directly for each fold (optimal K). Report the RMSE on respective validation sets.

## A)

### DATA pre-processing:

In the given dataset we check for if there is any null value and then we scaled the entire data.

Scaling formula used here : $z = \dfrac{x_i - \mu}{\sigma}$

```
1  data1 = pd.read_csv("Real estate\Real estate.csv")
2
3  # normalize_data = data1
   0.4s

1  data1=(data1-data1.mean())/data1.std()
   0.3s
```

We have also dropped the column 'No' from the data.

```
1  data1 = data1.drop(["No"], axis = 1)
   0.4s
```

### Implement Gradient Descent from scratch:

All formulas:

Cost function: $J(\theta)=$ $\dfrac{1}{2m}\displaystyle\sum_{1}^{m}(h(x^{(i)})-y^{(i)})^2$

$\theta_0,\theta_1,\ldots\ldots,\theta_n$ - weights

m : Total number of data samples

$x^{(i)}$: Input features for the $i^{th}$ data sample.

$h(x^{(i)})$: hypothesis

$y^{(i)}$: Actual output for $i^{th}$ data sample

Algorithm:



Code:

Side functions:

```
1  def plot_cost(costlist):
2      plt.title('Cost Function ')
3      plt.xlabel('No. of iterations')
4      plt.ylabel('Cost')
5      plt.plot(costlist)
6      plt.show()
```

```
1  def Y_actual(train_data):
2      Ytr_a = train_data.iloc[:,-1]
3      Ytr_a  = Ytr_a.values
4
5      return Ytr_a
6
7  def X_actual(train_data):
8      Xtr = train_data.iloc[:,:-1]
9      Xtr = Xtr.values
0
1      return Xtr
2
```

```python
# inplementing gradient descent
def helper_error(y_prediction , Ytr_a):  #giving us error= y_prediction - Y_actual
        return y_prediction - Ytr_a
```

Main function:

```python
def gradient_descent(x, y, m, alpha,epoch):
    cost_list = []   #to record all cost values to this list
    theta_list = []  #to record all theta_0 and theta_1 values to this list
    Ydash = []  # y predictions to store for testing later
    #  making theta
    np.random.seed(10)
    theta = np.random.rand(7)

    #adding 1 column in x
    x =np.insert(x,0,1,axis = 1)

    cost_list.append(20000000000000000)
    #declaring temp value later we will dlt this just to make it double from int

    for i in range(epoch) :
        y_prediction = np.dot(x, theta)
        Ydash.append(y_prediction)
        temp_error =helper_error(y_prediction , y)
        cost = 1/(2*m) * np.dot(temp_error.T, temp_error)   #cost = (1/2m)*sum[(y_prediction - Y_actual)^2
        cost_list.append(cost)
        theta = theta - (alpha * (1/m) * np.dot(x.T, temp_error))   # updating theta:  alpha * (1/m) * sum[error*x]
        # print(theta)
        theta_list.append(theta) #saving thetaa


    cost_list.pop(0)
    return Ydash, cost_list, theta_list
```

*Theta : weight of model*

```
1  def testing(X, theta,Yt_actual):
2      m = Yt_actual.size
3      X = np.insert(X,0,1,axis = 1)
4      y_prediction = np.dot(X, theta)
5      temp_error =helper_error(y_prediction , Yt_actual)
6      cost = 1/(2*m) * np.dot(temp_error.T, temp_error) # cost = (1/2m)*sum[(y_prediction - Y_actual)^2
7
8      return cost
```

We have initialized all the theta for the time randomly, then updating theta in every epoch

- *We have also Calculated the testing and training error for our model*

```
1 Y_preductionlist , costlist, thetas = gradient_descent(X_actual(train_data),Y_actual(train_data),Y_actual(train_data).size,0.0035,5500)
✓ 0.9s
```

Alpha = 0.0035

## Testing our model using test set

```
1  Yte_a = testin_data.iloc[:,-1]
2  Xte = testin_data.iloc[:,:-1]
3  Yte_a  = Yte_a.values
4  Xte = Xte.values
5  tr_theta = thetas[-1]
6  #  Mse error from testing set
7  train_error = testing(Xte,tr_theta,Yte_a)
8  print("Training error = " ,train_error)
9  # 0.24417512226623297 0.19545963519882276
```
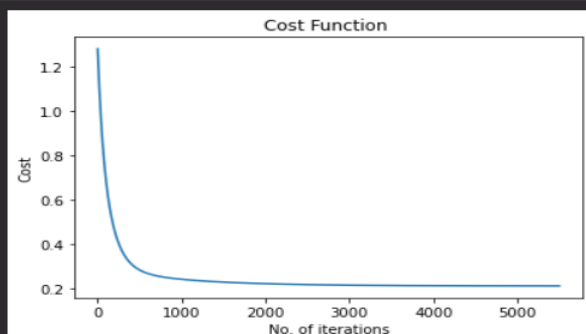[13]    ✓  0.5s

··· Training error =  0.19468835688259967

```
1  plot_cost(costlist)
✓ 0.2s                                                                   Pytho
```



Cost Function

## Part A)

First we have randomised the data sample

```
1  # randomize the data
2  kfoldData = data1
3  kfoldData = kfoldData.reindex(np.random.permutation(kfoldData.index))
4

0.6s                                                                    Python
```

```
1  # fixing index so that we can divide them into K parts
2  kfoldData = kfoldData.reset_index(drop = True)
```

Then we have done K-fold cross validation according to given question:

First we divided the data into K parts (K =5,4,3,2) and for every K we have made Train($n$) set and Test($m$) set {where n = k-1 and m =K-n}

Then we have Trained our model with the help of gradient descent on the Train set for every n and tested the model on Test set. We have then calculated the Root mean square error(rmse) with help of train model weights.

Rmse for all K-values

*Table: K_value vs Mean Rmse*

| K-value | Mean Rmse | |
|---------|-----------|---|
| 5 | Fold 1 | 0.598533 |
| | Fold 2 | 0.543886 |
| | Fold 3 | 0.567571 |
| | Fold 4 | 0.535383 |
| | Fold 5 | 0.937109 |
| | Mean:0.636797 | |
| 4 | Fold 1 | 0.628975 |
| | Fold 2 | 0.511595 |
| | Fold 3 | 0.535202 |
| | Fold 4 | 0.870449 |
| | Mean:0.636555 | |
| 3 | Fold 1 | 0.597461 |
| | Fold 2 | 0.549705 |
| | Fold 3 | 0.788908 |
| | Mean:0.6453584 | |
| 2 | Fold 1 | 0.5915280 |
| | Fold 2 | 0.7233616 |
| | 0.651257 | |

The optimal K-value is **K =4** that is **0.63655**


**Code:**

```
1  def k_fold_split(K) :
2
3
4      fold_data = []
5      i = 0
6      current = 0
7      add =int(len(kfoldData)/K)
8      count = add
9
10     for i in range(K-1):
11
12         fold_data.append(kfoldData.iloc[current:count])
13         current = count
14         count += add
15
16     fold_data.append(kfoldData.iloc[current:len(kfoldData)])
17
18     return fold_data
19
```

*Showing output code for Best Kmean:*

```
-------------> K =4

 1  fold_data = k_fold_split(4)
 2
 3  # for K  = 4
 4
 5  train1 = pd.concat([fold_data[0] , fold_data[1], fold_data[2]])
 6  test1 = fold_data[3]
 7  K1 = 0
 8
 9  train2 = pd.concat([fold_data[0] , fold_data[1], fold_data[3]])
10  test2 = fold_data[2]
11  K2 = 0
12
13  train3 = pd.concat([ fold_data[2], fold_data[0], fold_data[3]])
14  test3 = fold_data[1]
15  K3 = 0
16
17  train4 = pd.concat([ fold_data[2], fold_data[3], fold_data[1]])
18  test4 = fold_data[0]
19  K4 = 0
20
21
22  Y_preductionlist , costlist, thetas = gradient_descent(X_actual(train1),Y_actual(train1),Y_actual(train1).size,0.0035,5500)
23
24  K1 = testing(X_actual(test1),thetas[-1],Y_actual(test1))
25  K1 = (2*K1)**(0.5)    #Rmse = (2 cost(fn))^1/2
26  print("K1:",K1)
27
28  Y_preductionlist , costlist, thetas = gradient_descent(X_actual(train2),Y_actual(train2),Y_actual(train2).size,0.0035,5500)
29  K2 = testing(X_actual(test2),thetas[-1],Y_actual(test2))
30  K2 = (2*K2)**(0.5) #Rmse = (2 cost(fn))^1/2
31  print("K2:",K2)
32
33  Y_preductionlist , costlist, thetas = gradient_descent(X_actual(train3),Y_actual(train3),Y_actual(train3).size,0.0035,5500)
34  K3 = testing(X_actual(test3),thetas[-1],Y_actual(test3))
35  K3 = (2*K3)**(0.5)    #Rmse = (2 cost(fn))^1/2
36  print("K3:",K3)
37
38  Y_preductionlist , costlist, thetas = gradient_descent(X_actual(train4),Y_actual(train4),Y_actual(train4).size,0.0035,5500)
39  K4 = testing(X_actual(test4),thetas[-1],Y_actual(test4))
40  K4 = (2*K4)**(0.5)    #Rmse = (2 cost(fn))^1/2
41  print("K4:",K4)
42
43
44  # K value
45  print("Mean K value:",(K1+K2+K3+K4)/4)
46  # K1: 0.6289757784431346
47  # K2: 0.5115958606944053
48  # K3: 0.5352026365186768
49  # K4: 0.8704491695393924
50  # Mean K value: 0.6365558612989023

✓ 0.8s

K1: 0.6289757784431453
K2: 0.5115958606944057
K3: 0.5352026365186744
K4: 0.8704491695393949
Mean K value: 0.6365558612989051
```
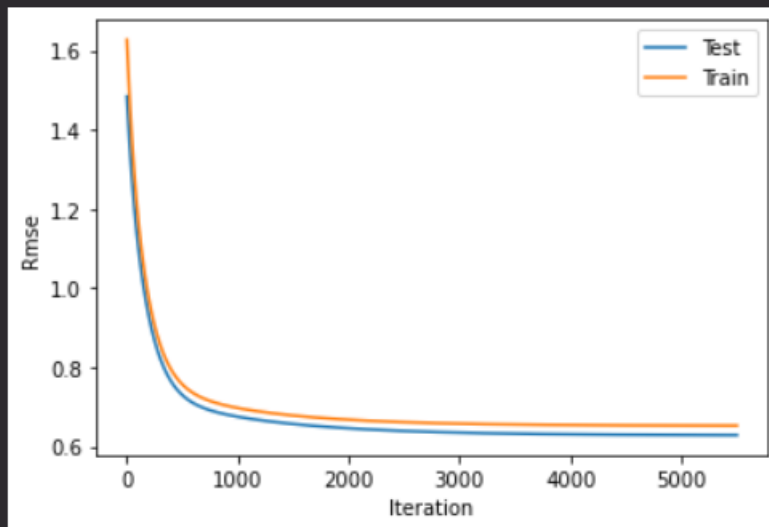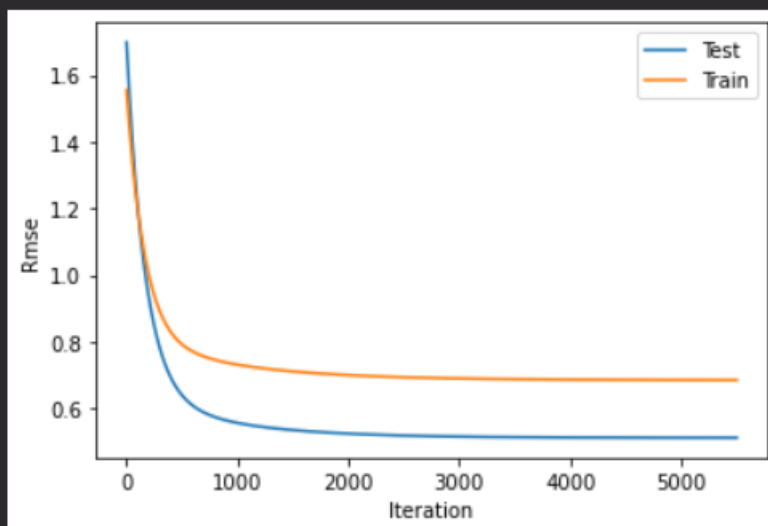
PART B

*RMSE vs Iteration graph for all models trained with optimal value of K=4 for K-fold*
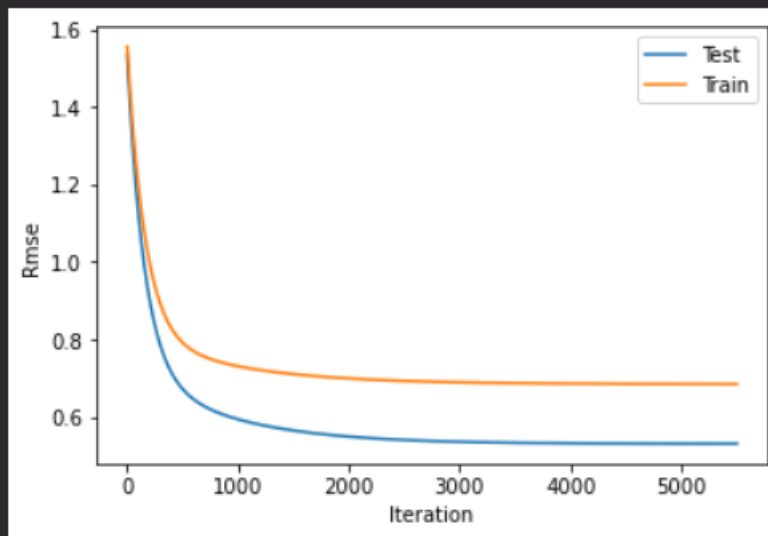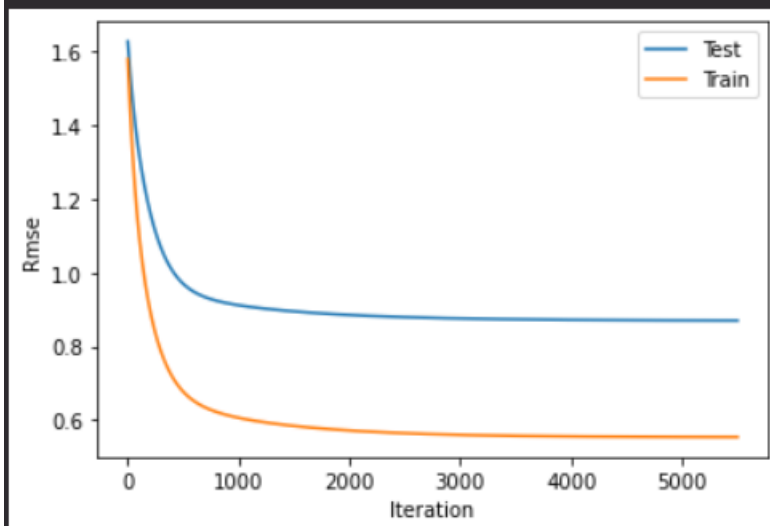
**RMSE vs Iteration Fold:1**



**RMSE vs Iteration Fold:2**

RMSE vs Iteration Fold:3



RMSE vs Iteration Fold:4

PART :C

L1 regularization:

*Formula used:*

Cost function

$$\mathcal{J}(\theta) = \frac{1}{2n} \sum_{i=1}^{n} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{i=0}^{n} |\theta_i|$$

$$\begin{cases} \frac{d}{d\theta}(\mathcal{J}) = \frac{1}{m} \sum_{i=1}^{n} (h_\theta(x^{(i)}) - y^{(i)}) + \lambda & \text{when } \theta_i \\ \frac{d}{d\theta}(\mathcal{J}) = \frac{1}{m} \sum_{i=1}^{n} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \lambda \end{cases}$$

Updating Θⱼ (way to update the theta inside code for L1)



$$\begin{cases} \theta_{j} = \theta_j - \text{cost function} - \lambda & \text{when } \theta_j \geq 0 \\ \theta_j = \theta_j - \text{cost function} + \lambda & \text{when } \theta_j \leq 0 \end{cases}$$

## Code

```
1   # implementing l1
2   def L1_helper_error(y_prediction , Ytr_a):
3           return y_prediction - Ytr_a
4
```

```
def L1_testing(X, theta,Yt_actual,lemda):
    m = Yt_actual.size
    X = np.insert(X,0,1,axis = 1)
    y_prediction = np.dot(X, theta)
    temp_error =L1_helper_error(y_prediction , Yt_actual)
    cost = 1/(2*m)*np.dot(temp_error.T, temp_error) + abs(lemda)*sum(theta)    # cost = (1/2m)*sum[(y_prediction - Y_actual)^2 + sum(theta)
    rmse = (2*cost)**(0.5)
    return rmse
```

```
def L1_gradient_descent(x, y, x_test,y_test,m, alpha,epoch,lemda,rmsevsval):
    #rmsevsval when =1 then act as function to caluclate rmse vs val graph
    rmse_list_Train = []    #to record all cost values to this list
    theta_list = []  #to record all theta_0 and theta_1 values to this list
    # Ydash = []  # y predictions to store
    #   making theta
    rmse_list_Test = []
    np.random.seed(10)
    theta = np.random.rand(7)

    #adding 1 column in x
    x =np.insert(x,0,1,axis = 1)

    rmse_list_Train.append(20000000000000000)    #declaring temp value later we will dlt this just to make it double from int

    for i in range(epoch) :
        y_prediction = np.dot(x, theta)
        # Ydash.append(y_prediction)
        temp_error =L1_helper_error(y_prediction , y)
        cost = 1/(2*m) * np.dot(temp_error.T, temp_error) + abs(lemda)*sum(theta)    #  cost = (1/2m)*sum[(y_prediction - Y_actual)^2 + lemda*sum(theta)
        # print(cost)
        rmse = (2*cost)**(0.5)

        rmse_list_Train.append(rmse)
        # print(theta)
```

```
    for i in range(len(theta)):# In maths terms we can say that we
        if theta[i]>0:          # are opening mod that is :|lemda|
            theta[i]-=lemda
        else:
            theta[i] +=lemda
    theta = theta - (alpha * (1/m) * np.dot(x.T, temp_error))   # updating theta:  alpha * (1/m) * sum[error*x] - |lemda|
    # print((alpha * (1/m) * np.dot(x.T, temp_error)))
    theta_list.append(theta) #saving thetaa


    if(rmsevsval == 1):
        rmse_list_Test.append(Re_testing(x_test,theta,y_test))



rmse_list_Train.pop(0)
if(rmsevsval== 1):
    return  rmse_list_Train , rmse_list_Test
else :
    return rmse_list_Train, theta_list
```
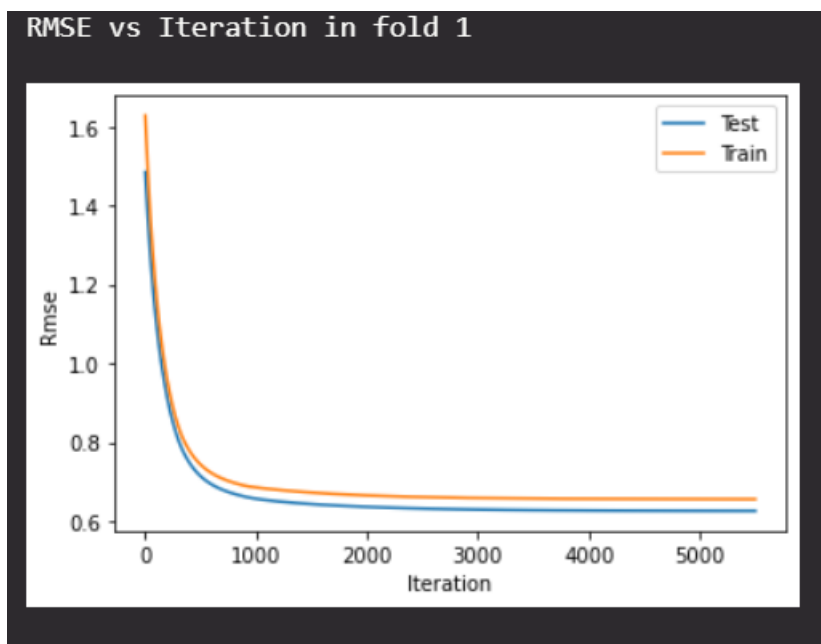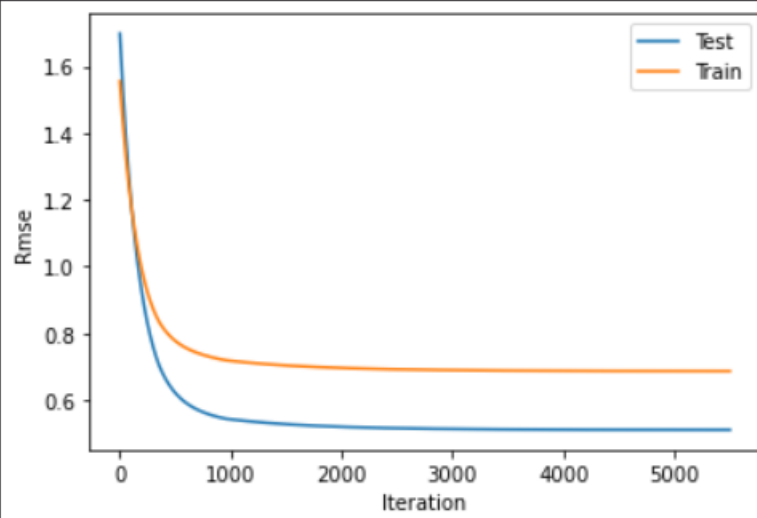
# Different values of lambda vs Mean RMSE:

**Table:**

| Lambda | Mean Rmse |
|---|---|
| 0.1 | 0.9516836 |
| 0.001 | 0.7142055 |
| 0.0001 | 0.625595 |
| 0.000051 | 0.6261506 |
| 0.0000059 | 0.6285677 |
| 0.00000069 | 0.6289269 |

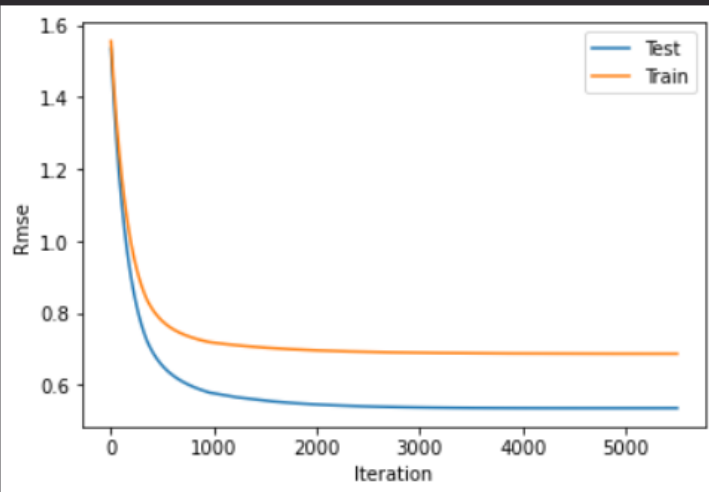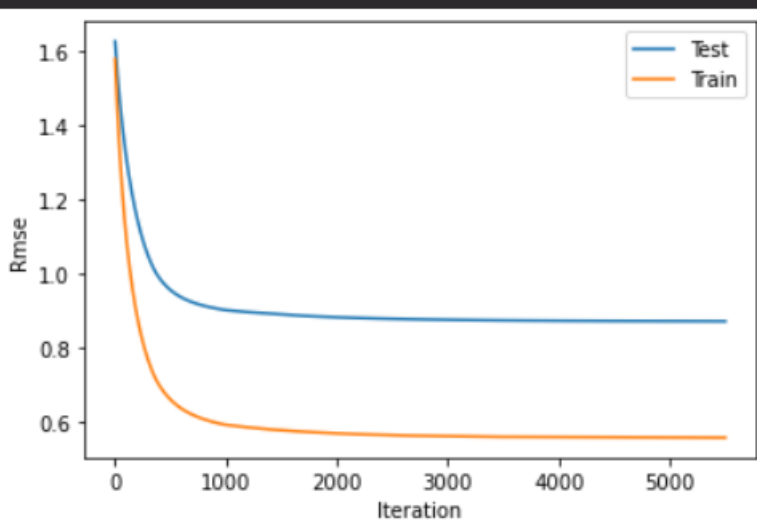*K-Fold Cross Validation Method with K = 4 and Lambda = 0.0001*

Graph :

RMSE vs Iteration in fold 2



RMSE vs Iteration in fold 3



RMSE vs Iteration in fold 4

Code for K = 4 and lambda = 0.0001

```
1  fold_data = k_fold_split(4)
2
3  # for K  = 4
4
5  train1 = pd.concat([fold_data[0] , fold_data[1], fold_data[2]])
6  test1 = fold_data[3]
7
8  train2 = pd.concat([fold_data[0] , fold_data[1], fold_data[3]])
9  test2 = fold_data[2]
10
11 train3 = pd.concat([fold_data[1], fold_data[0], fold_data[3]])
12 test3 = fold_data[1]
13
14
15 train4 = pd.concat([fold_data[1] , fold_data[2], fold_data[3]])
16 test4 = fold_data[0]
17
18
19 rmsetrain1, rmsetest1  = L1_gradient_descent(X_actual(train1),Y_actual(train1),X_actual(test1),Y_actual(test1),Y_actual(train1).size,0.0035,5500,0.0001,1)
20
21 rmsetrain2, rmsetest2  = L1_gradient_descent(X_actual(train2),Y_actual(train2),X_actual(test2),Y_actual(test2),Y_actual(train2).size,0.0035,5500,0.0001,1)
22
23 rmsetrain3, rmsetest3  = L1_gradient_descent(X_actual(train3),Y_actual(train3),X_actual(test3),Y_actual(test3),Y_actual(train3).size,0.0035,5500,0.0001,1)
24
25 rmsetrain4, rmsetest4  = L1_gradient_descent(X_actual(train4),Y_actual(train4),X_actual(test4),Y_actual(test4),Y_actual(train4).size,0.0035,5500,0.0001,1)
26
✓ 0.8s
```

```
1  print("RMSE vs Iteration in fold 1")
2  re_plot(rmsetest1,rmsetrain1)
3  print("RMSE vs Iteration in fold 2")
4  re_plot(rmsetest2,rmsetrain2)
5  print("RMSE vs Iteration in fold 3")
6  re_plot(rmsetest3,rmsetrain3)
7  print("RMSE vs Iteration in fold 4")
8  re_plot(rmsetest4,rmsetrain4)
9
```

# L2 regulaeization(Ridge regression)

Formulas used:



Cost function

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=0}^{n} \theta_{(j)}^2$$

$$\frac{dJ}{d\theta} = \frac{1}{M} \sum (h_\theta(x^{(i)}) - y_i^{(i)}) \, h(x)_j^{(i)} + 2\lambda \theta_j$$

Code:

```
def L2_gradient_descent(x, y, x_test,y_test,m, alpha,epoch,lemda,rmsevsval):
    rmse_list_Train = []   #to record all cost values to this list
    theta_list = []  #to record all theta_0 and theta_1 values to this list
    # Ydash = [] # y predictions to store
    #  making theta
    rmse_list_Test = []
    np.random.seed(10)
    theta = np.random.rand(7)

    #adding 1 column in x
    x =np.insert(x,0,1,axis = 1)

    rmse_list_Train.append(20000000000000000)    #declaring temp value later we will dlt this just to make it double from int

    for i in range(epoch) :
        y_prediction = np.dot(x, theta)
        # Ydash.append(y_prediction)
        temp_error =L2_helper_error(y_prediction , y)
        cost = 1/(2*m) * np.dot(temp_error.T, temp_error) + lemda*np.dot(theta.T,theta)    #  cost = (1/2m)*sum[(y_prediction - Y_actual)^2 + (lambda)*sum(theta)^2
        rmse = (2*cost)**(0.5)
        rmse_list_Train.append(rmse)
        # print(theta)
        theta = theta - (alpha * (1/m) * np.dot(x.T, temp_error)) - 2*lemda*sum(theta)  # updating theta:  alpha * (1/m) * 2(lambda)sum[error*x]

        theta_list.append(theta) #saving thetaa

        if(rmsevsval == 1):
            rmse_list_Test.append(Re_testing(x_test,theta,y_test))


    rmse_list_Train.pop(0)
    if(rmsevsval== 1):
        return  rmse_list_Train , rmse_list_Test
    else :
        return rmse_list_Train, theta_list
```

# Different values of lambda vs Mean RMSE:

| Lambda | Mean Rmse |
|---|---|
| 0.1 | 0.8888019 |
| 0.001 | 0.7704632 |
| 0.0001 | 0.6275338 |
| 0.000051 | 0.62808050 |
| 0.0000059 | 0.62885283 |
| 0.00000069 | 0.62885283 |

Best value of lamda for L2(**Ridge regression**) is 0.0001

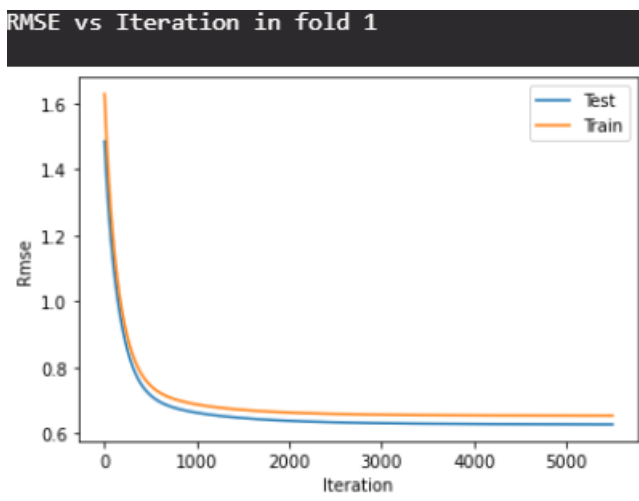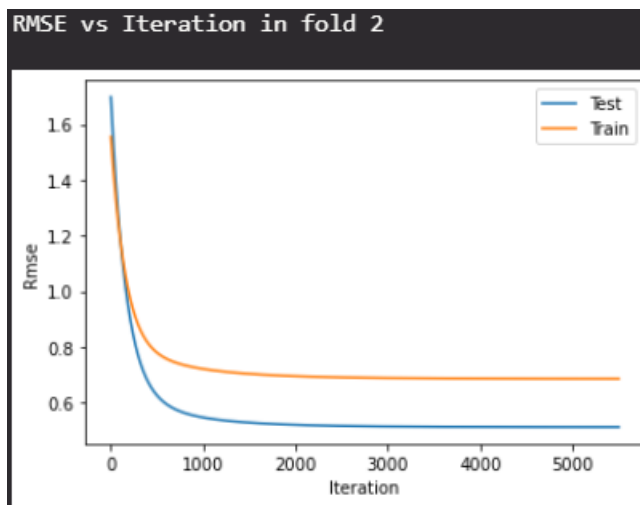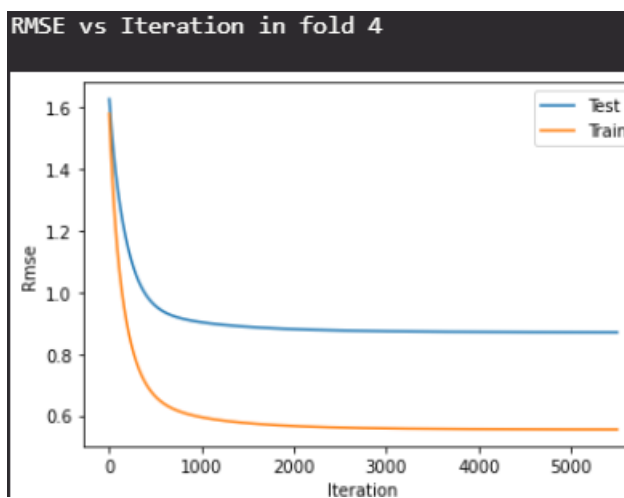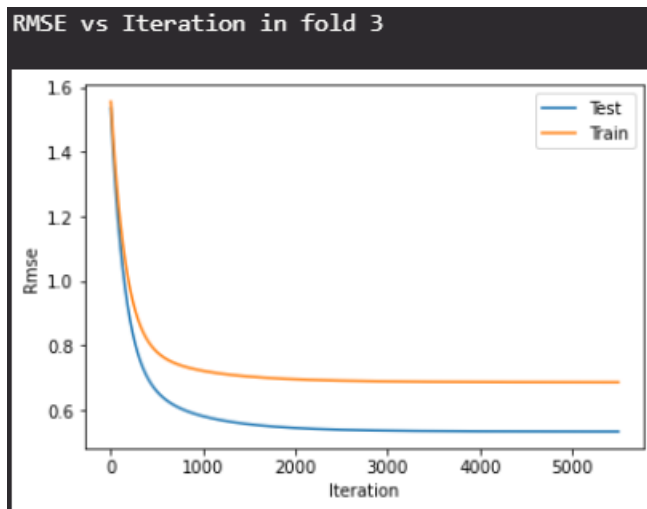L2 plot Rmse V/s iteration for K =4

 Code:

```
1  fold_data = k_fold_split(4)
2
3  # for K = 4
4
5  train1 = pd.concat([fold_data[0] , fold_data[1], fold_data[2]])
6  test1 = fold_data[3]
7
8  train2 = pd.concat([fold_data[0] , fold_data[1], fold_data[3]])
9  test2 = fold_data[2]
10
11 train3 = pd.concat([fold_data[1], fold_data[0], fold_data[3]])
12 test3 = fold_data[1]
13
14
15 train4 = pd.concat([fold_data[1] , fold_data[2], fold_data[3]])
16 test4 = fold_data[0]
17
18
19 rmsetrain1, rmsetest1  = L2_gradient_descent(X_actual(train1),Y_actual(train1),X_actual(test1),Y_actual(test1),Y_actual(train1).size,0.0035,5500,0.0001,1)
20
21 rmsetrain2, rmsetest2  = L2_gradient_descent(X_actual(train2),Y_actual(train2),X_actual(test2),Y_actual(test2),Y_actual(train2).size,0.0035,5500,0.0001,1)
22
23 rmsetrain3, rmsetest3  = L2_gradient_descent(X_actual(train3),Y_actual(train3),X_actual(test3),Y_actual(test3),Y_actual(train3).size,0.0035,5500,0.0001,1)
24
25 rmsetrain4, rmsetest4  = L2_gradient_descent(X_actual(train4),Y_actual(train4),X_actual(test4),Y_actual(test4),Y_actual(train4).size,0.0035,5500,0.0001,1)
26
```

Graph:

**RMSE vs Iteration in fold 3**



**RMSE vs Iteration in fold 4**



PART :**D**

*Normal equation :*

$$\theta = (X^TX)^{-1}.(X^TY)$$

*Code:*

```
1  def nor_equation_train(X,Y):
2      m =len(X)
3      X =np.insert(X,0,1,axis = 1)
4      #θ =(XTX)-1.( XTY)
5      theta = np.dot(np.linalg.inv((np.dot(X.T,X))),np.dot(X.T,Y))
6      Y_eq = np.dot(X,theta)
7      temp_error = helper_error(Y_eq,Y)
8      cost = 1/(2*m) * np.dot(temp_error.T, temp_error)
9      rmse = (2*cost)**(0.5)
10     return theta ,rmse
```

==As we know optimal value of K=4==

Training our model for K=4

Table for K-value vs RMSE

| Fold no. | Optimal parameters | | RMSE |
|---|---|---|---|
| 1 | *Parameter* | *Value* | 0.627966 |
| | $\theta_0$ | 0.0207766 | |
| | $\theta_1$ | 0.0877178 | |
| | $\theta_2$ | -0.2477262 | |
| | $\theta_3$ | -.04113045 | |
| | $\theta_4$ | 0.23591649 | |
| | $\theta_5$ | 0.22846858 | |
| | $\theta_6$ | -0.0130103 | |
| 2 | *Parameter* | *Value* | 0.511711 |
| | $\theta_0$ | $-2.8157*e^{-4}$ | |
| | $\theta_1$ | 0.011904 | |
| | $\theta_2$ | -0.22159 | |
| | $\theta_3$ | -0.42743 | |
| | $\theta_4$ | 0.22803 | |
| | $\theta_5$ | 0.20452 | |
| | $\theta_6$ | 0.002188 | |
| 3 | *Parameter* | *Value* | 0.5320325 |
| | $\theta_0$ | $-2.8157e^{-4}$ | |
| | $\theta_1$ | 0.19404 | |
| | $\theta_2$ | 0.22159 | |
| | $\theta_3$ | -0.42743 | |
| | $\theta_4$ | 0.22802 | |
| | $\theta_5$ | 0.20452 | |
| | $\theta_6$ | $2.18820e^{-3}$ | |
| 4 | *Parameter* | *Value* | 0.8691125 |
| | $\theta_0$ | -0.01102 | |
| | $\theta_1$ | 0.114846 | |
| | $\theta_2$ | -0.234727 | |
| | $\theta_3$ | 0.3940832 | |
| | $\theta_4$ | 0.2766535 | |
| | $\theta_5$ | 0.1821304 | |
| | $\theta_6$ | -0.019418 | |

3. (15 points) **Section C (Algorithm implementation using packages)**

   In this question, you are expected to understand and run Naive Bayes Algorithm.
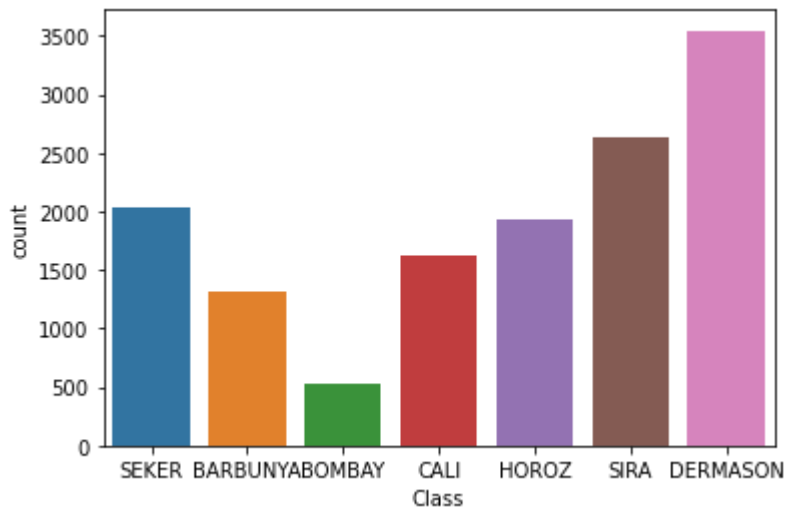
   Dataset: **Dry Bean Dataset**

   (a) (1 marks) For the given dataset, plot the class distribution and analyze.

   (b) (2 marks) Perform EDA (histograms, box plots, scatterplots, etc.) and give at least five insights on the data. Check the missing values in the dataset.

   (c) (3 marks) Use TSNE (t-distributed stochastic neighbor embedding) algorithm to reduce data dimensions to 2 and plot the resulting data as a scatter plot. Comment on the separability of the data.

   (d) (2 marks) Run the sklearn's implementation of Naive Bayes (Any 2 of your choice - refer **here**). Report Accuracy, Recall, and Precision. Comment on the results and their differences from the two implementations of Naive Bayes. (80:20 train test split)

   (e) (3 marks) Use Principal Component Analysis (PCA) to reduce the number of features and use the reduced data set for model training. Use values 4,6,8,10 and 12 for the number of components. Compare results (Accuracy, Precision, Recall, and F-1 score). (80:20 train test split)

   (f) (2 marks) Use Scikit-learn to plot the ROC-AUC curves and comment on the output.

   (g) (2 marks) Train your model using Sklearn's implementation of Logistic Regression, choose appropriate parameters, and comment on your choice. Compare the results with the ones obtained from Naive Bayes models. (80:20 train test split)

## Checking for any missing value In data or NA values

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 13611 entries, 0 to 13610
Data columns (total 17 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   Area             13611 non-null  int64
 1   Perimeter        13611 non-null  float64
 2   MajorAxisLength  13611 non-null  float64
 3   MinorAxisLength  13611 non-null  float64
 4   AspectRation     13611 non-null  float64
 5   Eccentricity     13611 non-null  float64
 6   ConvexArea       13611 non-null  int64
 7   EquivDiameter    13611 non-null  float64
 8   Extent           13611 non-null  float64
 9   Solidity         13611 non-null  float64
 10  roundness        13611 non-null  float64
 11  Compactness      13611 non-null  float64
 12  ShapeFactor1     13611 non-null  float64
 13  ShapeFactor2     13611 non-null  float64
 14  ShapeFactor3     13611 non-null  float64
 15  ShapeFactor4     13611 non-null  float64
 16  Class            13611 non-null  object
```

PART : **A**

Class distribution:
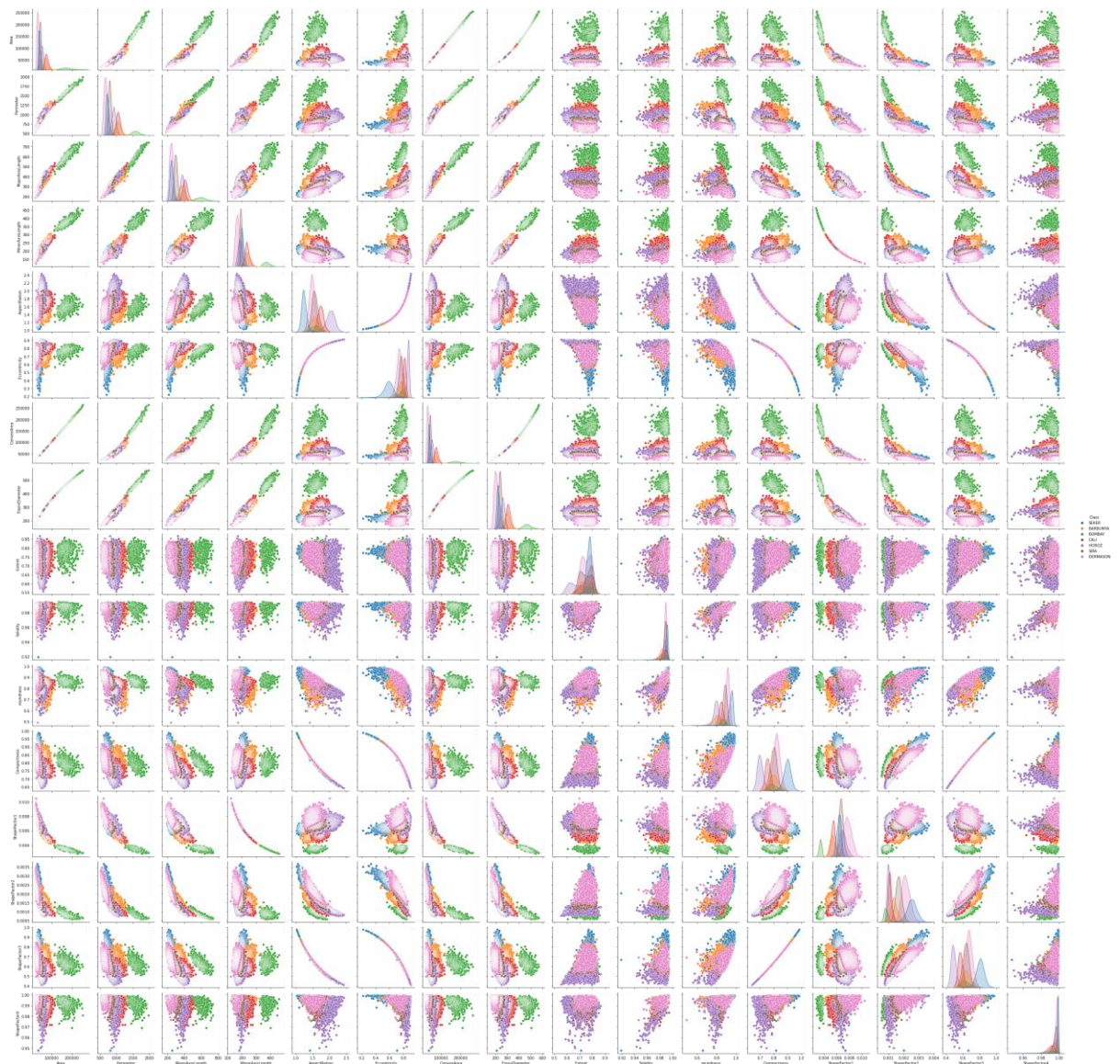


*Analysis*

From the above plot, we can say that :

1) Our data set classified dry beans into seven parts based on different features. Therefore it is multi-class classification
2) The Dermason dry bean has the highest count, meaning it has more sample data in the entire data set.
3) The Bombay dry bean has the lowest count, meaning it has least sample data in the entire data set

## PART : B

**1) Scatter plot:**
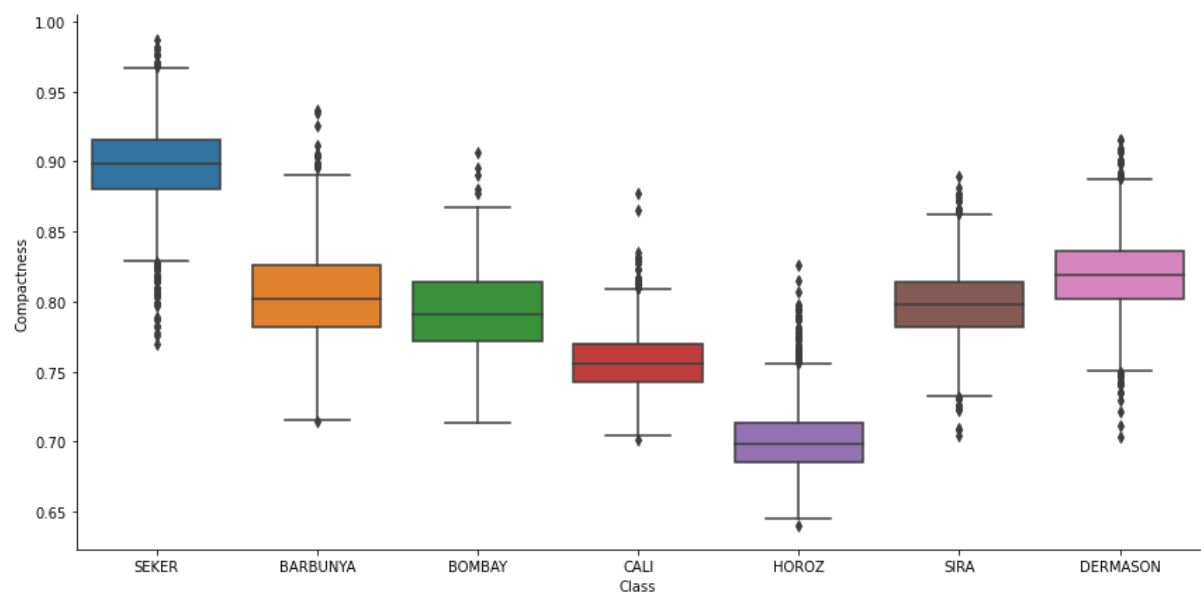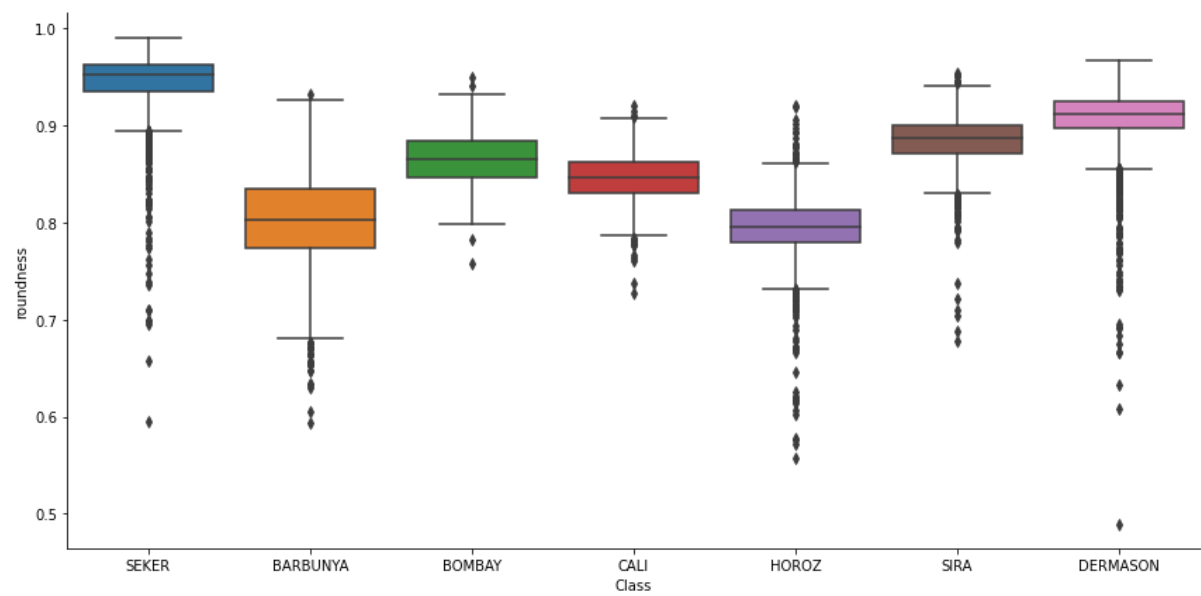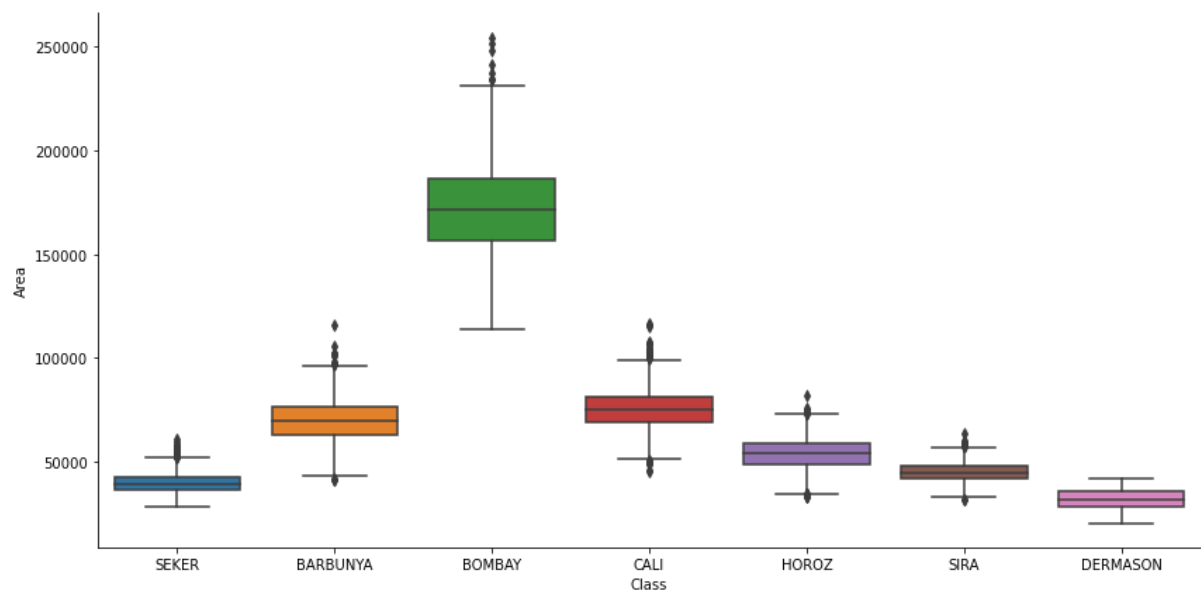
Code:
```
1  sea.pairplot(df,hue="Class")
```
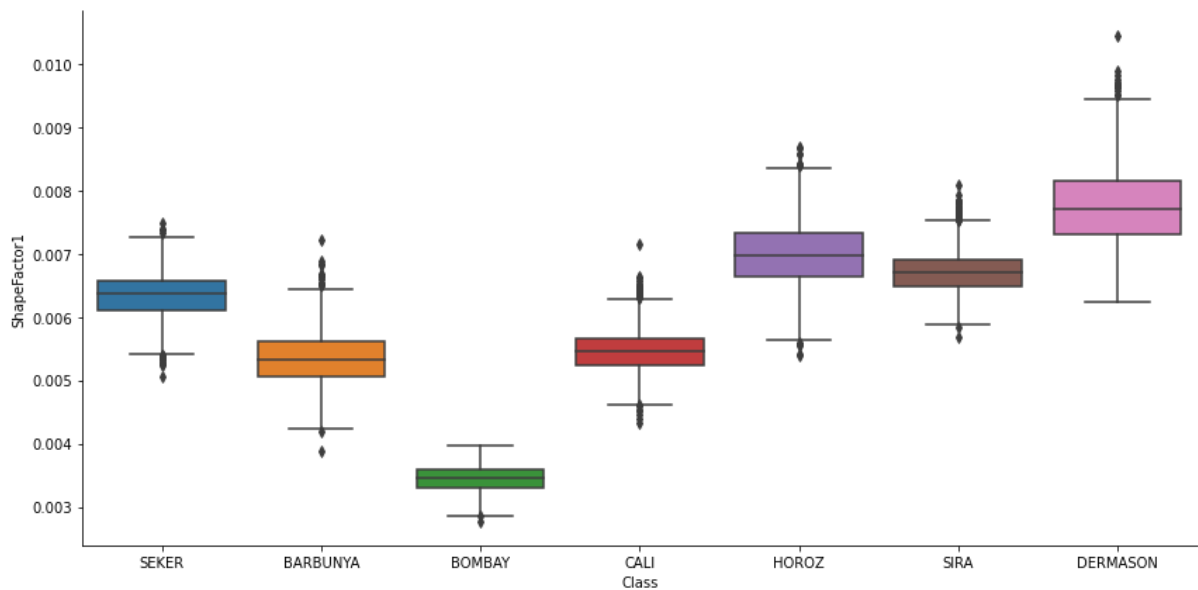
2) Box plots:

Code:

```
1  sea.catplot(x ="Class",y="Area",data=df,height=6,aspect=2,kind="box")
2  sea.catplot(x ="Class",y="roundness",data=df,height=6,aspect=2,kind="box")
3  sea.catplot(x ="Class",y="Compactness",data=df,height=6,aspect=2,kind="box")
4  sea.catplot(x ="Class",y="ShapeFactor1",data=df,height=6,aspect=2,kind="box")
3.3s
```
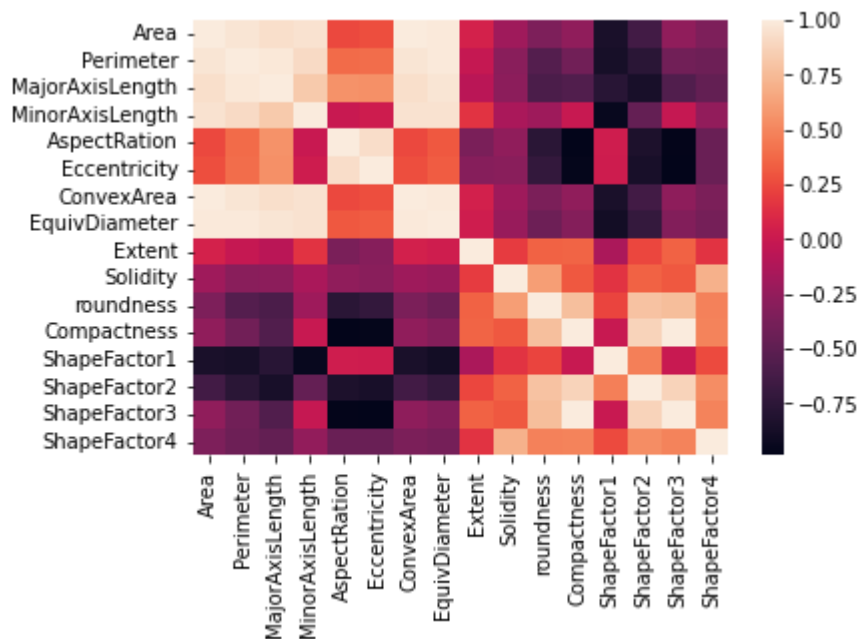
### 3) Corelation map

Code:
```
1  #corelation
2  corr = df.corr()
3  # temp=corr.style.background_gradient(cmap='co
4  sea.heatmap(corr,annot =False,)
5
```
✓ 0.8s



### Insights from the data:

- Using Boxplot, we tell about the distribution of data set, variance, and outliers of features in our data.
  1) Area , perimeter, MajorAxislength, MinorAxisLength has identical distribution
  2) Eccentricity,Solidity,Roundness have higher number of outliers

- Using scatter plot:
  1) Shape factor 2 and major axis length are hyperbolic functions in nature only in positive x and y axis.
  2) Perimeter and major axis length have linear relationship
  3) We also know the density distribution.


- Area, convex have very high positive correlation with each other.
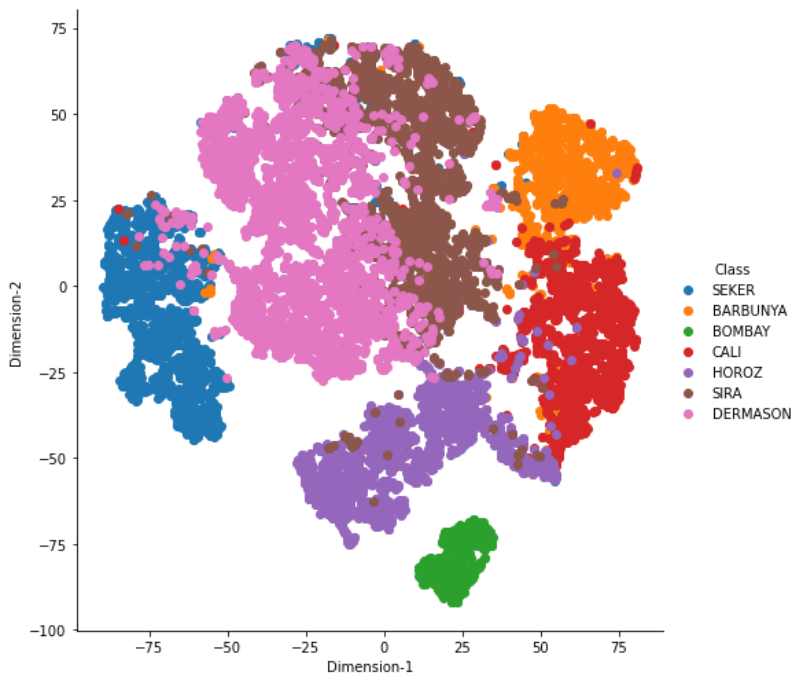- Minor axis length and shape factor1 have very high negative correlation

## PART :C

First we have done data processing

```
1  standardized_data = StandardScaler().fit_transform(df)
2  print(standardized_data.shape)
```

Code:

```
1  model = TSNE(2,random_state=0)
2  tsne_d = model.fit_transform(standardized_data)
3  tsne_d = np.vstack((tsne_d.T,Class)).T
4  tsne_df= pd.DataFrame(tsne_d,columns=("Dimension-1" , "Dimension-2" ,"Class"))
5
6  #plottting in 2d
7  sea.FacetGrid(data = tsne_df,hue="Class",height =7).map(plt.scatter,'Dimension-1','Dimension-2').add_legend()
58.8s
```



**Comment on the separability of the data.**

From the above plot, we can say:

- BOMBAY class have low correlation and no data is overlapping with other classes
- DERMANSON and Sira class have high correlation
- Cali class have many outliers as we can see it has many scatter dots far away from it

PART :**D**

*According to given question, I have implemented Bernoulli Naïve Bayes and Gaussian Naïve bayes*

➢ Split data into test and train set

```
X_train, X_test, y_train, y_test = train_test_split(tsne_df.iloc[:,:2], tsne_df.iloc[:,-1] ,test_size=0.2, random_state=0)
```

*Bernoulli Naïve Bayes*

*Code:*

```
1  bnlli = BernoulliNB()
2  model_b = bnlli.fit(X_train,y_train)
3  y_pred_train= bnlli.predict(X_train)
4  y_pred_test = bnlli.predict(X_test)
5
6
7  print("Accuracy on test set",accuracy_score(y_test,y_pred_test))
8  print("precision on test set",precision_score(y_test,y_pred_test,average="weighted"))
9  print("recall on test set",recall_score(y_test,y_pred_test,average="macro"))
```

Accuracy, precision ,recall are mention below

```
Accuracy on test set 0.4781491002570694
precision on test set 0.37312143502755596
recall on test set 0.38329053201827845
```

*Gaussian Naïve Bayes*

*Code:*

```
gauss = GaussianNB()
model_b = gauss.fit(X_train,y_train)
y_pred_train= gauss.predict(X_train)
y_pred_test = gauss.predict(X_test)

print("Accuracy on test set",accuracy_score(y_test,y_pred_test))
print("precision on test set",precision_score(y_test,y_pred_test,average="weighted"))
print("recall on test set",recall_score(y_test,y_pred_test,average="macro"))
```

Accuracy, precision ,recall are mention below

```
Accuracy on test set 0.9045170767535806
precision on test set 0.9045778178745367
recall on test set 0.9174189240812061
```

Result: The accuracy, precision and recall of *Gaussian Naïve Bayes*

Far more than *Bernoulli Naïve Bayes*

PART :**E**

Use Principal Component Analysis (PCA) to reduce the number of features and use the reduced data set for model training.

Code:

```python
1  from sklearn.decomposition import PCA
2  #already standardized_data-
3
4  def principal_c_a(n_comp , data):
5      pca = PCA(n_components = n_comp)
6      pca.fit(data)
7      data_pca = pca.fit_transform(data)
8
9      return data_pca
```
Assistant

```python
1  # A.t.o applying PCA
2  given_time = [4,6,8,10,12]
3  for i in given_time:
4
5      pca = principal_c_a(i,standardized_data)
6      X_train, X_test, y_train, y_test = train_test_split(pca, Class ,test_size=0.2, random_state=0)
7
8      model_b = gauss.fit(X_train,y_train)
9      y_pred_train= gauss.predict(X_train)
10     y_pred_test = gauss.predict(X_test)
11     print("for the numbers of components:", i)
12     print("Accuracy on test set",accuracy_score(y_test,y_pred_test))
13     print("precision on test set",precision_score(y_test,y_pred_test,average="weighted"))
14     print("recall on test set",recall_score(y_test,y_pred_test,average="macro"))
15     print("F1  on test set",f1_score(y_test,y_pred_test,average="macro"))
16     print()
17
```

*for the numbers of components outputs:*

```
for the numbers of components: 4
Accuracy on test set 0.8670583914799853
precision on test set 0.8654861093320928
recall on test set 0.8694964513745642
F1   on test set 0.8704054725617627

for the numbers of components: 6
Accuracy on test set 0.9015791406536908
precision on test set 0.9040559819325191
recall on test set 0.9150020522588272
F1   on test set 0.9152300314347179

for the numbers of components: 8
Accuracy on test set 0.9026808666911494
precision on test set 0.90828111840676
recall on test set 0.9186659952734894
F1   on test set 0.9189686359561994

for the numbers of components: 10
Accuracy on test set 0.8872567021667279
precision on test set 0.8990897236147188
recall on test set 0.9064773166690104
F1   on test set 0.9051115676407203

for the numbers of components: 12
Accuracy on test set 0.8868894601542416
precision on test set 0.8978695844729355
recall on test set 0.9065575456275566
F1   on test set 0.9059484232352736
```

## PART :G

using Sklearn's implementation of Logistic Regression,choose appropriate parameters

Code:

```python
from sklearn.linear_model import LogisticRegression
X = standardized_data
Y = Class
X_train, X_test, y_train, y_test = train_test_split(X, Y ,test_size=0.2, random_state=0)
lg = LogisticRegression(random_state=0,penalty='l2',max_iter=15000,solver='sag').fit(X_train

y_pred_train= lg.predict(X_train)
y_pred_test = lg.predict(X_test)

print("Accuracy on test set",accuracy_score(y_test,y_pred_test))
print("precision on test set",precision_score(y_test,y_pred_test,average="weighted"))
print("recall on test set",recall_score(y_test,y_pred_test,average="macro"))
```

Output:

```
Accuracy on test set 0.9280205655526992
precision on test set 0.9286837346038485
recall on test set 0.9360002219265073
```

1) I have  taken L2 penality so that model don't overfit
2) I have taken solver – "sag" because given data is multiclass classification. In multiclass classification only 4 solver available. Both "sag" and "saga" are best for normalized data and our data used here are normalized one only
3) Max_iter = 15000 because after numbers of trials we get best value for this


Result : The result of  Logistic Regression far more better than Naive Bayes models.