

**MACHINE LEARNING (CSE343/ECE343)**  
**ASSIGNMENT-1**  
**Sourav Goyal 2020341**  
**SECTION-A**

**ANS - (a)**

Assignment-1 (ML)

Section-A

A-① Given :- Model we have is linear regression,

To prove :- least square fit line always passes through the point  $(\bar{x}, \bar{y})$ , where  $\bar{x}$  and  $\bar{y}$  represent the arithmetic mean of the independent variables and dependent variables respectively.

∴ Cost function for simple linear regression is :-

$$J(\omega) = \underset{n}{\operatorname{argmin}} \sum_{i=1}^n (\omega_1 x_i - y_i + \omega_0)^2, \text{ where } \omega_0 \text{ is bias and } \omega_1 \text{ is feature weight}$$

→ Now in least square fit our cost function will be minimum so, we will derive our  $J(\omega)$  and put 0 because we have to find minima, i.e.,

$$\frac{dJ(\omega)}{d\omega} = \frac{dJ(\omega)}{d\omega_1} + \frac{dJ(\omega)}{d\omega_0}$$

$$\frac{dJ(\omega)}{d\omega_1} = \frac{1}{n} \sum_{i=1}^n (\omega_1 x_i - y_i + \omega_0) x_i \times 2 = 0 \quad \text{--- (1)}$$

$$\frac{dJ(\omega)}{d\omega_0} = \frac{1}{n} \times 2 \sum_{i=1}^n (\omega_1 x_i - y_i + \omega_0) = 0$$

$$\frac{dJ(\omega)}{d\omega_0} = \sum_{i=1}^n \omega_1 x_i - y_i + \sum_{i=1}^n \omega_0 = \sum_{i=1}^n \omega_1 x_i - y_i + n\omega_0 = 0$$

$$\therefore \omega_0 = -\frac{1}{n} \sum_{i=1}^n \omega_1 x_i - y_i \text{ and also from eq (1) we again get}$$

$$\omega_0 = -\frac{1}{n} \sum_{i=1}^n \omega_1 x_i - y_i$$

$\Rightarrow$  So, expanding above eq<sup>n</sup> :-

$$\Rightarrow \omega_0 = -\frac{1}{n} (\omega_1(x_1 + x_2 + \dots + x_n) + (-y_1 - y_2 - \dots - y_n))$$

$$\Rightarrow \omega_0 = -(\omega_1 \underbrace{(x_1 + x_2 + \dots + x_n)}_{n} + \underbrace{(-y_1 - y_2 - \dots - y_n)}_{n})$$

$$\Rightarrow \omega_0 = -(\omega_1 \bar{x} - \bar{y}), \text{ where } \bar{x} = \underbrace{x_1 + x_2 + \dots + x_n}_{n} \text{ and}$$

$$\bar{Y} = \underbrace{y_1 + y_2 + \dots + y_n}_{n}$$

$$\Rightarrow \omega_1 \bar{x} + \omega_0 = \bar{Y}.$$

$$\Rightarrow \bar{Y} = \omega_1 \bar{x} + \omega_0 \quad \text{--- (2)}$$

$$\Rightarrow \text{Now from hypothesis of linear regression we know, } \\ y_i^0 = \omega_1 x_i + \omega_0 \quad \text{--- (3)}$$

$$\Rightarrow \text{Comparing eq<sup>n</sup> (3) and (2), we get,} \\ \bar{x} = \bar{x}_i \quad \text{and} \quad \bar{Y} = \bar{y};$$

$\Rightarrow$  Hence, eq<sup>n</sup> (2) shows that least square fit line passes through Co-ordinate  $(\bar{x}, \bar{y})$  and this is our desired result. (Hence Proved).

ANS - (b)

A-② Given: We have 3 variables let's call it as  $x, y, z$ ,  
 So from question we know that 2 variables are highly correlated  
 with 3rd variable i.e.,  $x$  and  $y$  are correlated with  $z$ .  
 Is so, now we have to check whether  $x$  and  $y$  also highly correlated  
 or not.

ii) we will denote correlation factor as  $P_{xy}$ , it tells the correlation  
 between  $x$  and  $y$ .

iii) Now, we start with a claim that  $x$  and  $y$  will not always  
 highly correlated.

iv) formula of Partial Correlation is  $P_{xy.z} = \frac{P_{xy} - P_{xz} \cdot P_{yz}}{\sqrt{1-P_{xz}^2} \sqrt{1-P_{yz}^2}}$

v) Simplify above formula,

$$v) P_{xy} = P_{xy.z} (\sqrt{1-P_{xz}^2} \times \sqrt{1-P_{yz}^2}) + P_{xz} \cdot P_{yz}$$

vi) Now the range of a correlation factor is  $-1$  to  $1$ ,

$$v) \text{ so, } P_{xy} = P_{xz} \cdot P_{yz} \pm \left( \sqrt{1-P_{xz}^2} \right) \left( \sqrt{1-P_{yz}^2} \right)$$

$\therefore$  we have a range for  $P_{xy}$ , ie.,

$$P_{xy} \in \left[ P_{xz} \cdot P_{yz} - \left( \sqrt{1-P_{xz}^2} \right) \left( \sqrt{1-P_{yz}^2} \right), P_{xz} \cdot P_{yz} + \left( \sqrt{1-P_{xz}^2} \right) \left( \sqrt{1-P_{yz}^2} \right) \right]$$

vii) Now we take an example,

viii)  $x$  and  $z$  are highly correlated so we take  $P_{xz} = 0.9$  and similarly

for  $y$  and  $z$  we take  $P_{yz} = 0.85$

$$ix) P_{xy} = \left[ 0.765 - (0.435)(0.5267), 0.765 + (0.435)(0.5267) \right]$$

$$\Rightarrow P_{xy} = \left[ 0.335, 0.994 \right]$$

$\therefore$  Now we can see that we have a lower bound of  $0.335$  which is  
 moderate correlation.

$\therefore$  Our claim is correct that they will not always highly correlated

[Justified]

**Ans - (c)**

A-③ Given:  $x_1, x_2, \dots, x_n$  are iid with finite mean  $\mu$  and variance  $\sigma^2$ , i.e.,  $E(x_1) = E(x_2) = \dots = E(x_n) = \mu$  and  $\text{Var}(x_1) = \text{Var}(x_2) = \dots = \text{Var}(x_n) = \sigma^2$ .

To prove:  $\lim_{n \rightarrow \infty} M_n \rightarrow \mu$ , where  $M_n$  is the sample mean and  $\mu$  is our population mean, i.e.,  $P(|M_n - \mu| \geq \epsilon) \rightarrow 0$ , for any  $\epsilon > 0$

$$\therefore E(M_n) = E\left(\frac{x_1 + x_2 + \dots + x_n}{n}\right) = \frac{E(x_1) + \dots + E(x_n)}{n} \quad [ \because x_i \text{ are iid}]$$

$$\bar{M}_n = \frac{n\mu}{n} = \mu$$

$$\Rightarrow \text{Var}(M_n) = \text{Var}\left(\frac{x_1 + x_2 + \dots + x_n}{n}\right) = \frac{\text{Var}(x_1) + \text{Var}(x_2) + \dots + \text{Var}(x_n)}{n^2}$$

$$\Rightarrow \text{Var}(M_n) = \frac{\sigma^2 + \dots + \sigma^2}{n^2} = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}. \quad [ \because \text{by iid}]$$

$\Rightarrow$  Now by Chebyshev's inequality,

$$\Rightarrow P(|\bar{M}_n - \mu| \geq \epsilon) \leq \frac{\text{Var}(M_n)}{\epsilon^2}$$

$$\Rightarrow \therefore P(|\bar{M}_n - \mu| \geq \epsilon) \leq \frac{\sigma^2}{n\epsilon^2}$$

$\Rightarrow$  Now at  $n \rightarrow \infty$ ,

$$\therefore \lim_{n \rightarrow \infty} P(|\bar{M}_n - \mu| \geq \epsilon) \leq \lim_{n \rightarrow \infty} \frac{\sigma^2}{n\epsilon^2}$$

$$\therefore \text{Now } \lim_{n \rightarrow \infty} \frac{1}{n} \rightarrow 0$$

$$\therefore P(|\bar{M}_n - \mu| \geq \epsilon) \leq 0, \quad \epsilon \rightarrow 0$$

$\Rightarrow$  So, this proves that as  $n \rightarrow \infty$  the sample mean will equal to population mean because probability of  $> 0$  will get 0.

[Hence Proved]

**Pseudo Code :-**

We will take a range of numbers from 0 to 1000 so our population mean will be 500;

Arr = [] // this array will stores all the samples

For i in range 0 to 1000:

    For j in range (i, 1000, 3):

        if(array contains j): // if we already add a sample then we will not add it again  
            // so we will jump to next sample

        continue

        arr.append(j)

    A = Sample\_mean(arr) // Sample\_mean will be a function which gives the mean of samples currently we have in Arr.

If we compare A with population mean in 1 st iteration it will have a high difference but as our function works then we can see A will get closer to population mean that is 500 and after the program is over then our sample mean will be 500 exactly.

**Ans - (d)**

A - Q Derivation of maximum a posteriori (MAP) for Linear regression :-

for linear regression we have,  $y = \omega^T x + b$ .

As, from the definition of MAP we have,

$$P(\omega|D) = P(D|\omega) * P(\omega) \quad (\text{by Bayes theorem})$$

∴ So, we have to find most likely  $\omega$  given D and  $P(D)$  does not depend on  $\omega$  and therefore it doesn't affect in the optimization.

∴ Reduced form of MAP we have is  $P(D|\omega) * P(\omega)$

Now,

$$\underset{\omega}{\operatorname{argmax}} \ P(D|\omega) * P(\omega)$$

y

$$y = \omega^T x + b,$$

∴ Now, we know  $b \sim N(0, \sigma^2)$  (given),

$$\text{So, } y \sim N(\omega^T x, \sigma^2)$$

$$\underset{\omega}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | \omega) * P(\omega)$$

∴ Taking logarithm both sides

$$\rightarrow \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n \log P(y_i | \omega) + \sum_{i=1}^n \log P(\omega)$$

From Gaussian distribution,

$$P(y_i | \omega) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} (\omega^T x_i - y_i)^2}$$

$$\therefore \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (\omega^T x_i - y_i)^2 + \sum_{i=1}^n \log P(\omega)$$

Now  $\log \frac{1}{\sqrt{2\pi\sigma^2}}$  - h constant we can ignore because it doesn't affect our optimization,

$$\therefore \underset{\omega}{\operatorname{argmax}} \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (\omega^T x_i - y_i)^2 + \sum_{i=1}^n \log P(\omega) \right)$$

$$\therefore \underset{\omega}{\operatorname{argmax}} - \left( \frac{1}{2\sigma^2} \sum_{i=1}^n (\omega^T x_i - y_i)^2 - \sum_{i=1}^n \log P(\omega) \right)$$

$$\therefore \underset{\omega}{\operatorname{argmin}} \left( \frac{1}{2\sigma^2} \sum_{i=1}^n (\omega^T x_i - y_i)^2 - n \log P(\omega) \right)$$

Now if  $P(\omega)$  is constant then our result will be same as MLE but here  $\omega \sim \text{Normal}(0, \sigma_w^2) \rightarrow \text{Given}$

$$\therefore P(\omega) = \frac{1}{\sqrt{2\pi\sigma_w^2}} \times e^{-\frac{1}{2\sigma_w^2} (\omega)^2}$$

$$\therefore \underset{\omega}{\operatorname{argmin}} \left( \frac{1}{2\sigma^2} \sum_{i=1}^n (\omega^T x_i - y_i)^2 - n \log \frac{1}{\sqrt{2\pi\sigma_w^2}} + \frac{n\omega^2}{2\sigma_w^2} \right)$$

neglect again

∴ Our final result is :-

$$\boxed{\omega = \underset{\omega}{\operatorname{argmin}} \left( \frac{1}{2\sigma^2} \sum_{i=1}^n (\omega^T x_i - y_i)^2 + \frac{n\omega^2}{2\sigma_w^2} \right)}$$



## SECTION-B

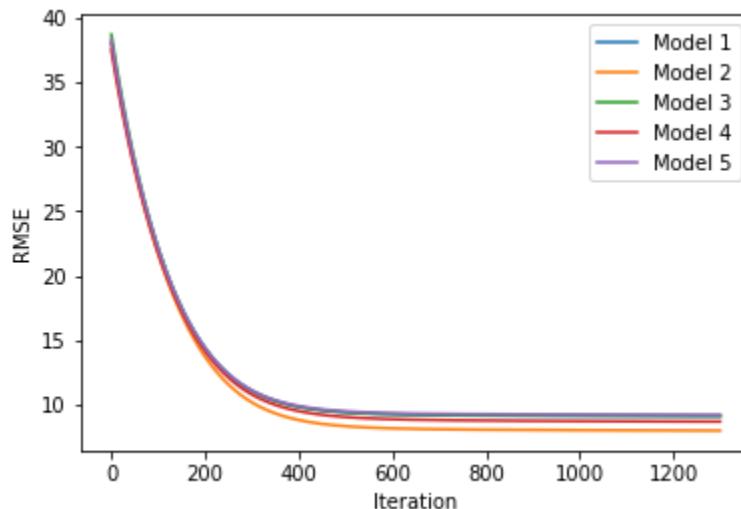
Ans - (a) Table of K value with their Mean RMSE :-

K value	Mean RMSE
2	8.909685034963294
3	8.882254992021346

4	8.84270311014231
5	8.729345090256059

Optimal value of K will be 5 because the mean rmse is lowest among all the mean rmse we get for other values of K.

(b) RMSE Vs Iteration plots for each model we get from Optimal value of K which is 5 :-



For training the models we use linear regression with learning rate = 0.006 and epoch = 1300.

Model Number	RMSE on Train Set	RMSE on Val set
1	9.074522911193775	7.835951654646455
2	7.962530168794994	11.747048375179345
3	9.111364677580731	7.506690888508186
4	8.667331486337238	9.461115568391167
5	9.188713471841982	7.095918964555147

Mean RMSE on Train Set = 8.800892543149743

Mean RMSE on Val set = 8.729345090256059

(c) For Lasso Regularization :-

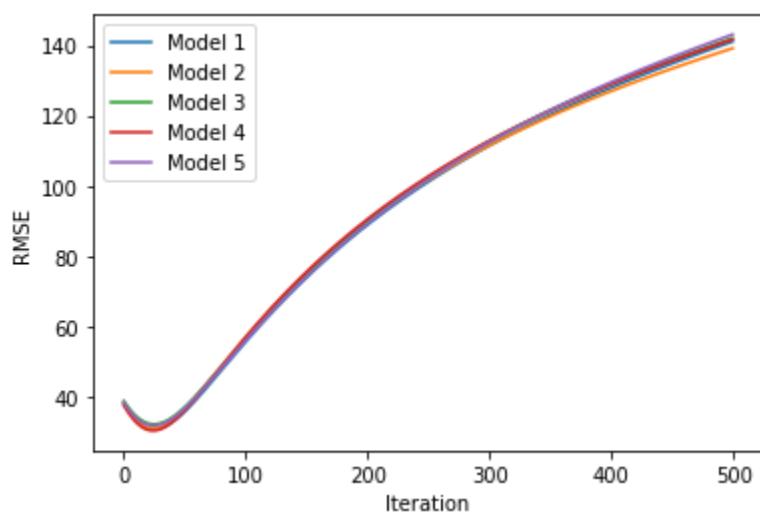
We Fix our learning rate = 0.006 and changes only epochs and lambda (Regularization Parameter)

S NO.	Lambda (Regularization Parameter)	Epochs	Mean RMSE On Train Set	Mean RMSE On Val Set
1	0.05	500	141.453963931	143.523323653

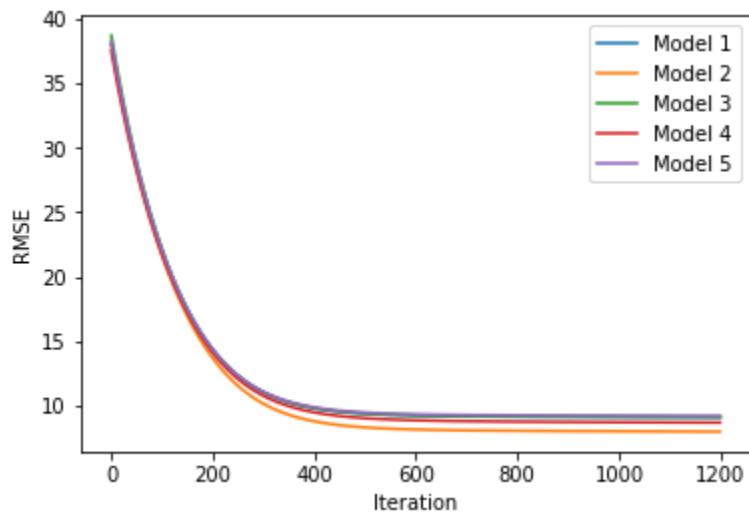
			5704	37116
2	0.0001	1200	8.82530507644 7638	8.75499655476 32
3	1	200	1834.26285969 01934	1843.34860396 726
4	0.2	600	606.608671257 1429	616.4692531121 339
5	0.009	700	31.1206143914 362	31.7639493888 7838

### Plot for Each Experiment :-

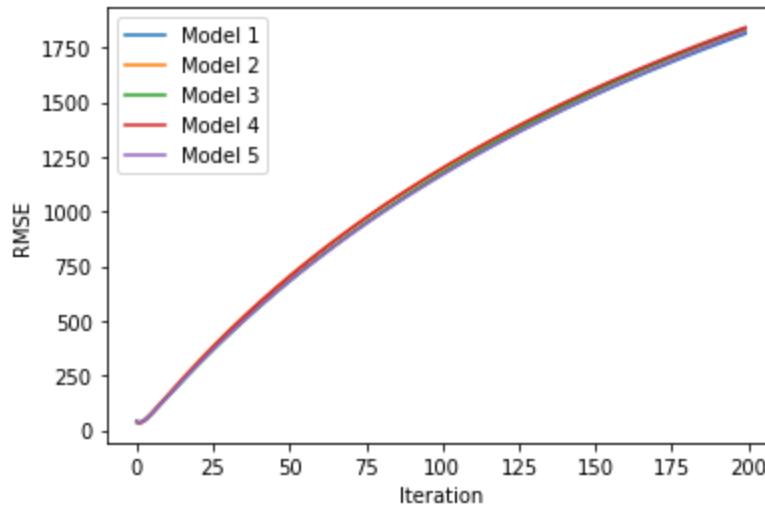
1.



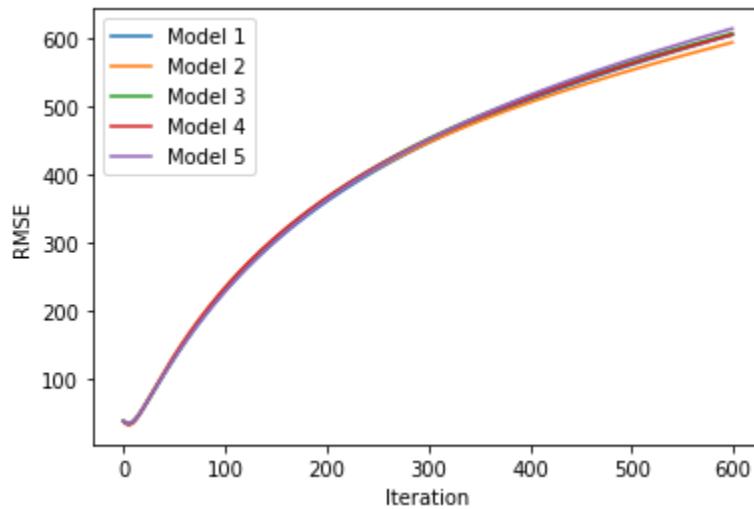
2.



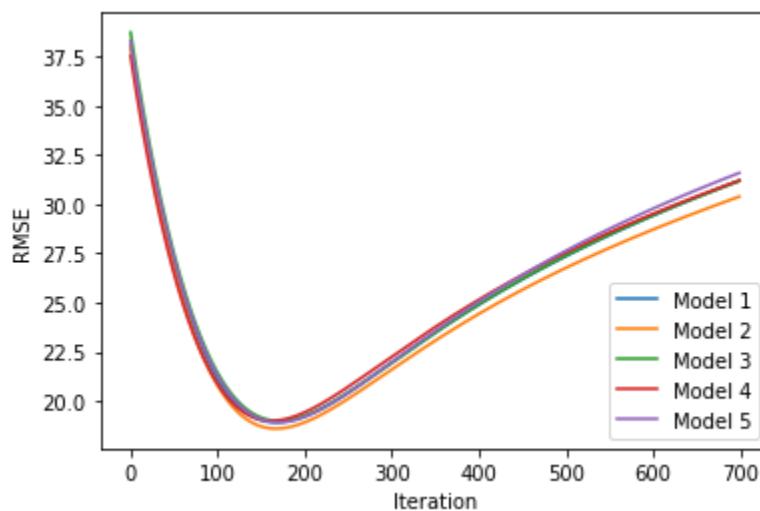
3.



4.

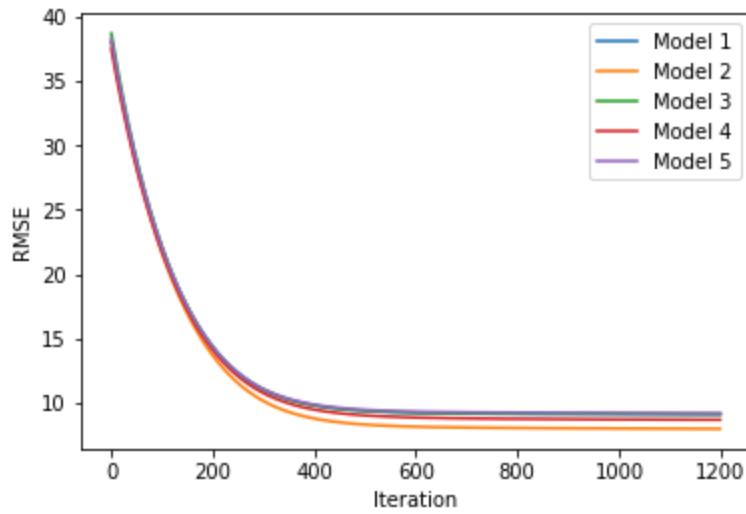


5.



**So we have an optimal solution for lambda = 0.0001 and RMSE on train set = 8.825305076447638 and RMSE on Val Set = 8.7549965547632**

**Plot for lambda = 0.0001 :-**



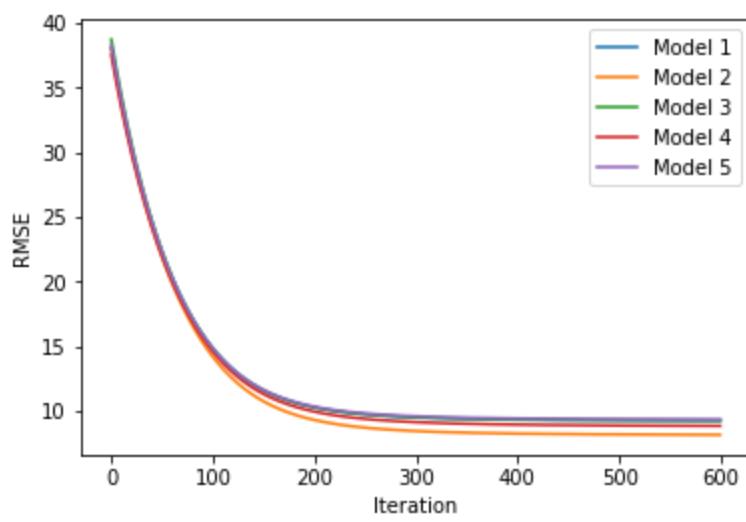
**For Ridge Regularization :-**

**We Fix our learning rate = 0.006 and changes only epochs and lambda (Regularization Parameter)**

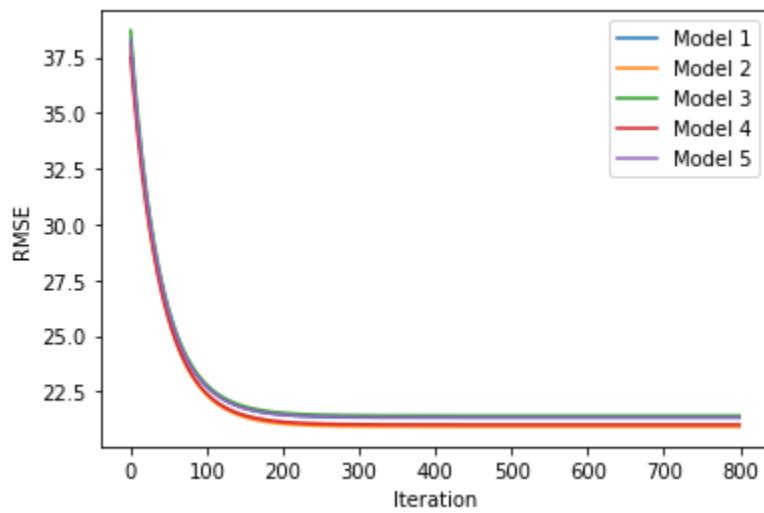
S.No	Lambda (Regularization Parameter)	Epochs	Mean RMSE On Train Set	Mean RMSE On Val Set
1	0.05	600	8.99991898949 606	8.93471306283 1266
2	1	800	21.20021871150 5393	21.20021871150 5393
3	0.0055	1350	8.76836452098 2155	8.69599974745 391
4	10	500	36.5961740832 4361	36.6241257777 0649
5	5	1000	33.5967869677 2319	33.6487654812 03924

**Plot for Each Experiment :-**

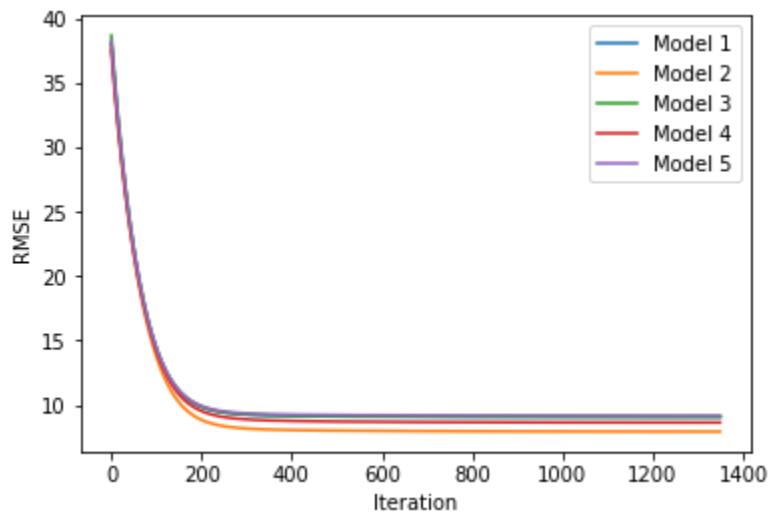
1.



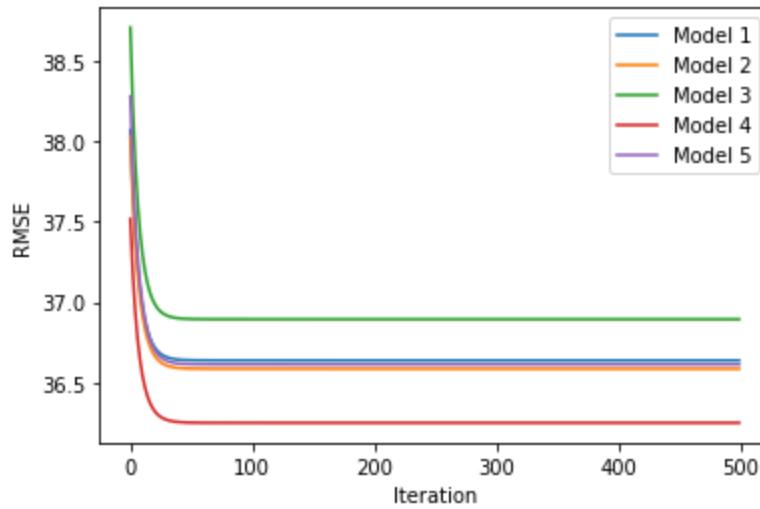
2.



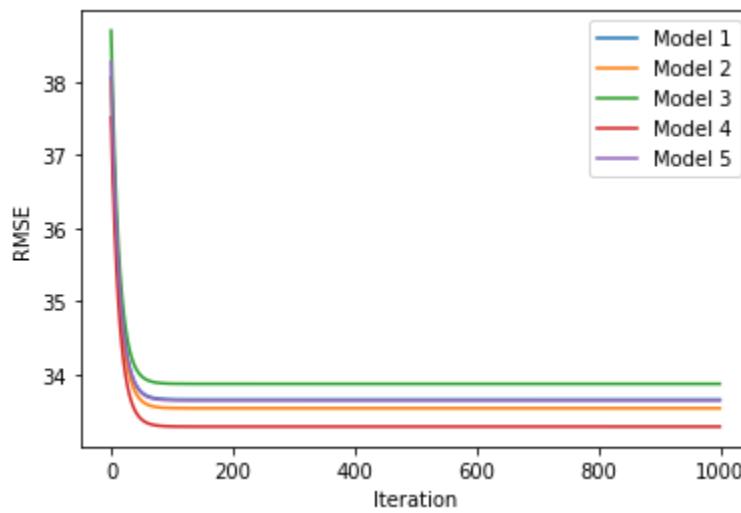
3.



4.

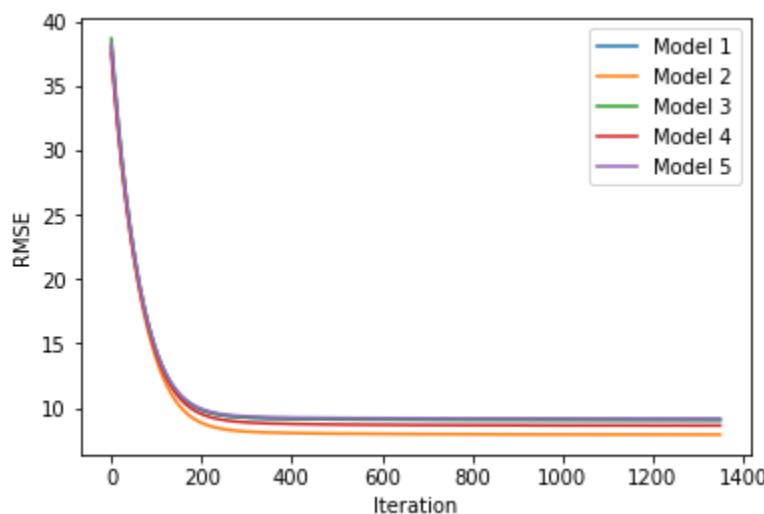


5.



**So we have an optimal solution for lambda = 0.0055 and RMSE on train set = 8.768364520982155 and RMSE on Val Set = 8.69599974745391**

**Plot for lambda = 0.0055 :-**



**(d) RMSE for each fold by normal equation :-**

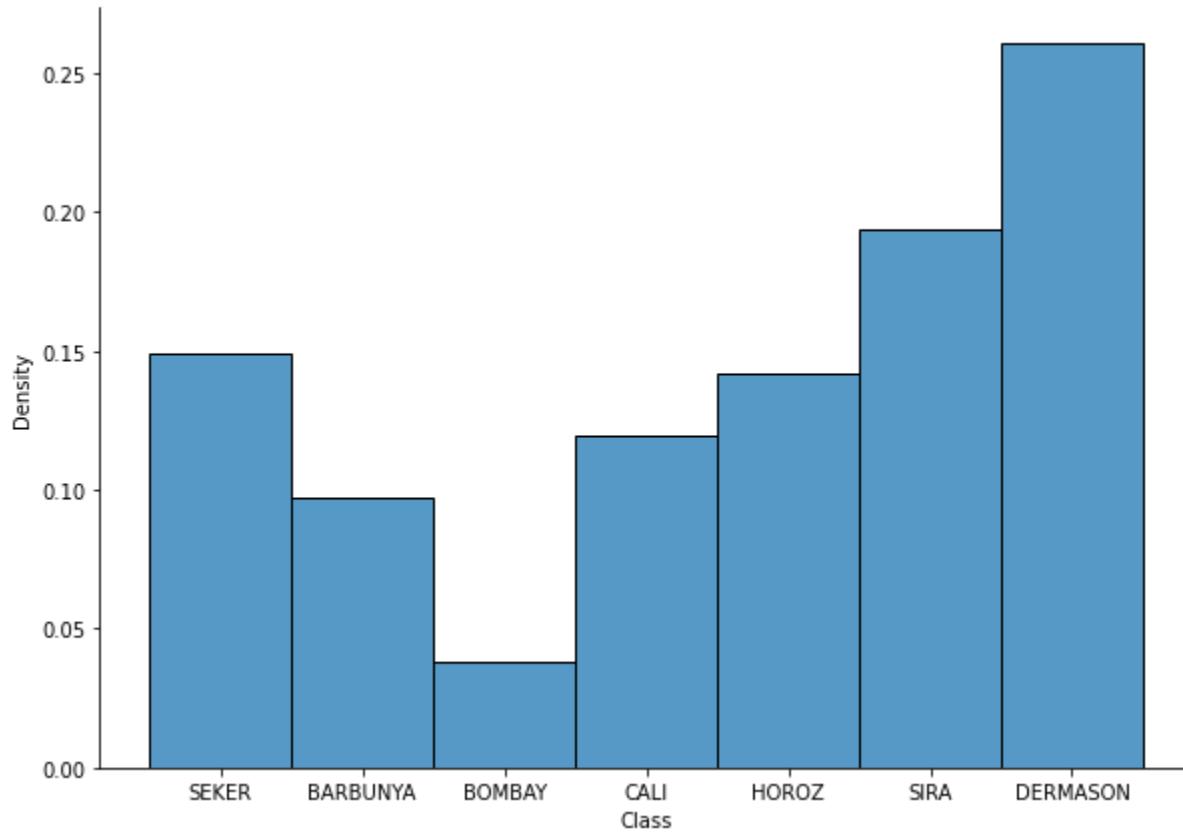
Model Number	RMSE
--------------	------

1	7.707739342790347
2	11.65785008339692
3	7.549243481223661
4	9.434235638731744
5	7.101061841935251

**Mean RMSE From Normal Equation = 8.690026077615583**

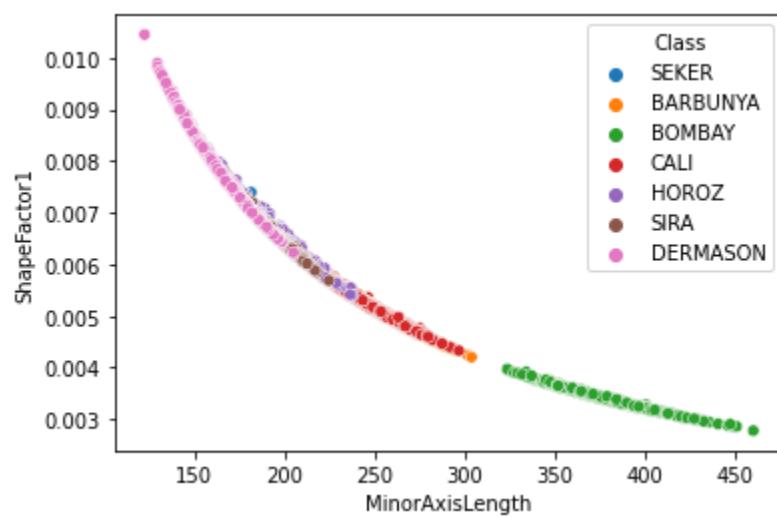
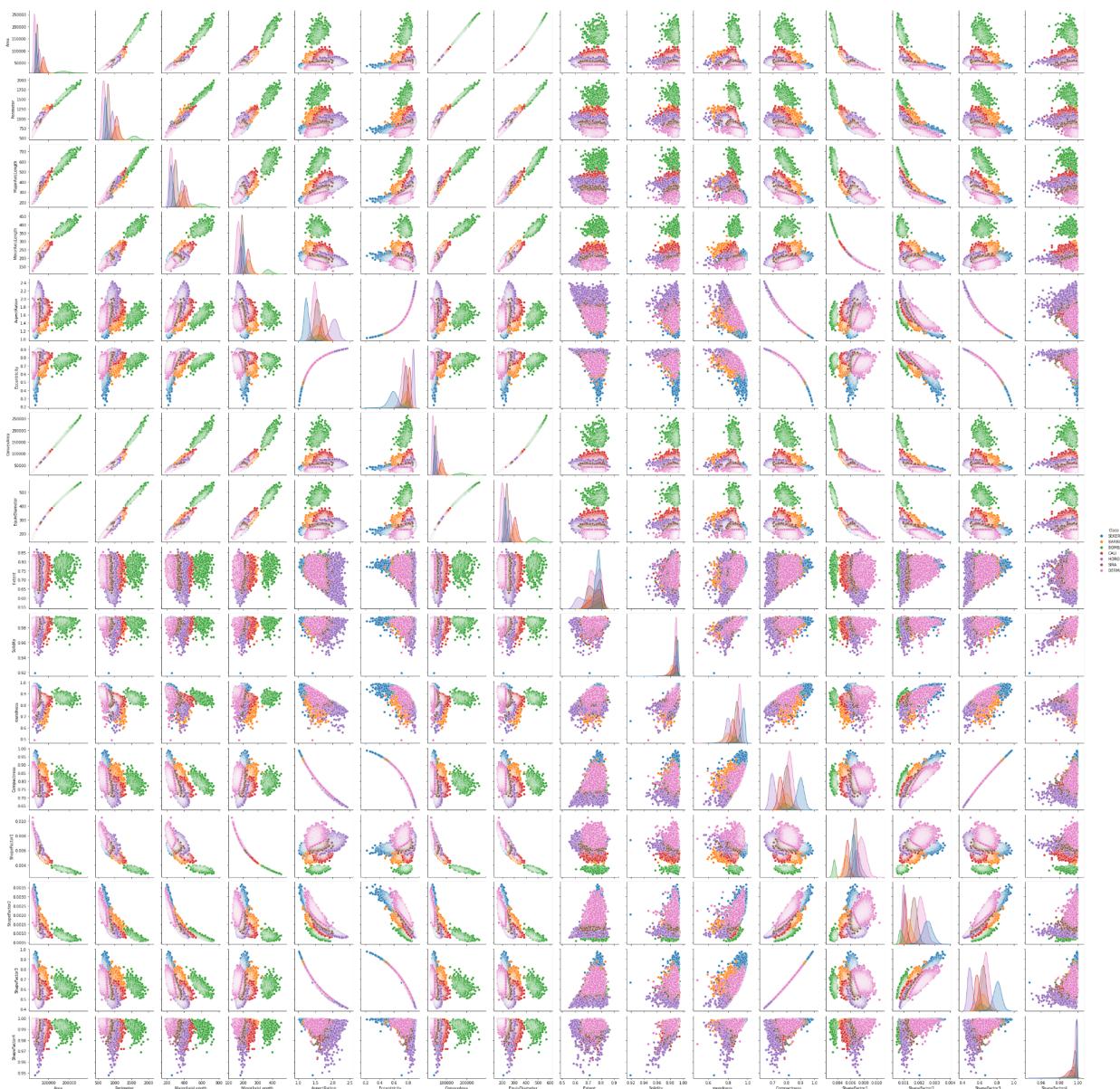
## SECTION-C

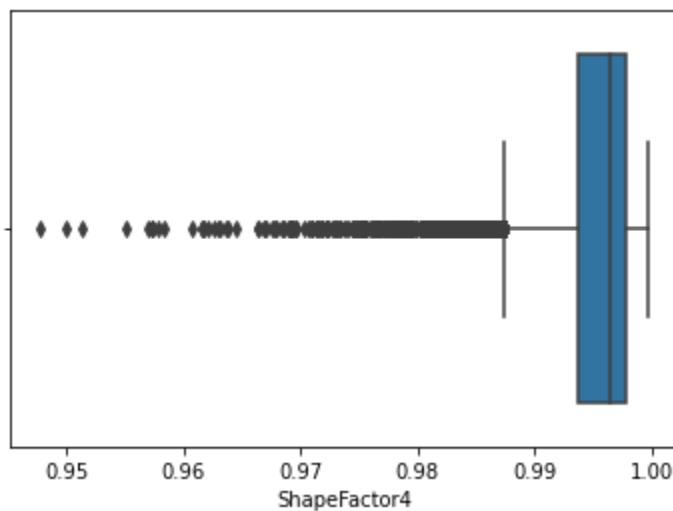
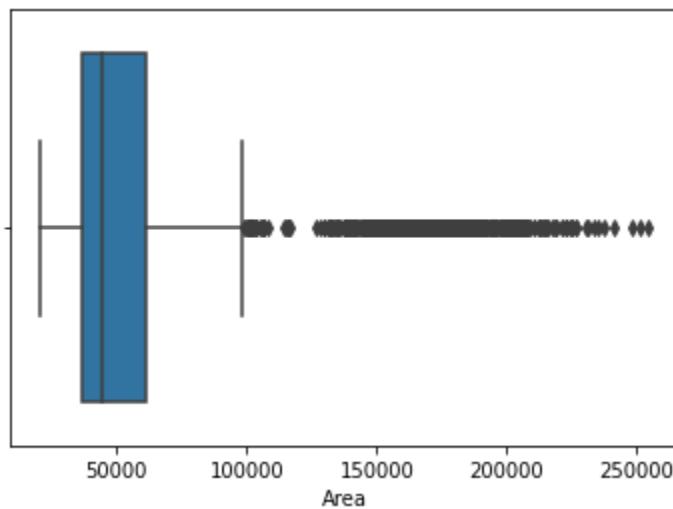
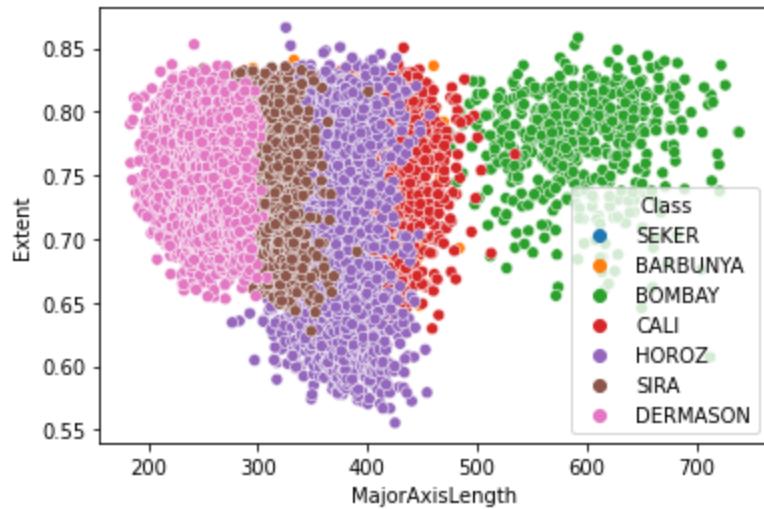
(a) Plot for Distribution which we get :-

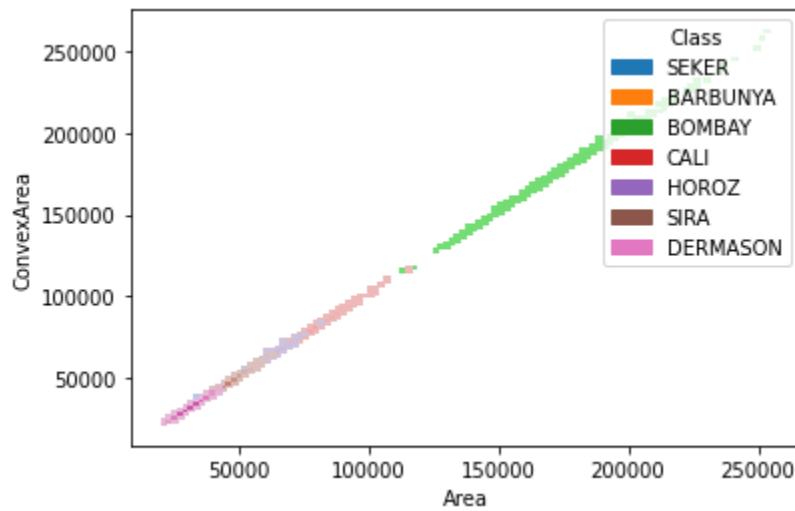
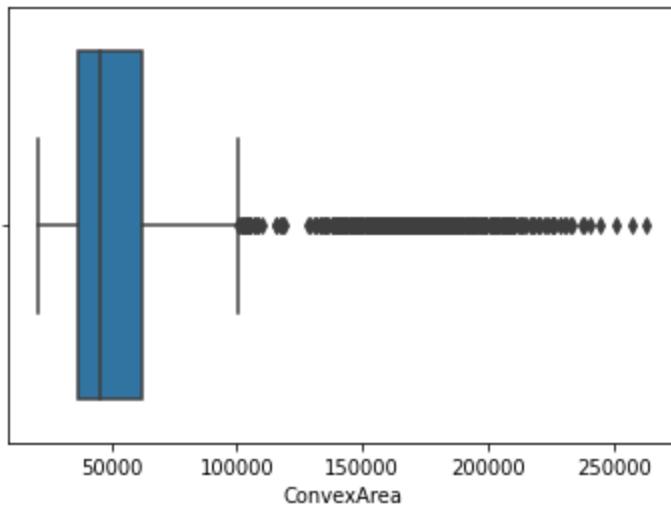


**Analysis :-** Demarson Class has the highest density than all other classes and Bombay class has the lowest density among all classes.

(b) After Plotting different types of graph like :-





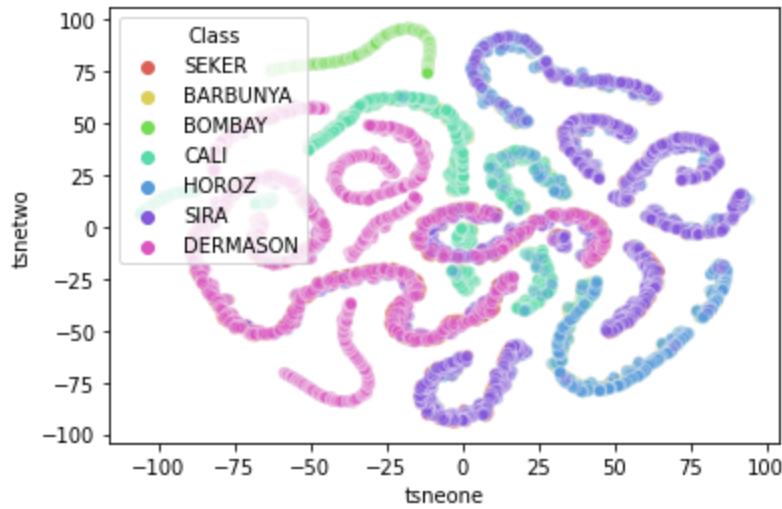


#### **Five insights on the data :-**

1. Bombay Class parameters are mostly non overlapping with other classes which means it is easy to classify from other classes.
2. There are many attributes which have a relation as we can see in the graph that some parameters make a linear relation some make exponential and etc.
3. There are many attributes which have many outliers which can affect our prediction.
4. Highest standard deviation seen in Convex Area.
5. ConvexArea and Area are the two parameters which have the bulky data in terms of their value.

**No missing value has been reported when we check for missing values in data.**

**(c) After using TSNE algorithm on our data we have a plot like this :-**



#### **Analysis on Separability of the data:-**

**Only Bombay Class has an non overlapping distribution with other classes otherwise all other classes are overlapped with each other.**

#### **(d) We chooses the GaussianNB and BernoulliNB for implementing Naive Bayes :-**

Recall for GaussianNB = 0.8968049944913699

Precision for GaussianNB = 0.8968049944913699

Accuracy for GaussianNB = 0.8968049944913699

Recall for BernoulliNB = 0.7245684906353287

Precision for BernoulliNB = 0.7245684906353287

Accuracy for BernoulliNB = 0.7245684906353287

GaussianNB has a higher score than BernoulliNB in all three types of score which shows that GaussianNB naive bayes model will perform better than BernoulliNB naive bayes model for a similar type of data.

#### **(e) We uses GaussianNB naive bayes model for calculating the scores after applying PCA with number of components = 4,6,8,10,12**

Accuracy for PCA with number of components 4 = 0.8670583914799853

Recall for PCA with number of components 4 = 0.8670583914799853

Precision for PCA with number of components 4 = 0.8670583914799853

F1 score for PCA with number of components 4 = 0.8670583914799853

Accuracy for PCA with number of components 6 = 0.9008446566287184

Recall for PCA with number of components 6 = 0.9008446566287184

Precision for PCA with number of components 6 = 0.9008446566287184

F1 score for PCA with number of components 6 = 0.9008446566287184

Accuracy for PCA with number of components 8 = 0.9030481087036357

Recall for PCA with number of components 8 = 0.9030481087036357

Precision for PCA with number of components 8 = 0.9030481087036357

F1 score for PCA with number of components 8 = 0.9030481087036357

Accuracy for PCA with number of components 10 = 0.8879911861917004

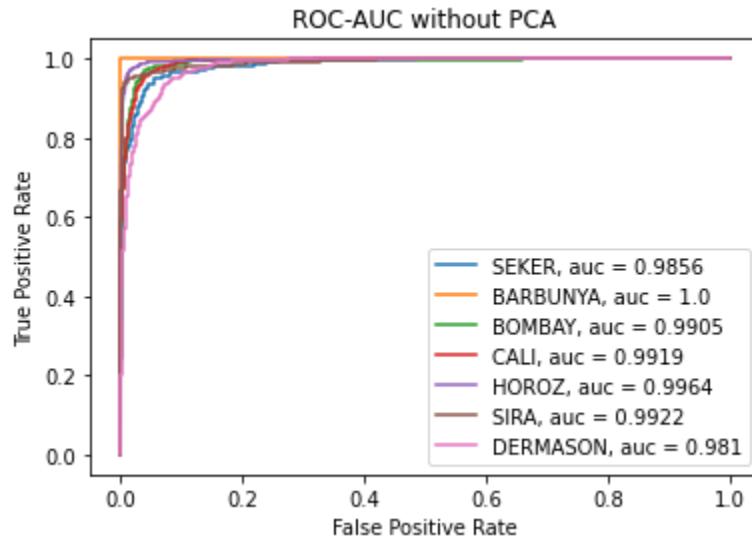
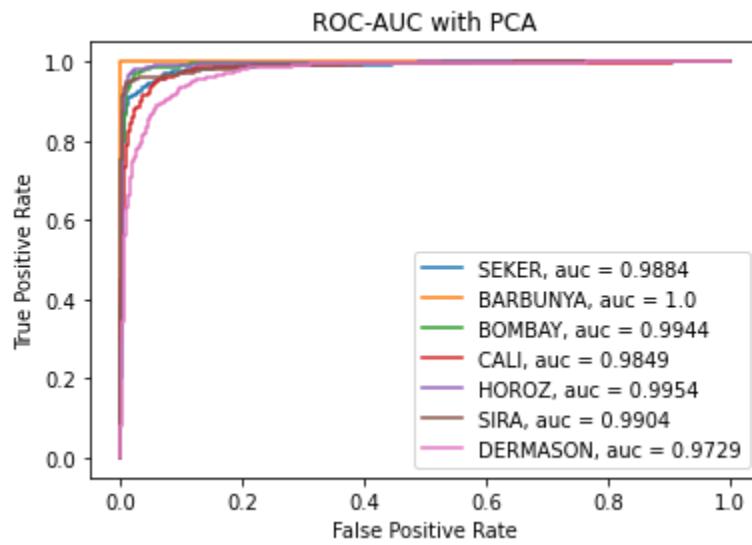
Recall for PCA with number of components 10 = 0.8879911861917004

Precision for PCA with number of components 10 = 0.8879911861917004  
F1 score for PCA with number of components 10 = 0.8879911861917004

Accuracy for PCA with number of components 12 = 0.8872567021667279  
Recall for PCA with number of components 12 = 0.8872567021667279  
Precision for PCA with number of components 12 = 0.8872567021667279  
F1 score for PCA with number of components 12 = 0.8872567021667279

**When we compare the scores which we get from all the models after applying different PCA values, we can see that for number of components = 8 we have the highest score in all four types of score from others which shows that we have better performance if we use PCA with 8 components.**

(f)



**From both the roc-auc curve (with pca and without pca) we can see that BARBUNYA Class has the highest true positive rate.**

(g) We will select only those features which hasn't any direct relationship between them like in (b) part we see that there are many parameters which are dependent linearly or exponentially or other types of relation they have so we will remove all that attributes from the data, while training and testing our Logistic Regression Model.

**Logistic Model Score = 0.9280205655526992**

**Naive Model Score = 0.8968049944913699**

We can see that the Logistic Model performs better after removing unnecessary features from the data than the Naive Bayes Model.

**Pair plot after removing dependent Features :-**

