

# Assignment 1 ML

Vaibhav Jaiswal (2020547)

## Question 1:

A) The proof can be viewed in the picture below :

Q.1 Least square fit line can be defined by  $y_i = w_0 + w_1 x_i$ , and  $y_i = \bar{y} + E$  ( $E \rightarrow \text{error}$ )

we have been given that the least square fit line passes through the point  $(\bar{x}, \bar{y})$ .

Based on the definition of  $\bar{y}$

$$\bar{y} = \frac{1}{n} \sum y_i = \frac{1}{n} \sum (y_{i1} + E)$$
$$\bar{y} = \frac{1}{n} (\sum y_{i1} + \sum E) \quad \text{This will become zero}$$

Since, expected error value is zero, we can say that  $\sum E = 0$ .

$$\bar{y} = \frac{1}{n} \left( \sum (B_0 + B_1 x_i) \right) \quad \text{constant for all.}$$
$$\bar{y} = \frac{1}{n} (\sum B_0 + \sum B_1 x_i) \quad [\sum B_0 = n \times B_0]$$
$$\bar{y} = \frac{n}{n} B_0 + B_1 \left( \frac{1}{n} \sum x_i \right). \quad \text{Since, } \frac{1}{n} \sum x_i = \bar{x}$$

∴ Thus, the point  $(\bar{x}, \bar{y})$  lies on given equation (I).

B) Thus, the best way to analyse how the correlation will function is to take an example. On taking an example of variables X, Y and Z. Now, Y is highly correlated to X and Z. But, X and Z are not relatively correlated to each other. Taking the example, that the correlation between Y and X is 0.7 and correlation between Y and Z is 0.75. Between X and Z, correlation is 0.05.

C)

PseudoCode :

1. Take a empty list x
2. Take a for loop for values ranging from 0, python
3. the value inside y will be generated randomly using random function in
4. Apply gaussian probability distribution formula and append it onto the list.
5. initialise another variable n
6. take another loop, and add the index of that loop into variable n.  $n = n + x[ ]$ .
7. We return n at the end.

Q.1 (c) WEAK LAW OF LARGE NUMBERS:
Given $n$ independent and identical distribution, $x_1, x_2, \dots, x_n$ . Standard error $\sigma$ , $E[x_i] = \mu$
$\bar{x} = \frac{(x_1 + x_2 + x_3 + \dots + x_n)}{n}$
$\lim_{n \rightarrow \infty} P[ \bar{x} - \mu  \geq \epsilon] = 0$
Let $x_1, x_2, x_3, \dots, x_n$ be random variables.
$E[x_i] = \mu$ so $\bar{x} = \frac{(x_1 + x_2 + x_3 + \dots + x_n)}{n}$
$\bar{x} = \frac{n\mu}{n} = \mu$
$\text{var}(x_i) = \sigma^2$
$M_n = \frac{1}{n} \sum_{i=1}^n x_i$
$E[M_n] = E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} E[x_1 + x_2 + \dots + x_n]$
$= \frac{n\mu}{n} = \mu$ .

8	$\text{var}(\bar{x}) = \frac{\text{var}(x_1 + x_2 + \dots + x_n)}{n^2}$
9	$= \frac{\sigma^2(n)}{n^2} = \frac{\sigma^2}{n}$
10	using chebyshev's inequality,
11	$P( \bar{x} - \mu  \geq \epsilon) \leq \frac{\text{var}(\bar{x})}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2}$ , $\forall \epsilon > 0$ .
12	Since, $n \rightarrow \infty$ ,
13	$\lim_{n \rightarrow \infty} P[ \bar{x} - \mu  \geq \epsilon] = \lim_{n \rightarrow \infty} \frac{\sigma^2}{n\epsilon^2} = 0$ .

D)

VK 35 • 238-127	AUGUST	26
THURSDAY		
M T W T F S S M T W T F S S 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30		
0.1 (d) $y = w^T x + \epsilon$ (as taught in class).		
$\epsilon \sim N(0, \sigma^2)$ . (noise follows normal distribution) $y \sim N(w^T x, \sigma^2)$ .		
Probability of gaussian distribution, when calculating probability of each response variable.		
$P(y x, w) = N(y   w^T x, \sigma^2)$ . $= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(y - w^T x)^2}{\sigma^2}}$		
Given data $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ .		
Assuming, gaussian prior distribution over the weight $w$ .		
$P(w) = N(w   0, \lambda^{-1} I)$ $= \frac{1}{\sqrt{2\pi\lambda}} e^{-\frac{1}{2} \frac{\lambda w^T w}{\lambda}}$		
$\log P(w D) = \log \frac{P(w)P(D w)}{P(D)} = \log \frac{P(w)}{P(D)} + \log \frac{P(D w)}{P(D)}$		

8 Maximum - A - Posterior solution:

9

10  $P(w | \vec{x}_i, y_i) = \frac{\sum_i P(y_i | \vec{x}_i, w) P(w)}{P(y_i | x)}$

11  $\text{argmax}_w \prod_{i=1}^n P(w | \vec{x}_i, y_i)$

12  $= \text{argmax}_w \frac{\prod_i P(y_i | \vec{x}_i, w) P(w)}{P(y | x)}$

2 Taking log - posterior:

$w_{MAP} = \text{argmax}_w \log P(w | D)$

3  $= \text{argmax}_w \left\{ \sum_i \log \left( \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left( \frac{y_i - w^T \vec{x}_i}{\sigma} \right)^2} \right) + \log P(w) \right\}$

4

5

6  $= \text{argmax}_w \left\{ \sum_i \frac{1}{2} \left( \frac{y_i - w^T \vec{x}_i}{\sigma} \right)^2 + \frac{1}{2} \log(2\pi) - \frac{\lambda}{2} w^T w \right\}$

7 Simply ignoring Independence assumption

8  $= \text{argmin}_w \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T \vec{x}_i)^2 + \frac{\lambda}{2} w^T w$ .

## Question 2 :

Part A →

The optimal value of K is at K = 5. The values of K with the Average RMSE can be seen in the table given below.

K	RMSE Means for K
2	1.0000046405180678
3	0.9999995260347649
4	0.9999995260347649
5	0.9999995260347649

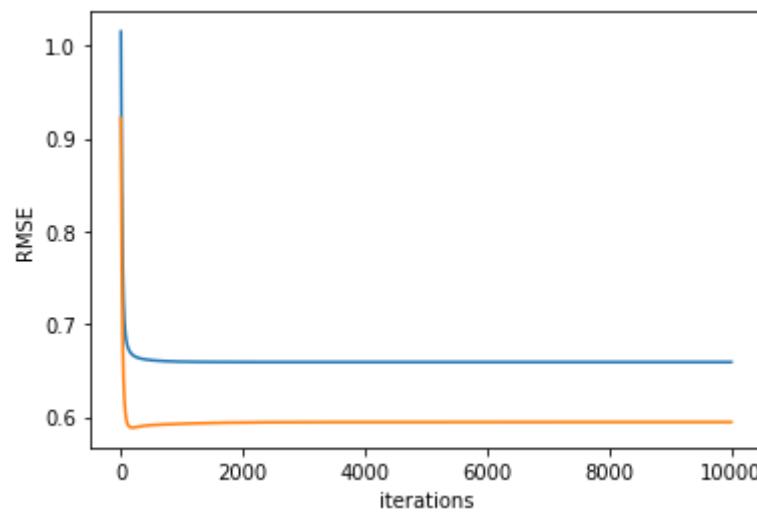
```

1 RMSEKvalues = []
2 for i in range(2,6):
3     print(i)
4     RMSEKvalues.append(train(kfold(df,i),i))
5
6 print(RMSEKvalues)

[1.0000046405180678, 0.9999995260347649, 1.0001641050721342, 0.9995186621625347]

```

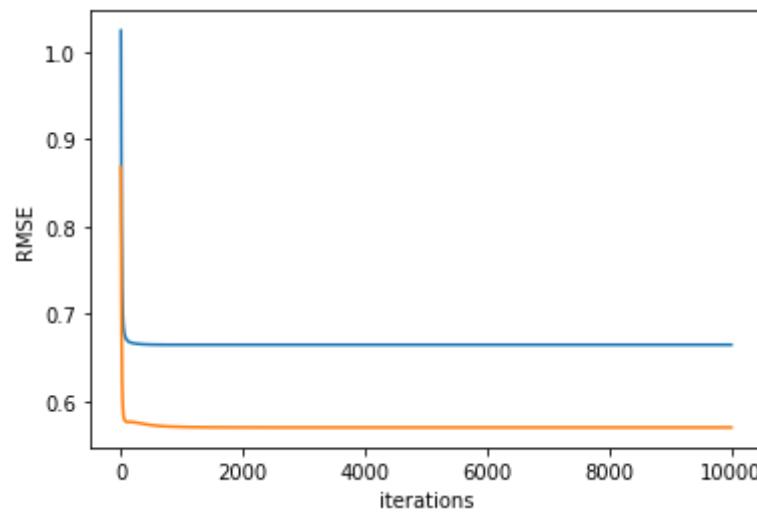
Part B →



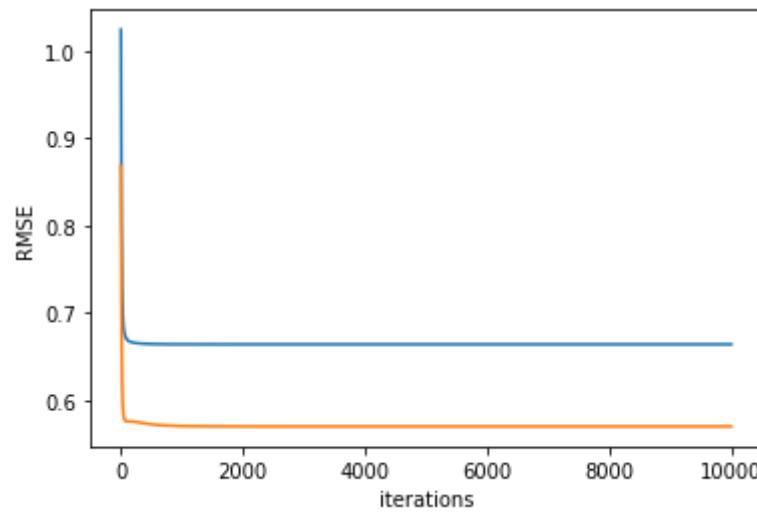
### Part C →

Tried different values of alpha, the values of alpha (or a comes) to be equal to 0.01

Lasso Regression :



Ridge Regression :



### Part D →

The normal form directly gives the optimal parameter *wparameter* for the equation. This will be informed for the optimal value of K. The optimal Value is 5.

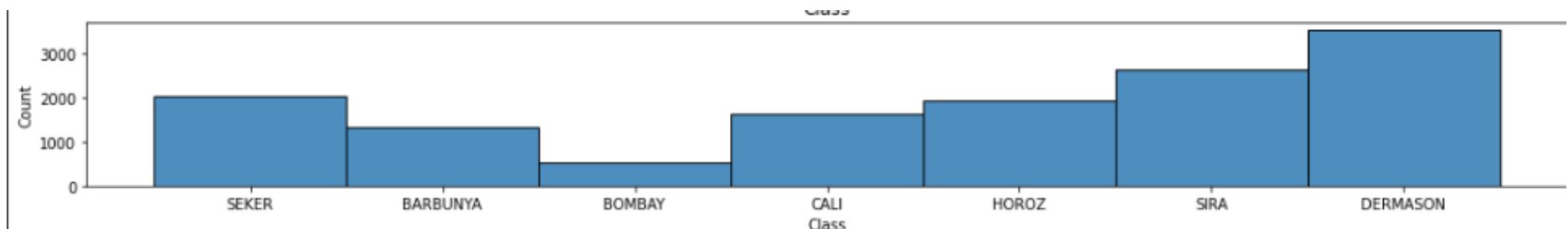
For the 5 validation sets generated using K = 5, we get the set as shown below.

```
[1.045289278556397,
 1.0738903617321636,
 1.1330584321885546,
 1.259860933441022,
 1.0302762014828666]
```

## Question 3 :

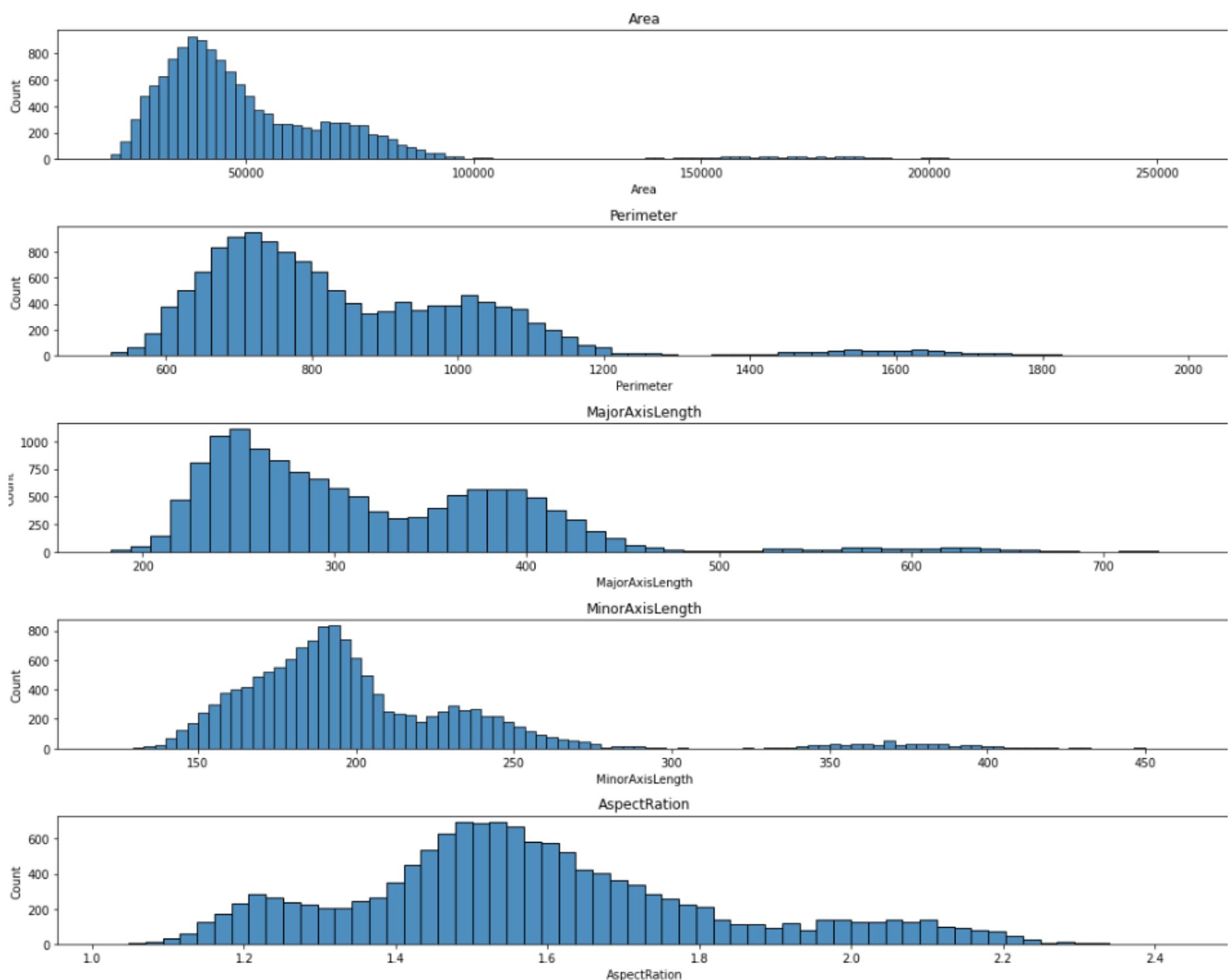
### Part A →

The class distribution for different dry beans set was plotted



It was observed that the dermason class of dry beans had a greater number of count in the dataset. The count for the different types is not equally divided, and is unequally distributed. Other than this more distributions were made on other features

```
'Area' 'Perimeter' 'MajorAxisLength' 'MinorAxisLength' 'AspectRatio'
'Eccentricity' 'ConvexArea' 'EquivDiameter' 'Extent' 'Solidity'
'roundness' 'Compactness' 'ShapeFactor1' 'ShapeFactor2' 'ShapeFactor3'
'ShapeFactor4']
```



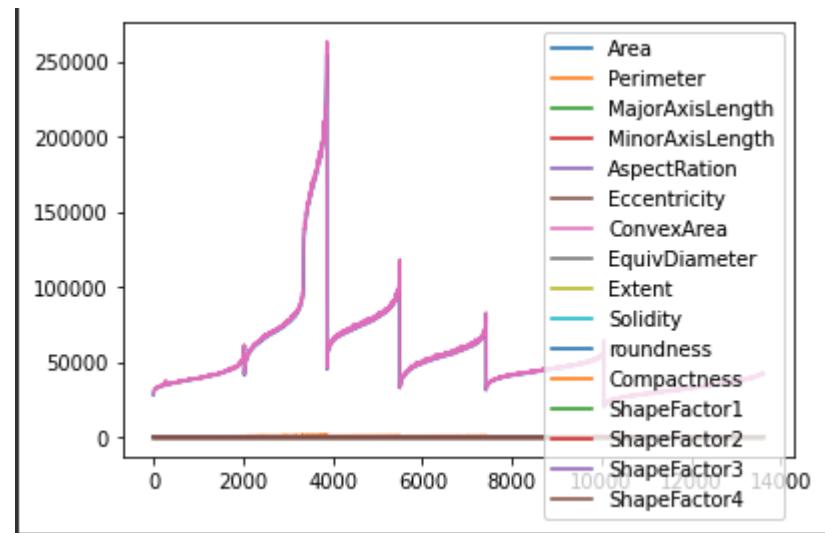
The observation was that for most of the classes the values are distribution is skewed in nature. Data is not equally distributed.

## Part B (Exploratory Data Analysis) →

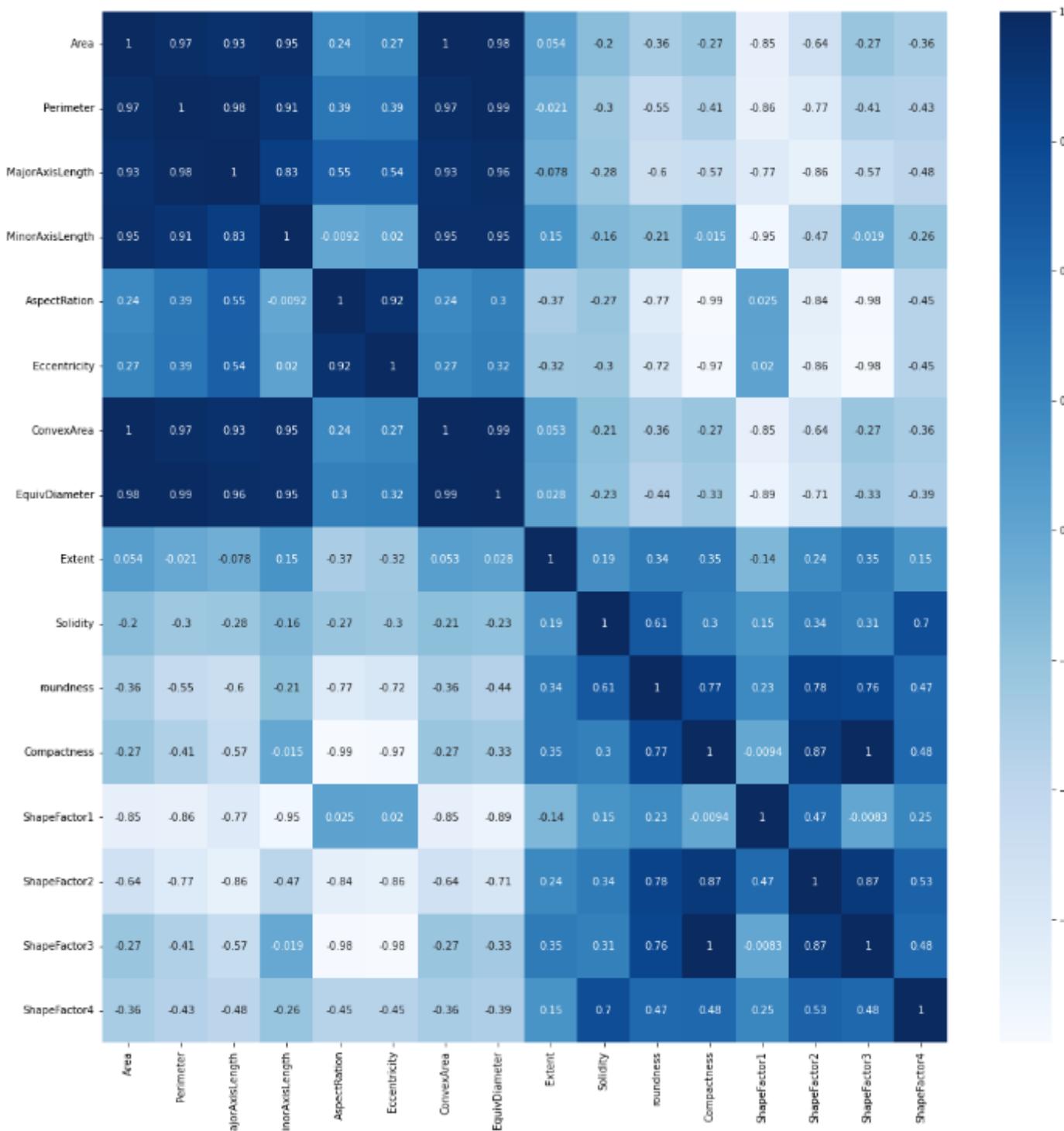
### Data Insight 1 : Plotting →

On plotting the data, it can be seen that the data is skewed, and is not consistent. The convex area and area, are related can also be also observed from the graph itself. The plotting for various graphs are consistent as

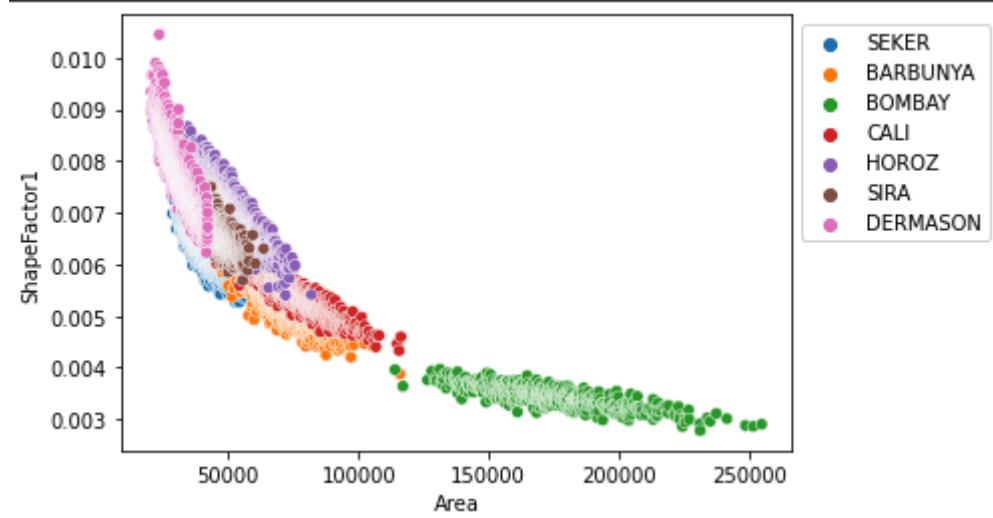
well. The max values can be seen from the graph as well.



### Data Insight 2 : Correlation Map →



- On Using heatmap, correlation in the data can be clearly observed. It is well understood that features like convex area, and area will be correlated. Features like perimeter, major axis length, convex area, equivdiameter are related to each other. Other than this, shapefactor 3, 2 Shape and Compactness have a strong correlation. Roundness and Shape factor 2,3 also have a strong correlation.
- Similar observations became more apparent for specific conditions when using scatter plot. This also shows it in respect of various types of classes.

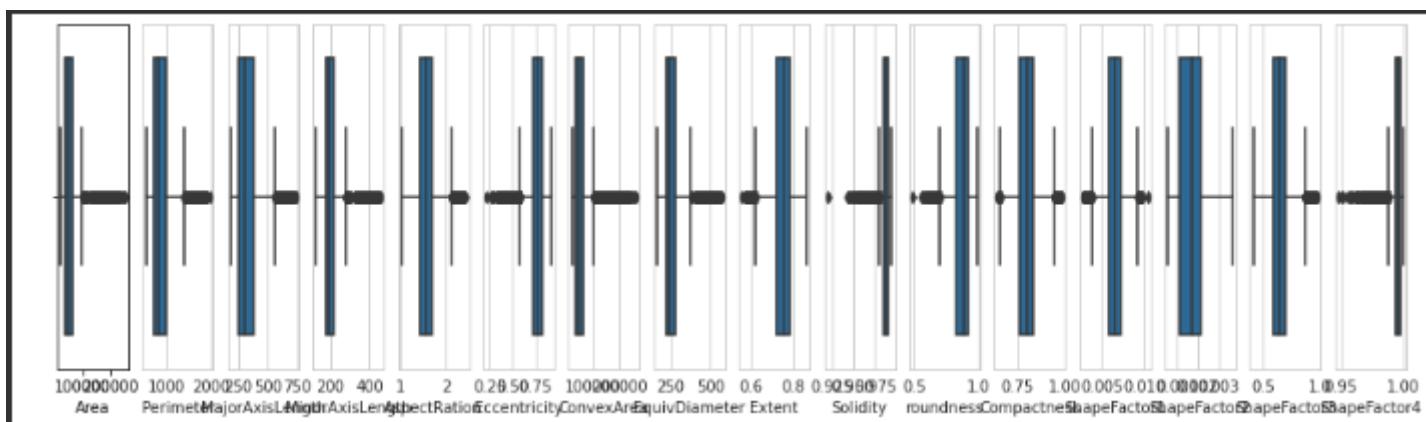


### Data Insight 3 : Plotting Data Distribution →

Using Histogram, we can see the different distributions of data and its classes. It is observed that the data is skewed in nature. Specifically shape factor 4 and more are highly skewed in nature.

### Data Insight 4 : Box Plot →

The box plot shows that the data has a massive number of outliers, which is represented by the black dots in the graph. The features like extent, still have reduced number of outliers, whereas things like shape-factor have massive number of features.

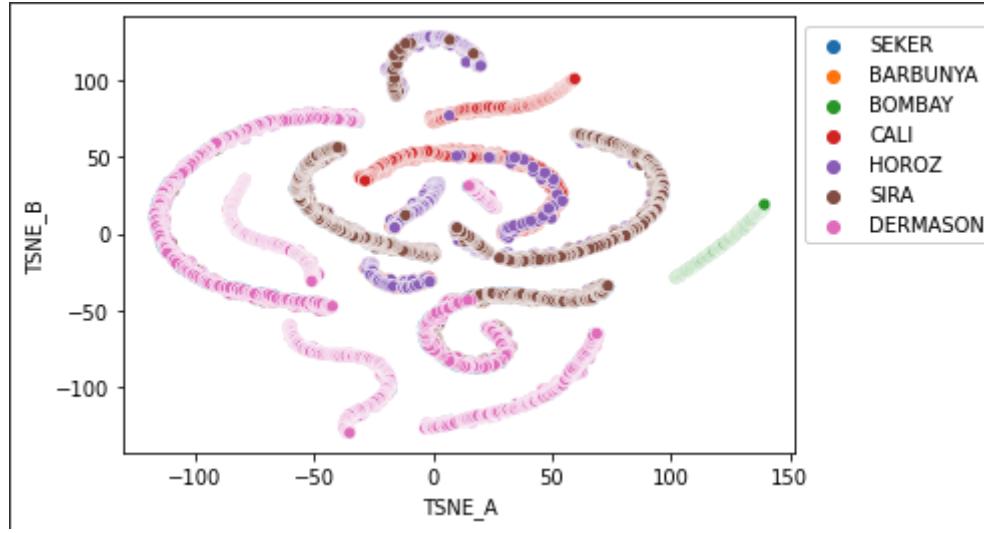


### Data Insight 5 : Exploring the classes and data →

Using commands like df.head and df.info(), df.shape, on the command, what we get the data, which gives things like the dimensions. The non-null value reflects that there are NO MISSING VALUE in the dataset. It can also be observed that all the values are in the form of Dtype of float64, things like area and convex area are in the form of int64. And one is in the form of object. The Data is of the dimensions 13611 entries to 17 columns. Memory usage is 1.8+ MB.

```
<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   AspectRatio       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  
dtypes: float64(14), int64(2), object(1)
memory usage: 1.8+ MB
```

### Part C (T-SNE) →



On using a perplexity value of 50 (since the data is decently large), and the for a number iterations = 4000. The data is separable, with maximum of that of dermason and bombay. The values still have outliers, that overlap with other classes of dry beans

#### Part D (Naive-Bayes Algorithm) →

Using the sklearn, we train different naive bayes algorithms, and then report precision, recall, f-score and accuracy of different values. Note that we have used the 80:20 split itself.

**For Gaussian Naive Bayes :**

	precision	recall	f1-score	support
BARBUNYA	0.66	0.52	0.59	246
BOMBAY	1.00	1.00	1.00	101
CALI	0.69	0.81	0.75	325
DERMASON	0.89	0.83	0.86	716
HOROZ	0.81	0.78	0.80	404
SEKER	0.69	0.72	0.71	406
SIRA	0.73	0.79	0.76	525
accuracy			0.77	2723
macro avg	0.78	0.78	0.78	2723
weighted avg	0.78	0.77	0.77	2723

**For Multinomial Naive Bayes:**

	precision	recall	f1-score	support
BARBUNYA	0.64	0.65	0.64	246
BOMBAY	1.00	1.00	1.00	101
CALI	0.77	0.77	0.77	325
DERMASON	0.88	0.83	0.85	716
HOROZ	0.85	0.79	0.82	404
SEKER	0.76	0.76	0.76	406
SIRA	0.72	0.81	0.76	525
accuracy			0.79	2723
macro avg	0.80	0.80	0.80	2723
weighted avg	0.80	0.79	0.79	2723

**For Bernoulli Naive Bayes :**

	precision	recall	f1-score	support
BARBUNYA	0.00	0.00	0.00	246
BOMBAY	0.00	0.00	0.00	101
CALI	0.00	0.00	0.00	325
DERMASON	0.26	1.00	0.42	716
HOROZ	0.00	0.00	0.00	404
SEKER	0.00	0.00	0.00	406
SIRA	0.00	0.00	0.00	525
accuracy			0.26	2723
macro avg	0.04	0.14	0.06	2723
weighted avg	0.07	0.26	0.11	2723

Since, the data is categorical and more skewed in nature. The gaussian naive bayes is able to give some good accuracy of 77 percent. Since, the data is multinomial (different classes of beans have been given), hence using a Bernoulli distribution is not able to give a good accuracy (26 percent) (barely above baseline random) The precision is only present for dermason due to its larger value. Multinomial naive bayes performs the best with 79 percent accuracy.

### Part E (Principal Component Analysis) →

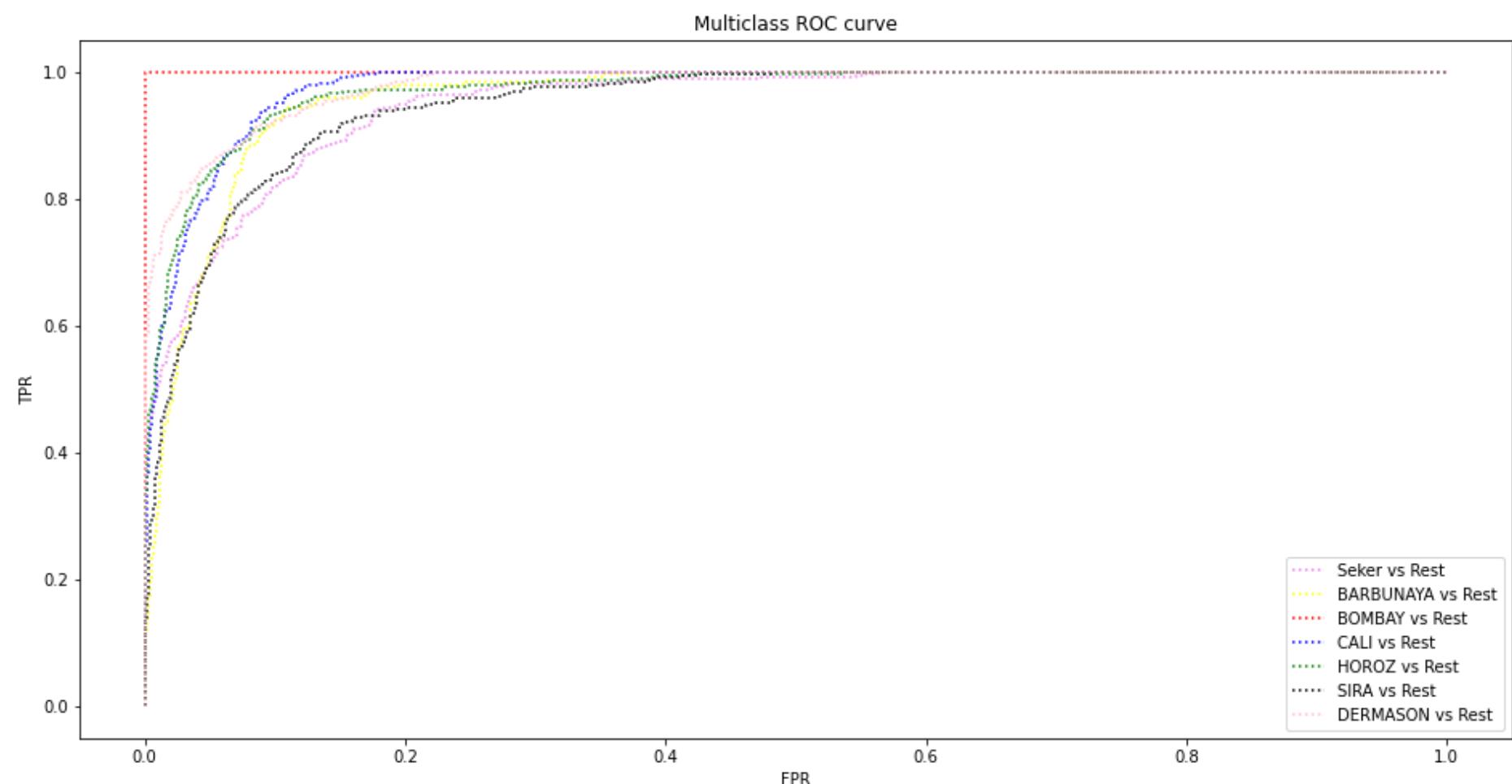
We can do the PCA for different components, 4,6,8,10,12. This is shown in this way : Using the gaussian naive bayes distribution, we will make the comparison.

S.No.	N Components	Accuracy Score (round off)	Precision Score	Recall Score	F1 Score
1	4	89.05	0.89	0.89	0.89
2	6	88.68	0.89	0.88	0.88
3	8	89.05	0.89	0.89	0.89
4	10	89.4	0.89	0.89	0.89
5	12	89.9	0.89	0.89	0.89

Accuracy has increased steadily as soon as the number of iterations are increased, with a drop and then an increase from 4 → 6 → 12 features. This will give the values required.

### Part F (ROC-AUC curve) →

Since, ROC\_AUC curves helps us figure the complete fitting of the classification model. We can say that the TPR values are higher than FPR values for the given scenario, that is why we get a ROC like that. (FPR → false positive rate) (TPR → true positive rates)



### Part G (Logistic Regression) →

The number of iterations used were = 10,000. This was so that the model gets enough iterations to reach the point of convergence. The alpha (learning rate is selected automatically). The algorithm was also throwing “failed to converge error on using lower number of iterations.” Increasing the iterations increases accuracy, till a single point. Increasing its enormously, does not return comparatively greater amount of accuracy.

The class chosen was multi-class due to a simple reason that the multinomial loss is minimised across the entire probability distribution. The newton\_cg solver gives a better accuracy with the l2 loss, function. l1 loss would have performed better, since the data has higher amount of outliers, but the newton\_cg covers up by giving better values in total. The accuracy comes out to be 92 percent. Which is far better than the naive bayes models trained on data on which PCA was used.

Model has greater accuracy than gaussian naive bayes, with 2

percent better accuracy. This can be credited to logistic regression and sigmoid function being used in it

	precision	recall	f1-score	support
BARBUNYA	0.95	0.91	0.93	246
BOMBAY	1.00	1.00	1.00	101
CALI	0.93	0.94	0.94	325
DERMASON	0.92	0.91	0.91	716
HOROZ	0.96	0.96	0.96	404
SEKER	0.94	0.93	0.94	406
SIRA	0.83	0.87	0.85	525
<b>accuracy</b>			0.92	2723
<b>macro avg</b>	0.93	0.93	0.93	2723
<b>weighted avg</b>	0.92	0.92	0.92	2723

S.No.	Penalty	Solver	multi_class	Accuracy
1	l2	newton_cg	multlinomial	92
2	l1	liblinear	—	90.1
3	l2	sag	—	67
4	l2	saga	—	62