

# Machine Learning Assignment - 1 report

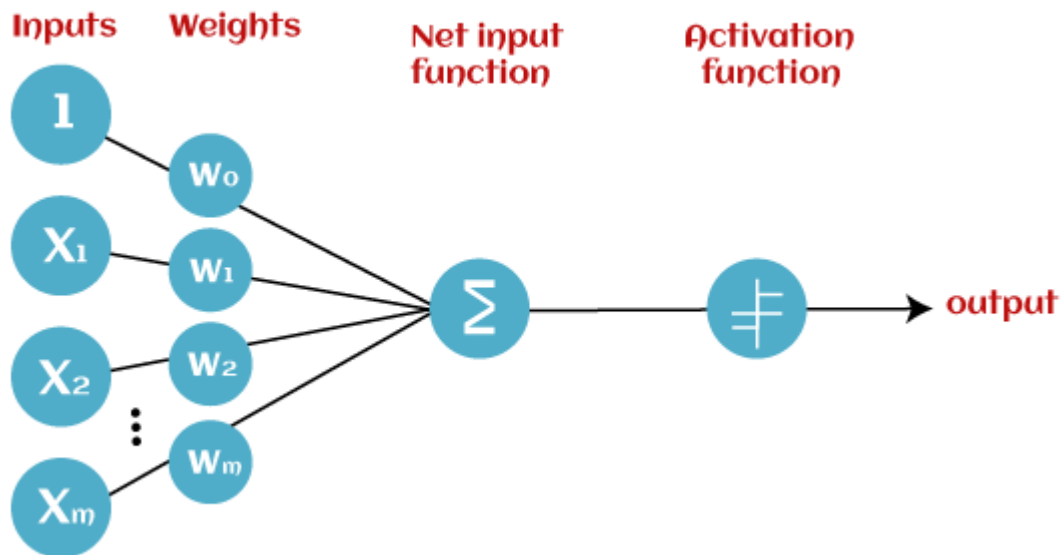
By:

Udit Gupta (ID:2020B4A72368H)

Aditya R Goyal (ID:2019B2A81443H)

Nimish Agrawal (ID: 2020B3A71857H)]

## Perceptron Algorithm:



The perceptron algorithm executes with the help of following steps:

**Initialize the weights:** The perceptron weights are initialized with small random values.

**Present an input:** The perceptron is presented with an input vector  $x$ , which is multiplied by the weights to produce the output  $y$ .

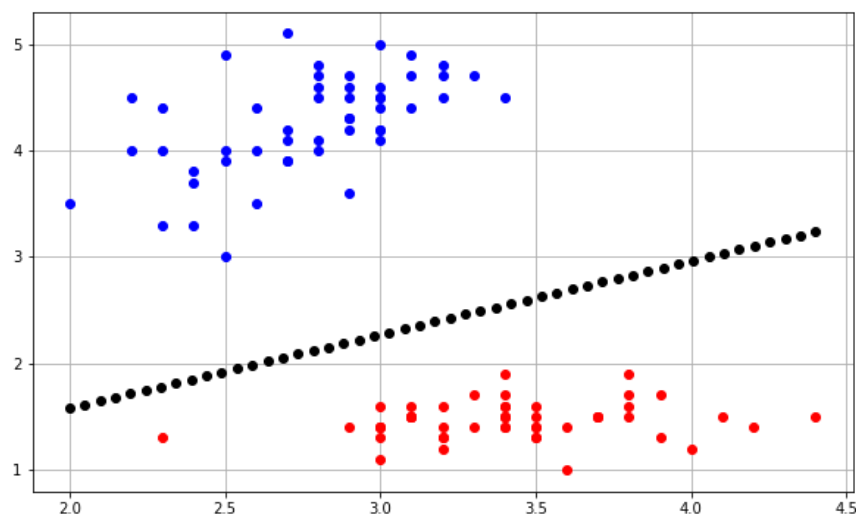
**Compute the error:** The error is the difference between the actual output  $y$  and the desired output  $d$ .

**Update the weights:** The weights are updated using the formula  $w = w + \alpha(d - y)x$ , where  $w$  is the weight vector,  $\alpha$  is the learning rate,  $d$  is the desired output,  $y$  is the actual output, and  $x$  is the input vector.

**Repeat steps 2-4:** Steps 2-4 are repeated for each input-output pair in the training data until the weights converge and the network can accurately predict the outputs for new inputs. The learning rate is an important hyperparameter in the Perceptron algorithm. It determines how much the weights are adjusted during each iteration. A high learning rate can cause the weights to oscillate and prevent convergence, while a low learning rate can cause the weights to converge slowly.

The Perceptron algorithm is used for binary classification tasks, where the goal is to separate the input data into two classes. The decision boundary of the perceptron is a linear function that separates the input data into two regions. If the input data is not linearly separable, the Perceptron algorithm may not converge and may produce inaccurate results. In such cases, more complex neural networks such as multi-layer perceptron's or convolutional neural networks are used.

The Perceptron Algorithm was used to test whether the data corresponding to the two types of tumors - benign and malignant - are linearly separable.



## Implementation:

The class perceptron has been defined and some instance variables have been initialized, including weights, bias, learning rate, and epochs (iteration).

The weights and bias have been initialized with 0 and get updated with each subsequent iteration according to the following mathematical equation:

$$w_i = w_0 + \eta * y_i * x_i$$

$$b_i = b_0 + \eta * y_i$$

where  $\eta$  denotes the learning rate.

we have assumed  $\eta=0.01$

A predict function has also been used governed by the following equation:

$$predict(x_{test}) = activation(w_i \cdot x_i + b_i)$$

### Pseudo Code:

**input:** A training set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$   
**initialize:**  $\mathbf{w}^{(1)} = (0, \dots, 0)$   
**for**  $t = 1, 2, \dots$   
    **if**  $(\exists i \text{ s.t. } y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle \leq 0)$  **then**  
         $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + y_i \mathbf{x}_i$   
    **else**  
        **output**  $\mathbf{w}^{(t)}$

### Inferences:

The given dataset is Not linearly separable.

One possible way for checking whether the data is linearly separable could have been to plot the various features taken two at a time and check.

PM1: Perceptron Model 1 is the first model we build using the perceptron algorithm on our dataset. It is trained on the original dataset, assuming that the features are in their original order. The algorithm updates the weights based on the order of the training examples until convergence, and the resulting model is used to classify new examples.

**PM1 accuracy: 0.9202127659574468**

PM2: Perceptron Model 2 is like PM1, but the order of the training examples is shuffled randomly before training. This can lead to different weight updates and a different final model compared to PM1.

**PM2 accuracy: 0.8936170212765957**

PM3: Perceptron Model 3 is a new model that we build using the perceptron algorithm, but the dataset is first normalized this time. Normalization involves scaling the features to have zero mean and unit variance. This can help the algorithm converge faster and result in a better model.

**PM3 accuracy on normalized data: 0.973404255319149**

PM4: Perceptron Model 4 is trained on the same dataset as PM1, but this time the order of the features in each example is randomly permuted. This means that the same features will have different weights than in PM1, which can lead to a different final model.

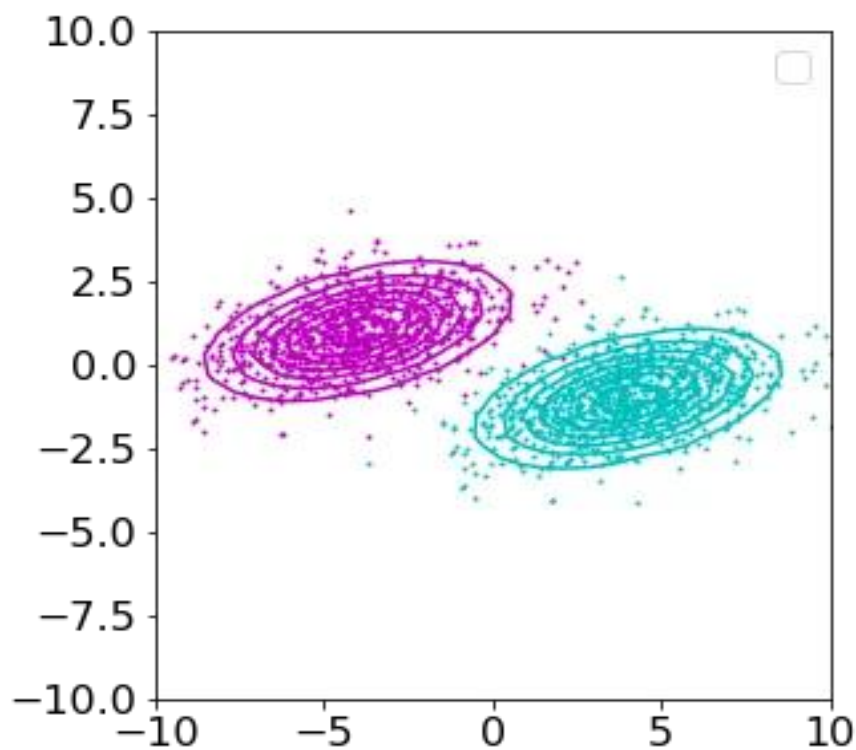
**PM4 Accuracy: 0.9202127659574468**

In summary, all four models are trained using the same perceptron algorithm, but they differ in their input data (PM3 is normalized), input order (PM2 shuffles the training examples), and input feature order (PM4 permutes the feature order). These differences can lead to different weight updates and resulting models, and we can compare the performance of each model on a test set to evaluate their effectiveness.

As we can observe PM1 is approximately same as PM4 as the order of features is changed randomly, it does not affect the linear separability of the data. The Perceptron algorithm only depends on the order of the training examples, not on the order of the features. Therefore, the model PM4 would be very similar to PM1.

## Fischer LDA:

Fisher's linear discriminant can be used to make a classifier that learns with help. If the classifier is given data that has been labeled, it can find a set of weights that can be used to draw a decision limit and sort the data. Fisher's linear discriminant tries to find the vector that separates the classes of projected data as much as possible. "Separation" can be hard to define. Fisher's linear discriminant does this by making sure that the gap between the projected means is as big as possible and that the projected within-class variance is as small as possible



## Fisher's criterion

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

Where  $m_1 - m_2 = \mathbf{w}^T(\mathbf{m}_1 - \mathbf{m}_2)$  and  $s_k = \sqrt{\sum (\mathbf{w}^T \mathbf{x}_n - m_k)^2}$

Find the weights vector that maximizes the above equation to maximize Fisher's criterion. Getting the most out of this math is the same as what we saw before. Maximizing means making the numerator (the distance between predicted means) as big as possible and to make the denominator as small as possible. (The within-class variance). We can rewrite Fisher's criterion by putting the formulae for the means and covariance in place of the means and covariance

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}}$$

Where  $S_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$

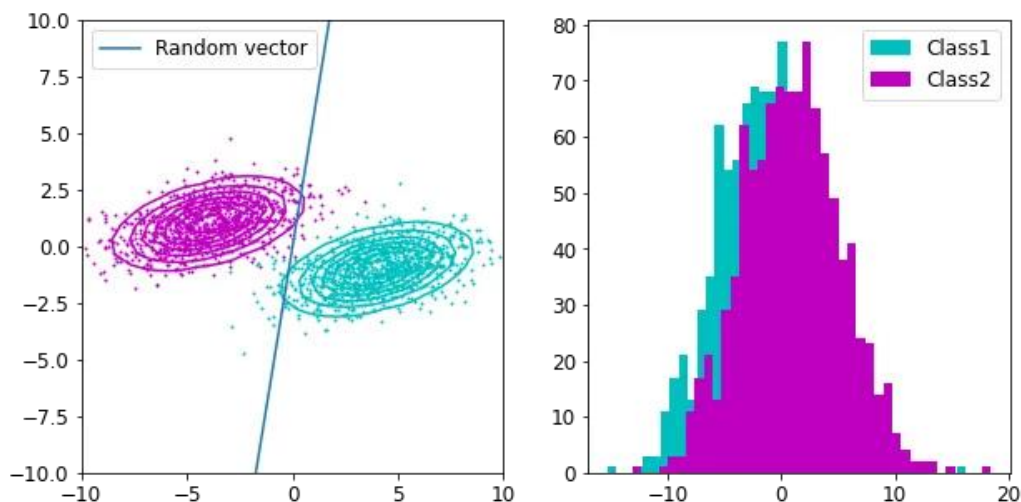
And  $S_w = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^T$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{(\mathbf{w}^T S_w \mathbf{w})(2 S_B \mathbf{w}) - (\mathbf{w}^T S_B \mathbf{w})(2 S_w \mathbf{w})}{(\mathbf{w}^T S_w \mathbf{w})^2} = 0$$

$$\Leftrightarrow S_w \mathbf{w} = \alpha S_B \mathbf{w}$$

$$\Leftrightarrow S_w \mathbf{w} = \alpha (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w}$$

$$\Leftrightarrow \mathbf{w} \propto S_w^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

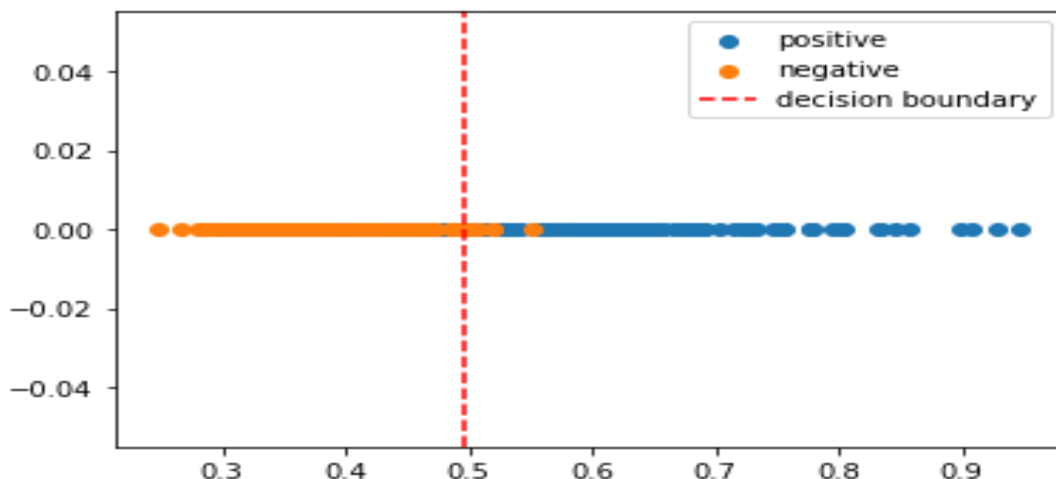


## Methodology

### FDLA -1

1. First, we imported the dataset and checked for any Nans and filled all such cells with 0s.
2. The diagnosis column of the dataset was *one-hot encoded*, and columns of diagnosis and ID were dropped.
3. The dataset which would become the training set was normalized and then the one-hot encoded diagnosis-2 was added to it.
4. We then wrote functions for the following:
  - a. Train-test splitting
  - b. Divided the X and y into training and testing halves.
  - c. Then we segregated the training data into positive class and negative class
  - d. Calculated Sw\_0 and Sw\_1
  - e. Calculated mean and variance and mean difference.
  - f. Solved the equations for roots.
  - g. Plotted the roots of the equation.
  - h. The actual threshold value is the root between two means.
  - i. That x is now our threshold value.
  - j. Values lying left of the threshold values will be classified as Malignant and values lying toward the right will be classified as Benign.

## Results



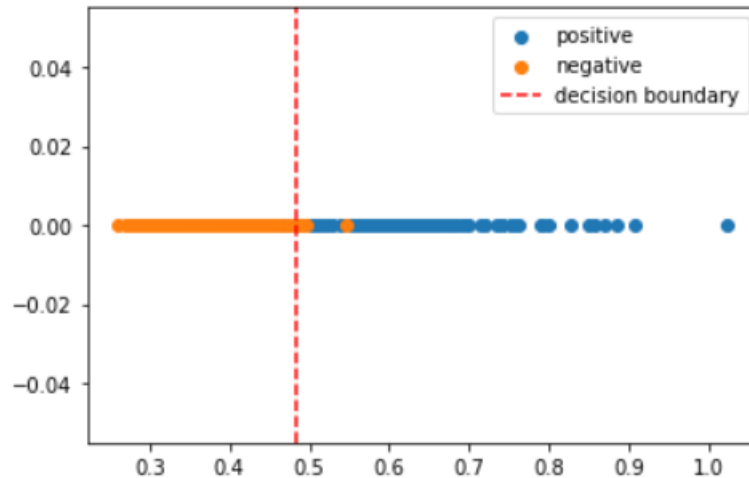
Threshold value(th\_val) = 0.4955243461707635

the threshold value we have obtained, th\_val is our separating hyperplane. for  $x < \text{th\_val}$ , the value is classified as Malignant and for  $x > \text{th\_val}$  it's classified as Benign.

Testing accuracy obtained = 93.617

## FDLA -2

Fisher's Linear Discriminant model (FLDM2) on the shuffled training data and found the decision boundary in the univariate dimension using a generative approach. We can now compare the performance of FLDM1 and FLDM2 on the testing data.



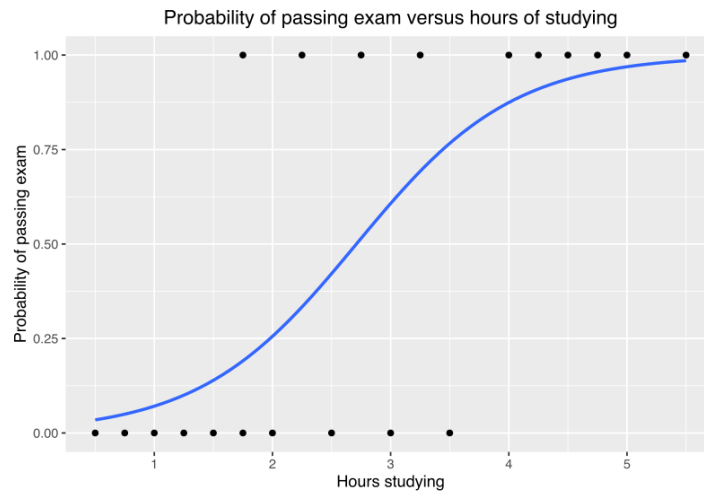
Threshold value(th\_val) = 0.48334336

Testing accuracy obtained = 0.93617021



# Logistic Regression:

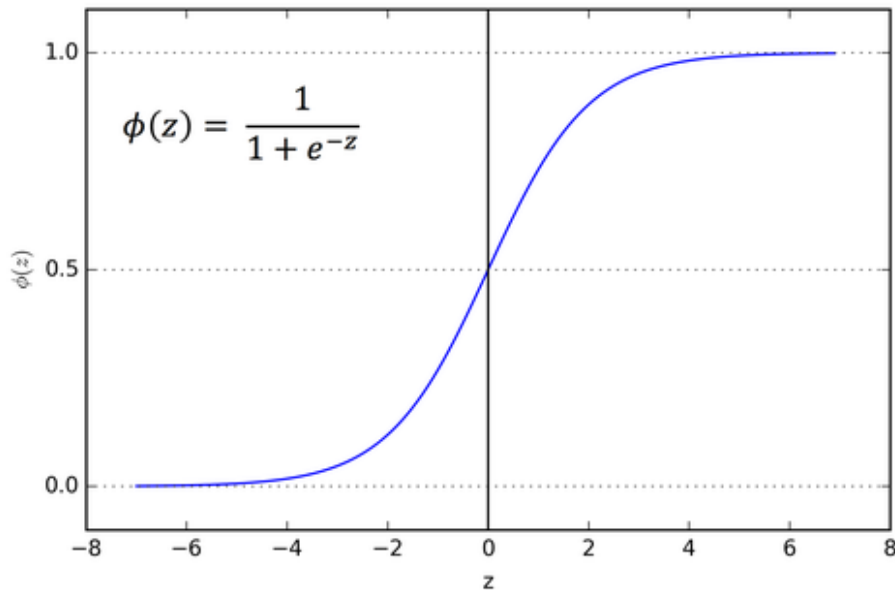
It is a classification technique used when the target variable is *categorical* (example 0 or 1).



## Sigmoid Function:

The hypothesis is  $\Rightarrow Z = Wx + B$  and the output is supposed to be one of 0 or 1

The probability is defined by  $\Rightarrow$  Sigmoid function:



If 'Z' goes to infinity, Y(predicted) will become 1 and if 'Z' goes to negative infinity, Y(predicted) will become 0.

The hypothesis' output is considered as the estimated probability. This is used to infer how confident can predicted value be actual value when given an input X.

### **In mathematical format:**

$h(x) = P(Y=1 | x) \Rightarrow$  probability of Y given x and  $P(Y=0 | x) = 1 - P(Y=1 | x)$

Data is fit into linear regression model, which then be acted upon by a logistic function predicting the target categorical dependent variable.

### **Decision boundary:**

To predict which class a data belongs, a **threshold** can be set. Based upon this threshold, the obtained estimated probability is classified into classes.

If  $(y_{\text{pred}} \geq 0.5) \Rightarrow$  then classify 1 else 0.

Decision boundary can be linear or non-linear. Polynomial order can be increased to get complex decision boundary.

### **Cost Function:**

The cost function is decided by the NLL (or the negative log likelihood).

$$\begin{aligned} \text{Cost}(h_{\theta}(x), Y(\text{actual})) &= -\log(h_{\theta}(x)) \text{ if } y=1 \\ &\quad -\log(1 - h_{\theta}(x)) \text{ if } y=0 \end{aligned}$$

To converge this **gradient descent** is used  $\Rightarrow$  which will only converge if the function is a convex one!

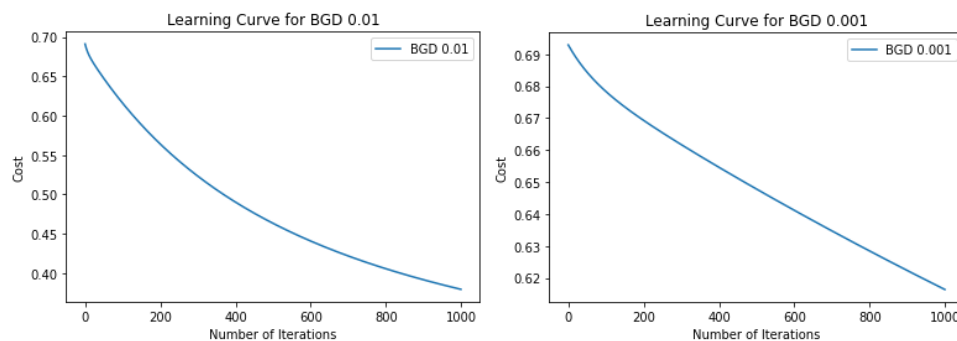
Negative function is because when we train, we need to maximize the probability by minimizing loss function. Decreasing the cost will increase the maximum likelihood assuming that samples are drawn from an identically independent distribution.

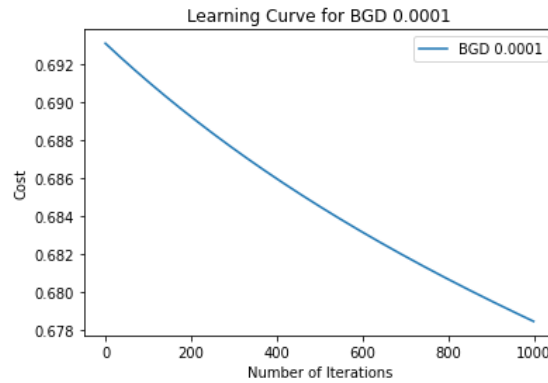
## LR-1:

### Methodology:

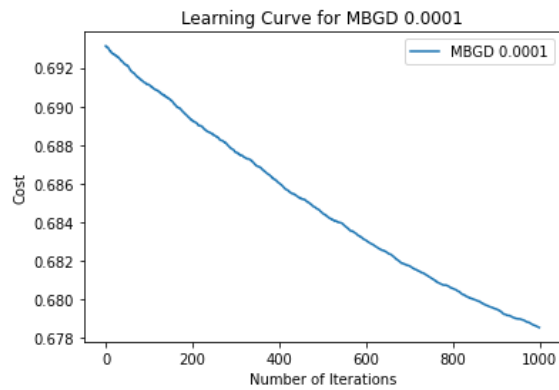
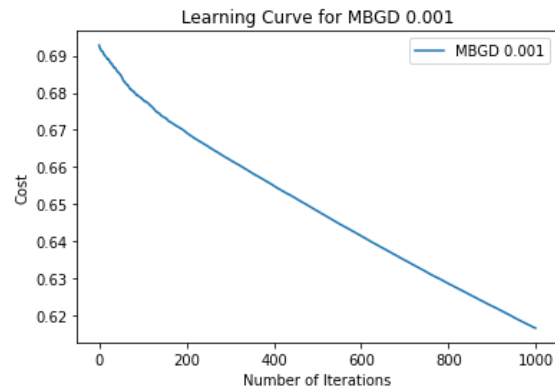
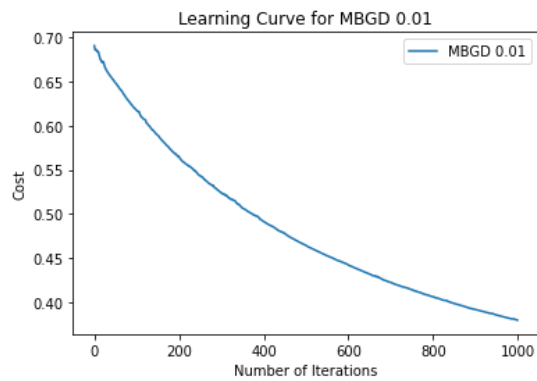
1. First, we imported the dataset and checked for any NaNs and filled all such cells with 0s.
2. Diagnosis column of the dataset was *one-hot encoded*, and columns of diagnosis and ID were dropped.
3. The dataset which would become the training set was normalized and then the one-hot encoded diagnosis-2 was added to it.
4. We then wrote functions for the following:
5. Train-test splitting
6. Divided the X and y into training and testing halves.
7. Then we converted them to numpy arrays
8. Lastly, we wrote a class for the logistic regression model which contains the following functions:
  - a. Initialization
  - b. Sigmoid
  - c. Compute\_cost
  - d. Batch\_gradient\_descent
  - e. Stochastic\_gradient\_descent
  - f. Mini\_batch\_gradient\_descent
  - g. Predict
  - h. Plot\_cost

**Plots for Learning Rates =0.01,0.001,0.0001 for batch gradient descent.**

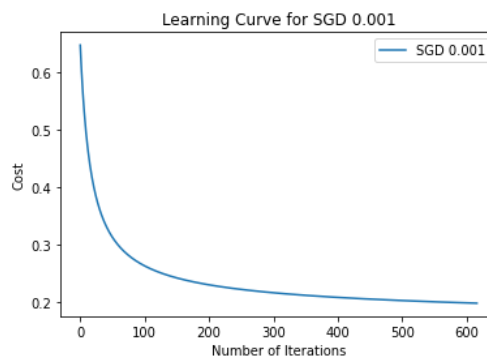
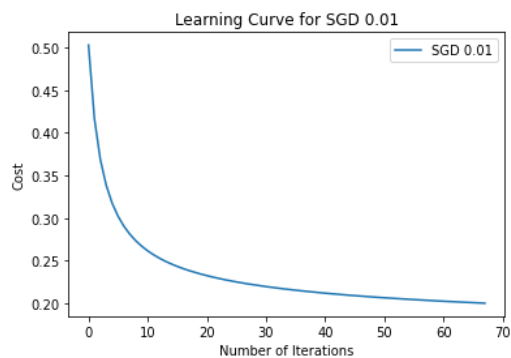


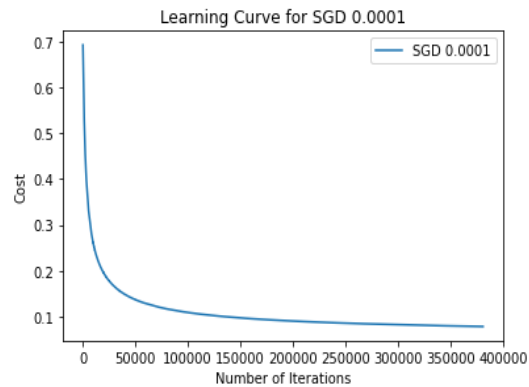


**Plots for Learning Rates =0.01,0.001,0.0001 for Mini batch gradient descent:**



**Plots for Learning Rates =0.01,0.001,0.0001 for Stochastic gradient descent :**



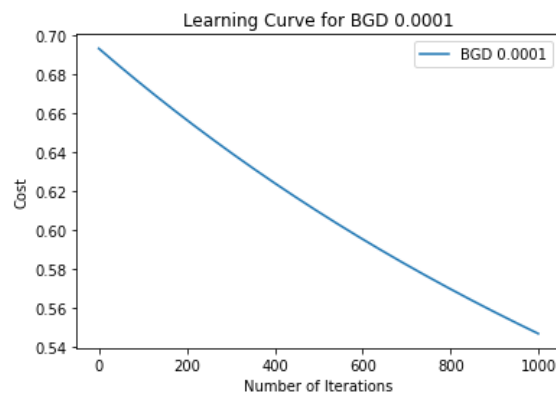
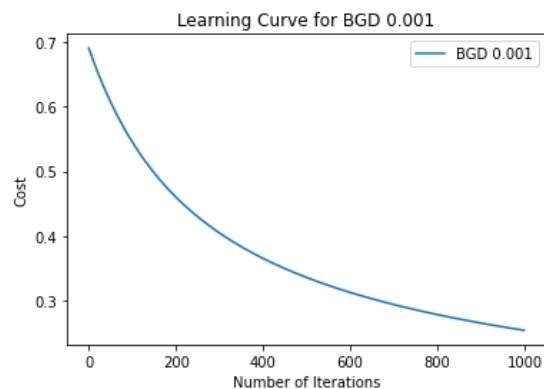
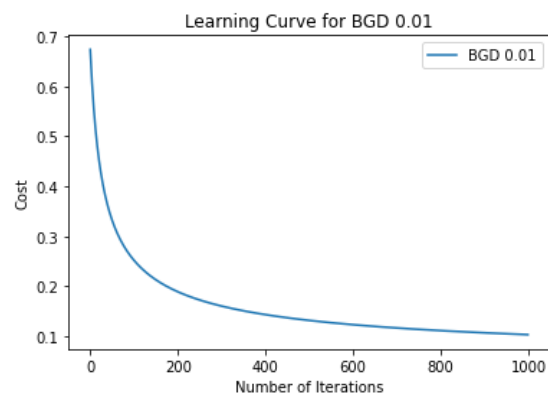


## LR2:

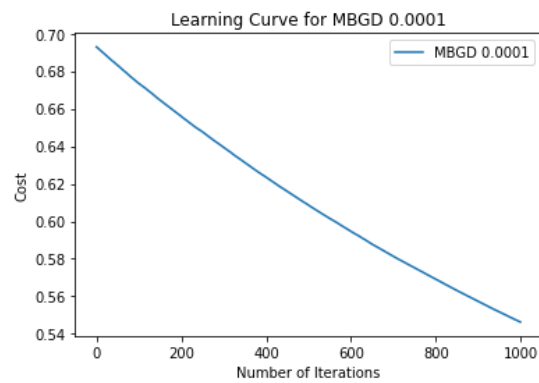
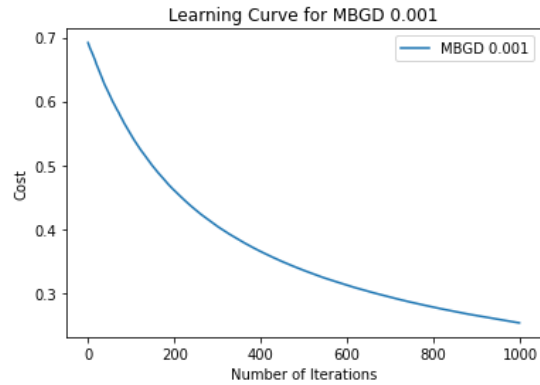
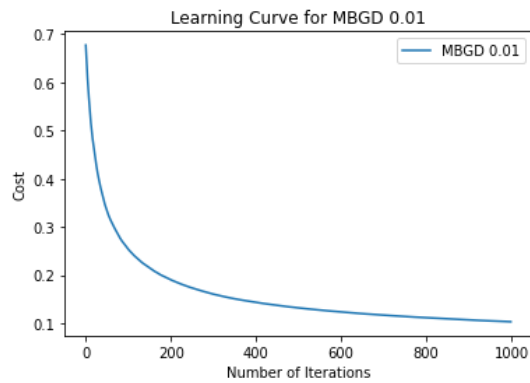
### Methodology

Same procedure as LR1 was followed, instead all the empty values were replaced with their means for continuous data and mode for discrete and X set was normalized.

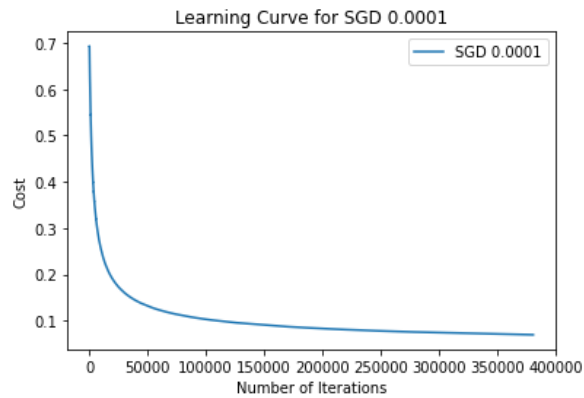
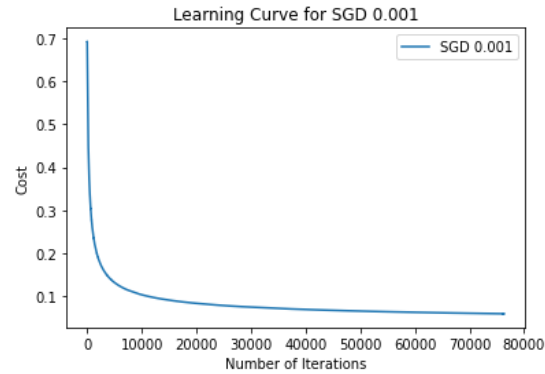
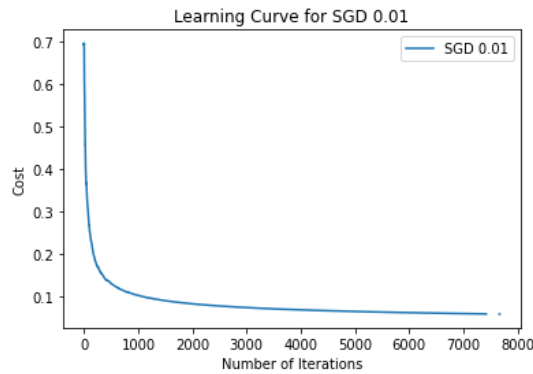
### Plots for Learning Rates =0.01,0.001,0.0001 for Batch gradient descent:



### Plots for Learning Rates =0.01,0.001,0.0001 for Mini Batch gradient descent:



**Plots for Learning Rates =0.01,0.001,0.0001 for Stochastic gradient descent :**



## COMPARATIVE STUDY:

Sr.no	Random state value	Accuracy (FLDA-1 )	Accuracy (FLDA-2 )	Accuracy (LR-1)	Accuracy (Logistic-LR2-> stochastic gradient descent)	Accuracy (PM-1)	Accuracy (PM-3)	Accuracy (PM-4)
1	1	0.952127 65957446 81	0.952127 65957446 81	0.888297 87234042 56	0.984042 55319148 94	0.9202127 65957446 8	0.941489 36170212 77	0.920212 76595744 68
2	2	0.936170 21276595 75	0.936170 21276595 75	0.888297 87234042 56	0.973404 25531914 9	0.9202127 65957446 8	0.973404 25531914 9	0.920212 76595744 68
3	3	0.946808 51063829 79	0.946808 51063829 79	0.813829 78723404 25	0.973404 25531914 9	0.8776595 74468085 1	0.952127 65957446 81	0.877659 57446808 51
4	4	0.936170 21276595 75	0.936170 21276595 75	0.904255 31914893 62	0.968085 10638297 87	0.8563829 78723404 3	0.984042 55319148 94	0.856382 97872340 43
5	5	0.914893 61702127 66	0.914893 61702127 66	0.877659 57446808 51	0.962765 95744680 85	0.8351063 82978723 4	0.984042 55319148 94	0.835106 38297872 34
6	6	0.904255 31914893 62	0.904255 31914893 62	0.877659 57446808 51	0.936170 21276595 75	0.8882978 72340425 6	0.960851 06382978 7	0.888297 87234042 56
7	7	0.941489 36170212 77	0.941489 36170212 77	0.867021 27659574 47	0.962765 95744680 85	0.8776595 74468085 1	0.965085 10638297 87	0.877659 57446808 51
8	8	0.936170 21276595 75	0.936170 21276595 75	0.893617 02127659 57	0.962765 95744680 85	0.9202127 65957446 8	0.968081 10638297 87	0.920212 76595744 68
9	9	0.941489 36170212 77	0.941489 36170212 77	0.893333 33333333 33	0.973404 25531914 9	0.8829787 23404255 3	0.952127 65957446 81	0.882978 72340425 53
10	10	0.946808 51063829 79	0.946808 51063829 79	0.882978 72340425 53	0.962765 95744680 85	0.5638297 87234042 5	0.957446 80851063 83	0.563829 78723404 25
Mean		0.935638 30	0.935638 30	0.878695 035	0.971808 5106	0.854255 32	0.964893 62	0.854255 32
SD		0.014163 33	0.014163 33	0.0250310 3	0.0125874 04	0.100423 50	0.013287 23	0.100423 50

As we can see **Logistic-stochastic gradient descent** give us best result among all the algorithm since it is efficient in handling multiple feature