

INTRODUCTION

- **Supervised Learning:**

- **Regression Problem:** When we have a continuos value in our output then it knows as regression problem.
- **Classification Problem :** There are two types of classification problems 1) **Binary classification** and 2) **Multiclassification** problem

- **Unsupervised learning:**

- **Clustering:** In Clustering, we don't have output feature or dependent variable, it's grouping algorithm ,**not classification problem**. We make clusters that describes info. For example, based on the data apple have customer who afford expensive products and cheap products. We can make two clusters, one who purchased expensive one consider as rich group of people and other one consider as middle class people. when they launch new products as per the price , it shows or give recommendation to the particular customers. Let's say apple launch i watch series 9 and i watch SE ,So based on the customer previous orders it suggest them either buy Series 9 or SE model.

1. Linear Regression

Equation of Linear regression is

$h(x) = \theta_0 + \theta_1 x$ where θ_1 is slop or coefficient, and θ_0 is intercept which means at what point this line meet at y axis.

- The main aim of Llinear regression is to create best fit line so that the distance between predicted values and actual value would be minimal as much as we can. We use Cost function to find out distance between predicted and actual values,

Cost Function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h(X_i) - Y_i)^2$$

Where, $h(X_i)$ is predicted values and Y_i is actual values, and $1/m$ is use to get average cost function , $1/2$ will help us in derivative purpose.

So, this is know as **squared error function**

What we need to solve

We need to minimize the above cost function by adjusting the θ_0 and θ_1 .

Gradient Descent:

Gradient descent will play an important role to ensure that we get right slop value. In Gradient descent , we get global minima and that is the best fit line for linear regression where the distance between actual values and predicted values are minimal.

So, How do we find all points for gradient descent and get global minima? Well, for that we will use convergence algorithm

What convergence algorithm says that,

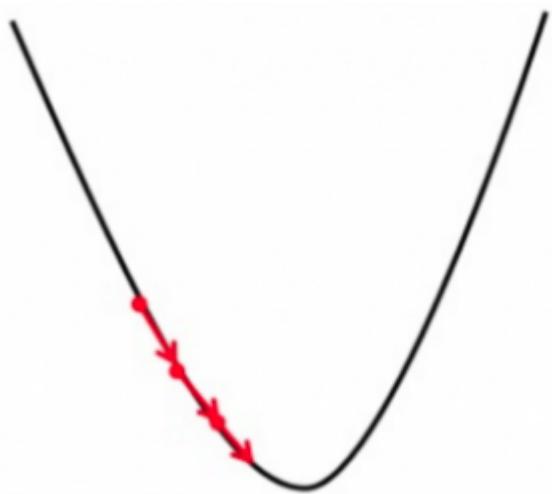
Repeat until convergence

{

$$\theta_j := \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta_0, \theta_1)$$

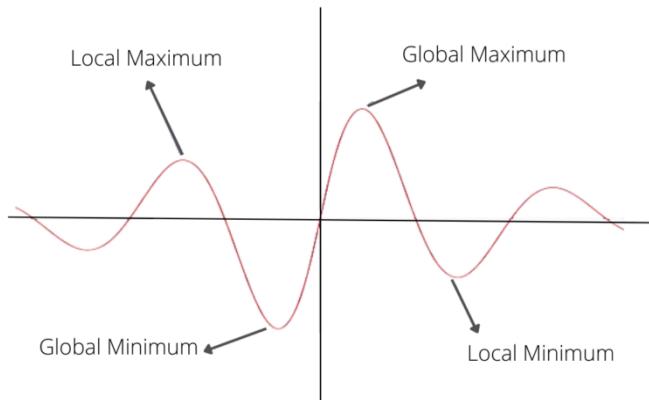
}

Here, α is Learning Rate , Alpha here play major role as if we take alpha value small like $\alpha = 0.01$, it takes small steps to move towards the optimal minima not even too small value , while if we take alpha value huge like $\alpha = 1$, the updation of θ_1 value will jump and never reach the global minima.



Learning rate when we use small alpha value

Now, what if we have local minima in our Cost function



At local minima our Equation would be liked,

$$\theta_1 := \theta_1 - \alpha(0), \text{ as } \text{slop}(\frac{\delta}{\delta \theta_j} J(\theta_0, \theta_1)) \text{ will be 0 at local minima,}$$

So, $\theta_1 := \theta_1$.

However, with the respect to the Gradient Descent and equation that we are using ,we don't get stuck in local minima.

Note:

What we do if we get local minima in Linear Regression, so Answer is that the Cost function we use ,will definitely not give us a local minima.

★ Convergence Algorithm.

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2$$

Here, convergence theorem

If $J=0$

$$\Rightarrow \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_0} \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)$$

$$= \frac{\partial}{\partial \theta_0} \sum_{i=1}^m (h(x_i) - y_i)$$

because, $h(x_i) = \theta_0 + \theta_1 x_i$, so

$$\frac{\partial}{\partial \theta_0} h(x_i) = 1$$

$$J=0 \Rightarrow \boxed{\frac{\partial}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)}$$

★ When $J=1$

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_1} \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i) x_i$$

$$\therefore \frac{\partial}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i) \cdot x_i$$

As, $h(x_i) = \theta_0 + \theta_1 x_i$

$$\frac{\partial}{\partial \theta_1} h(x_i) = x_i$$

$$\boxed{J=1 \Rightarrow \frac{\partial}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i) \cdot x_i}$$

Convergence Algorithm

So,

Repeat until Convergence

{

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h(X_i) - Y_i)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h(X_i) - Y_i) \cdot (X_i)$$

}

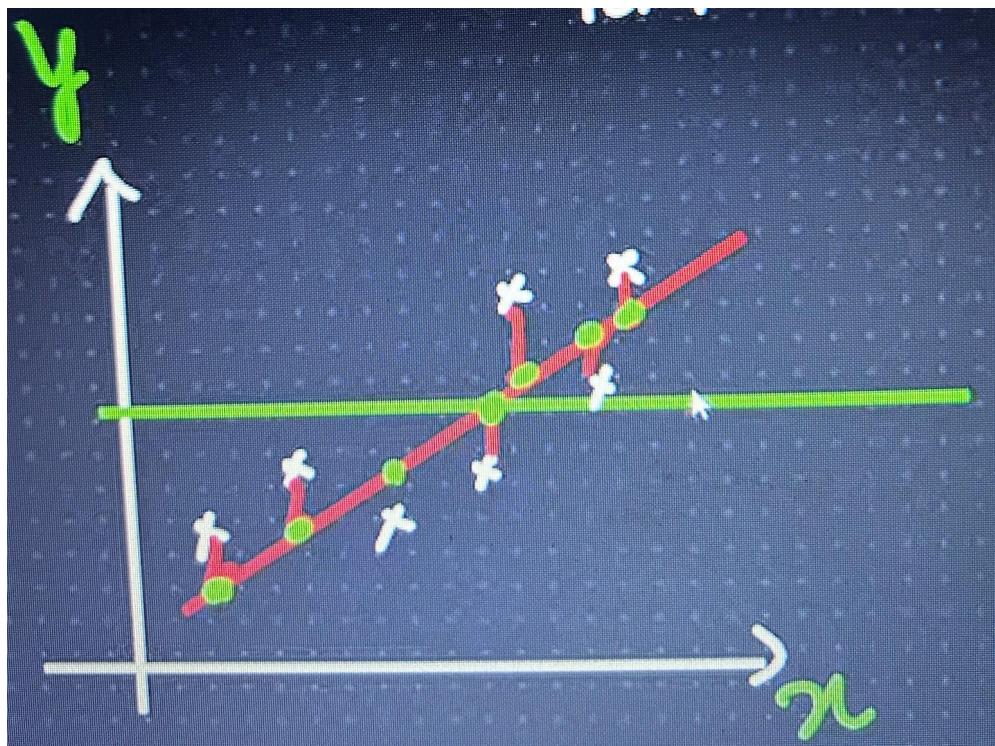
Performance Metrics:

We use two metrics to check the performance of model 1) **R²** and 2) **Adjusted R²**

R²:

$$R^2 = 1 - \frac{SS_{residual}}{SS_{total}} = 1 - \frac{\sum(Y_i - Y_i^-)^2}{\sum(Y_i - Y^-)^2}$$

Here, Y_i^- is $h(x)$, and Y^- is mean of y



$\sum(Y_i - Y^-)^2$ will have higher values than $\sum(Y_i - Y_i^-)^2$, So $1 - \frac{\text{Low value}}{\text{High value}}$ after calculation of $\frac{\text{Low value}}{\text{High value}}$, we get smaller value then we minus it from 1 and eventually we get higher number. In conclusion, we can say that our model work perfectly and have good R².

Let's say we have features that are correlated than in that scenario our R^2 will increase. Even though we have features that are not correlated , R^2 will also go up. When we have features that highly correlated our R^2 will be increased tremendously and have less correlated features our R^2 will be small increased.

For example, let's say we have features like bedroom, location that are correlated with price so here R^2 increases, and when we add another feature like gender which is not correlated at all although R^2 will increase but we should not have been picked this model, instead we consider the first model where all features are correlated. So, in order to prevent this situation we use Adjusted R^2.

R^2 Adjusted:

$$R^2 \text{ Adjusted} = 1 - \frac{(1-R^2)(N-1)}{N-P-1}, \text{ where , N is total number of samples and } P \text{ is Features or Predictors}$$

As P value is getting increased, our R^2 adjusted will be decreased to adjust R^2.

Ridge and Lasso Regression:

Overfitting and Underfitting:

When our model works well on **TRAINING DATA (also known as low bias)** and don't work well on **TEST DATA (High Variance)** that is called **Overfitting**, Whereas **Underfitting** means when not only on **Training data (High bias)**, but also on **testing data (High variance)** our model doesn't work well.

Ridge Regression (L2 Regularization):

Whenever we get overfitting in Linear Regression, we use L2 regularization in order to prevent Overfitting.

Let's consider here, $J(\theta) = 0$

$$\text{So, } J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h\theta(x_i) - y_i)^2$$

With the help of this equation if we get $J(\theta_0, \theta_1) = 0$, then we need to use L2 regularization, and the equation of L2 is

$$(h\theta(x_i) - y_i)^2 + \lambda(slop)^2$$

Here, lambda is a hyperparameter, we will find lambda through cross validation

Lasso Regression (L1 Regularization):

L1 Regularization is also used to prevent **Overfitting**, here the Equation of L1 is

$$(h\theta(x_i) - y_i)^2 + \lambda|slop|$$

Apart from Overfitting, it is also useful for feature selection, let's take example here,

$$h\theta(x) = y^- = \theta_0 + \theta_1x_1 + \theta_2x_2 + \theta_3x_3 + \theta_4x_4 + \dots + \theta_nx_n$$

So, now we take a mode of all slopes like this,

$$|\theta_0 + \theta_1 + \theta_2 + \theta_3 + \theta_4 + \dots + \theta_n|$$

As we go ahead whichever features do not play amazing role its slope value is small, here L1 will neglect that features.

In conclusion, we do two major things in Lasso regression 1) **To prevent Overfitting, and 2) Assist to select an important features and neglect the less important the less important features**, While we use Ridge regression **for only preventing Overfitting**.

Assumption of Linear Regression:

- 1) If our **features are normal or gaussian distribution**, Our model will be **trained well**. However, if our features does not follow normal or gaussian distribution, we will apply kind of mathematical equation in order to convert them into either normal or gaussian distribution.
- 2) **Standardization** is nothing but Scaling data with the help of **Z-score** in order to make **mean(μ) value 0** and **Standard deviation(σ) value 1**. If we want to **increase training time of model or optimize the model** it is inevitable to use standardization.
- 3) Linear Regression works well with respect to linearity.
- 4) We need to check **Multi-collinearity**.

Let's say we have three features **X1,X2,X3** and output feature **Y**, now let's say **X1 and X2 are 95% correlated**, also these are **highly correlated with output feature Y**, then it's not necessary to use both features,so we can drop one of two features.

Hence, **Multicollinearity** is also solved with the help of **Variation inflation factor**.

2. Logistic Regression (Classification)

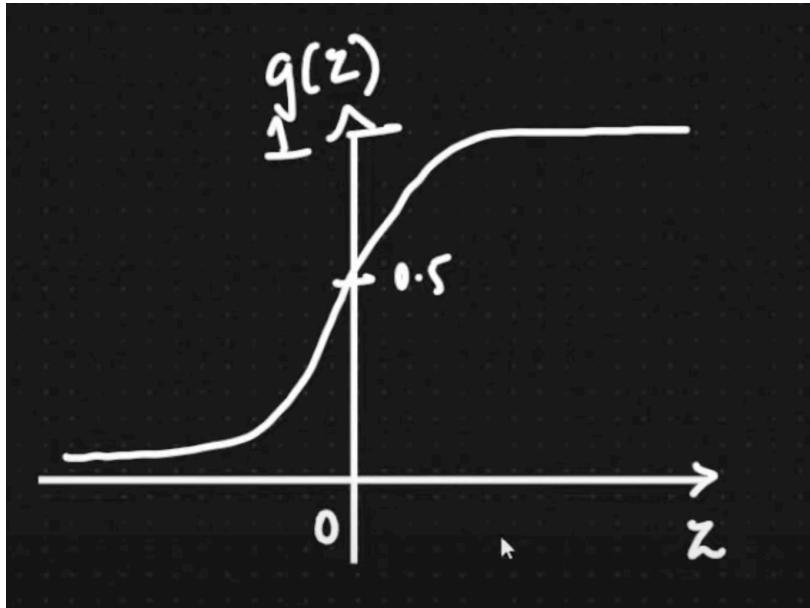
- This Regression is helpful when we have binary classification like yes or no , 0 or 1 , True or False.
- In order to solve this type of problem we can use **Linear Regression**, however, we might face problems when we get outliers. Our maxima and minima values are 1 and 0 respectively, and sometimes due to outliers we may get value grater than 1 or less than 0 in that case we cann't use Linear Regression.
- We apply mathematical formula on top of the linear regression to make a **Squash line at the beginning and in the end**.

$h\theta(x) = \theta_0 + \theta_1 X$, This is normal linear regression function in order to make logistic function we simply add **another mathematical function to make squash line**, but first let's say $Z = \theta_0 + \theta_1 X$

So, $h\theta(x) = g(z)$,

$$h\theta(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 X)}}$$

So, this is **Sigmoid** or **Logistic** function.



We can make assumption that is , $g(z) \geq 0.5$, When we have $z \geq 0$

If we use above $h_\theta(x)$ equation in linear regression's cost function then we may have to face non-convex function where we get lots of local minima So, there is a different way to handle this problem,

Logistic Regression Cost function

$$J(\theta_0) = \begin{cases} -\log(h_\theta(x)) & y=1 \\ -\log(1-h_\theta(x)) & y=0 \end{cases}$$

$$h_\theta(x) = \frac{1}{1+e^{-(\theta_0 + \theta_1 x)}}$$

$$J(\theta_0) = -\frac{1}{2m} \sum_{i=1}^m (y^i \log(h_\theta(x^i)) + (1-y^i) \log(1-h_\theta(x^i)))$$

$$h_\theta(x^i) = \frac{1}{1+e^{-\theta_0 x^i}}$$

Repet until convergence

$$\left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \\ \end{array} \right.$$



$J(\theta)$ is the cost function, we use convergence algorithm till we get global minima as usual.

Now, how to calculate accuracy with the help of confusion matrix

Untitled (7/7) - Scribble link

Actual

| | | | | |
|-------|-------|-----|-----------|----------------------|
| x_1 | x_2 | y | \hat{y} | \rightarrow Output |
| - | - | 0 | 1 | |
| - | - | 1 | 1 | |
| - | - | 0 | 0 | |
| - | - | 1 | 1 | |
| - | - | 1 | 1 | |
| - | - | 0 | 1 | |
| - | - | 1 | 0 | |

Pred

| | |
|---|---|
| 1 | 2 |
| 3 | 1 |

Predicted

| | | |
|----|----|------------------|
| 1 | 0 | Actual |
| TP | FP | Confusion matrix |
| FN | TN | |

Accuracy = $\frac{TP + TN}{TP + FP + FN + TN}$

$$= \frac{3+1}{3+2+1+1} = \frac{4}{7}$$

$$= 0.57 = \underline{\underline{57\%}}$$

A video frame showing a man with a beard and headphones, likely the lecturer, speaking.

However, we can't always depend on accuracy, for example if we have same 900 data out of 1000 so definitely, we get 90% accuracy. To handle this issue , we use

- 1) Precision
- 2) Recall
- 3) F-score

Precision

From all the **predicted** positive values how many are actually **true** or **Positive**.

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

For example, in spam mail identification, it should be identify as spam in that scenario **False positive** should be reduced and **FP** talks in a better way for spam classification, So here our priority is **FP**

Recall

Definition of recall is that from all actual positive values how many predicted true values, so the formula of recall is

Recall = $\frac{\text{TP}}{\text{TP} + \text{FN}}$, Our target in recall that is we should reduce **False Negative (FN)**

For example, if a person has a cancer model , value should be 1 otherwise it will be being trouble if it shows 0 ,So here our priority is **FN**.

F-score

We use F-score when positive and negative both are an important. For example, Stock market crash tomorrow, here , from people point of view it's positive as they can sell there stocks, while as per the industry it's negative view. In this type scenario ,we use **F-score**.

So, It's necessary to give focus on both **False Positive (FP)** and **False Negative (FN)**.

$$\text{F-score} = \frac{(1 + \beta^2) (\text{Precision} * \text{Recall})}{\beta^2 * (\text{Precision} + \text{Recall})}$$

Case - 1

If **False Positive (FP)** and **False Negative (FN)** are an important ,

$$\beta = 1 , \text{F1-score} = \frac{(1 + 1) (\text{Precision} * \text{Recall})}{(\text{Precision} + \text{Recall})}$$

Case - 2

If False Positive is more important than False negative, we use

$$B = 0.5, F0.5 - \text{score} = (1 + 0.5^2) (\text{Precision} * \text{Recall}) / 0.5^2 * (\text{Precision} + \text{Recall})$$

Case - 3

If False negative is more important than False positive, we use

$$B = 2, F2 - \text{score} = (1 + 2^2) (\text{Precision} * \text{Recall}) / 2^2 * (\text{Precision} + \text{Recall})$$