




Follow us!

Share it! 

# LINEAR REGRESSION

LINEAR REGRESSION IS A WAY TO EXPLAIN THE RELATIONSHIP BETWEEN DEPENDENT VARIABLE AND ONE OR MORE EXPLANATORY VARIABLES USING A STRAIGHT LINE.

## LINEAR REGRESSION FORMULA



$$y = mx + c$$

**y** is the criterion variable  
**x** is the predictor variable  
**c** is the constant/intercept  
**m** is the regression coefficient

Linear regression can therefore, **predict** the value of (y) when only the (x) is **known**. It doesn't depend on any **other** factors.



## INTUITIVE EXPLANATION



Linear regression is one of the **simplest** and most commonly used data **analysis** and predictive **modelling** techniques. The linear regression aims to find an **equation** for a continuous **response** variable known as (y) which will be a **function** of one or more variables (x).

## ① Linear Regression

Regression  $\Rightarrow$  Linear regression  $\rightarrow y = mx + c$

Lasso regression (the tugline for linear regression)

Ridge regression is best fit line

Elastic Search (Lasso + Ridge)

## \* ML Project

① Collect data

② Analysis data (EDA)

$\rightarrow$  ~~Missing~~ Numerical data

$\rightarrow$  categorical data

$\rightarrow$  outliers

$\rightarrow$  distribution data

$\rightarrow$  Imbalance data

$\rightarrow$  balance data

③ Preprocess

④ ML Algo (model)

⑤ Evaluate

Regression

$R^2$

Adjusted  $R^2$

(adjusted  $R^2$ )

Classification

Accuracy

Error

F1

Precision

Recall

⑥ Retrained [If not getting proper result/accuracy]

[If I am getting low accuracy]

## Math

① Linear Algebra

② Calculus

③ Probability

④ Statistics

⑤ Geometry

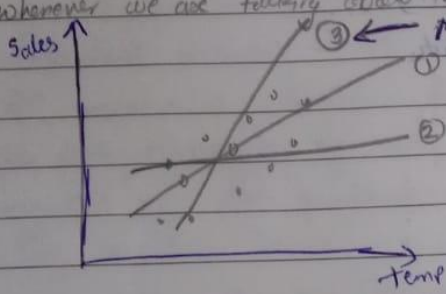
Linear Eq, System of Linear Eq, Angle, distance and length, Eigen value/Eigen vector, matrix

Differentiation, Partial diff, Maxima-minima, Integration

## \* Linear Regression

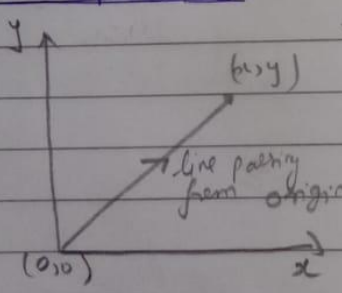
$$y = mx + c \quad \text{[Equation of line]}$$

→ whenever we are talking about LR, we have to find a best fit line (having lowest error) minimize error (Best fit line) This is called OLS (Ordinary least sq) (minimize the error)



$$y = mx + c \quad \dots \text{[Point Slope formula]}$$

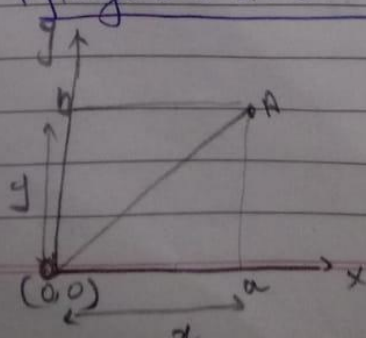
## • Vector Representation



$$\begin{aligned} \Rightarrow \hat{x} + y\hat{j} &= [x \ y] \\ \Rightarrow [x \ y] &= [x_1 \ x_2 \ x_3] \times [y_1 \ y_2 \ y_3] \end{aligned}$$

column vector  $\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \times 1 \times n$

## \* Pythagorean theorem



$$H^2 = B^2 + L^2$$

$$OA^2 = OA^2 + OB^2$$

$$OA = \sqrt{OA^2 + OB^2}$$

$$OA = \sqrt{x^2 + y^2}$$

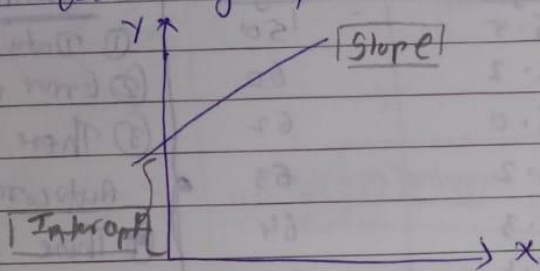
Euclidean



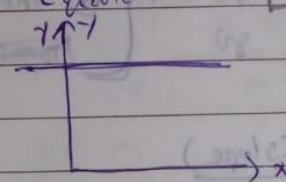
$$x_1, x_2 = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Euclidean Distance

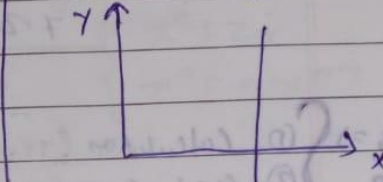
• General Equation of particular line  $\Rightarrow ax + by + c = 0$



• Horizontal Equation  $\Rightarrow y = x$



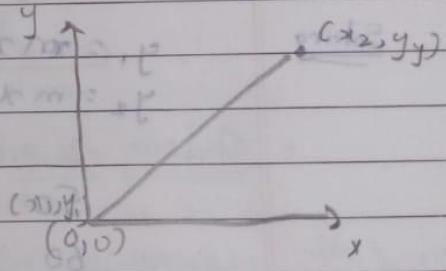
• Vertical Equation  $\Rightarrow x = y$



$$\text{slope } (m) = \frac{y_2 - y_1}{x_2 - x_1} \quad \text{or} \quad \frac{\Delta y}{\Delta x}$$

• If my line is going through origin

$$y = mx$$



- (1) point slope  $\Rightarrow y = mx + c$
- (2) General Equation of line  $\Rightarrow ax + by + c = 0$
- (3) Horizontal Equation  $\Rightarrow y = x$
- (4) Vertical Equation  $\Rightarrow x = y$
- (5) Slope  $(m) \Rightarrow m = \frac{y_2 - y_1}{x_2 - x_1}$  or  $\frac{\Delta y}{\Delta x}$
- (6) line passing through origin  $\Rightarrow y = mx$

$$-7 + 0 = -7$$

$$-7 + (0.3) \cdot 2 = -6.4$$

$$* y = mx + c$$

Dataset  $\rightarrow$  Based on Height we will predict weight

Height	Weight
5.5	50
5.7	60
6.0	62
6.2	63
6.3	64
6.8	69
7.1	70
7.2	80

Autocorrelation = Correlation values by itself

\* Assumptions (for LR)

- ① Data should be linear
- ② Error value tends to be 0
- ③ There should not be any autocorrelation between error
- ④ There should not be multicollinearity between columns
- ⑤ There should not be heteroscedasticity

Steps  $\Rightarrow$  for LB

- ① Calculation ( $y = mx + c$ )  $m = \frac{y}{x}$  (rate of changes)
- ② find Error (MSE)  $c \rightarrow$  Intercept
- ③ optimize (GD)  $\frac{\partial}{\partial m} \frac{1}{2} \sum (y_i - \hat{y}_i)^2 = 0$

Solving

$$y_1 = mx_1 + c$$

$$y_2 = mx_2 + c$$

choosing 2 points from the dataset, randomly

$$50 = m \times 5.5 + c \quad \text{--- ①}$$

$$60 = m \times 5.7 + c \quad \text{--- ②}$$

$$60 = 5.7m + (50 - 5.5m)$$

$$60 = 50 + 5.7m - 5.5m$$

$$60 = 50 + 0.2m$$

$$60 - 50 = 0.2m$$

$$10 = 0.2m$$

$$m = \frac{10}{0.2}$$

$$m = 50$$

$$c = 75$$

$$y = mx + c$$

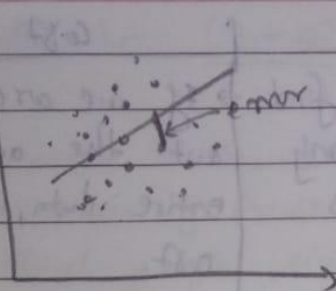
$$y = 50x + 75$$

$$y = 50(6.2) + 75$$

$$\Rightarrow y = 199$$



## \* Representation



$$\text{Error} = \sum_{i=1}^n (y_i - \hat{y}_i) \quad \text{where, } \boxed{\hat{y}_i = mx_i + c}$$

• Differentiation is a rate of change of a particular point

$$\frac{dy}{dx} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x}$$

$$\begin{cases} x^2 = 2x \\ x^n = nx^{n-1} \end{cases}$$

## Optimization formulae

### Gradient descent

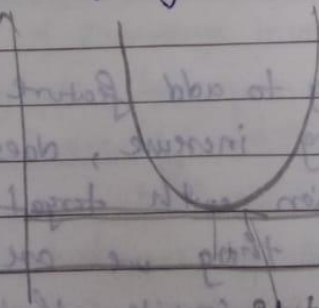
$$\text{slope} \rightarrow m_{\text{new}} = m_{\text{old}} - \eta \frac{d(\text{error})}{dm}$$

$$\text{Intercept} \rightarrow c_{\text{new}} = c_{\text{old}} - \eta \frac{d(\text{error})}{dc}$$

use of  $\eta$   $\Rightarrow \eta$  = we use Learning Rate ( $\eta$ ) for optimizing this  $m$  value for controlling the flow of this particular point. we don't want sudden change in our value we want minute change in our value

## GD

(Loss)  $y$



global minimum

- ~~Loss~~ <sup>Loss</sup>  
→ If we are going to find out different w.r.t only one point, hence it is called loss.
- <sup>Cost</sup>  
→ If we are going to find out the average w.r.t this entire data, hence this is called cost.

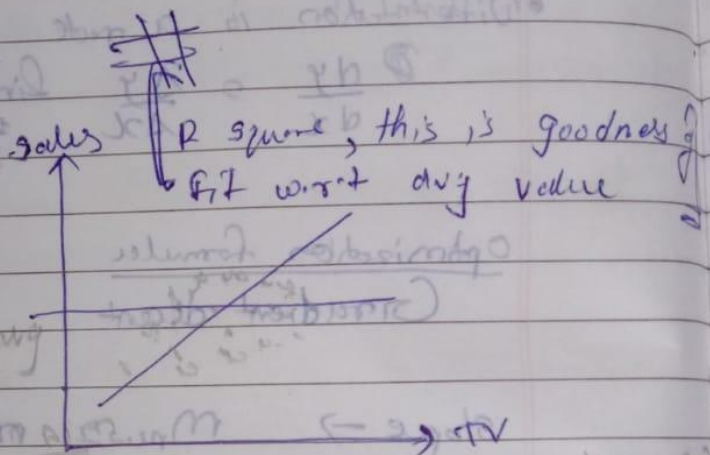
### • Homoscedasticity

$$y = m x + C$$

$\uparrow$  coefficient       $\uparrow$  intercept

$$R^2 = 1 - \frac{RSS}{TSS}$$

$$= \frac{TSS - RSS}{TSS}$$



### • OLS (Ordinary Least Square)

→ we need to find out optimal parameter  $m$  and  $C$  means we need to find out this best fit line w.r.t our data.

→ It is just a representation of this linear model.

### • Adj $R^2$

→ As we are going to add feature in Dataset. So my  $R^2$  value getting increase, doesn't matter if having any relation with target feature.

→ So, for solving this thing we are calculating Adj  $R^2$ .

→ Over there only we consider that feature which having any relation with target column.



• 
$$\text{Adj } R^2 = 1 - \frac{(1 - R^2)(N - 1)}{N - P - 1}$$

where,  $R^2$  = Sample R-Square

$P$  = Number of Predictors

$N$  = Total sample size

### • Multi - Collinearity

(more than one) (we are going to check dependency, we are going to check linearity)

→ if we have correlated between this two variables

TV	News	radio	Sales
$X_1$	$X_2$	$X_3$	$Y$

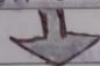
→ In multicollinearity, if we are having correlation between this two variable, this independent variable ~~it should not be~~.

→ In short, correlation between ~~variable~~ this independent variable we are going to ~~check~~ check the correlation.

Imp → More than one Linear relation in an Equation, hence the name multi-collinearity.

Eg:-

$\text{Salary} = a * (\text{years of Experience}) + b * (\text{age}) + c$  (A typical LR Equation)  
where, "a and b are coefficients" and "c is constant"



Here, age and years of Experience are correlated, since for a person as ~~they~~ the years of Experience increases, his/her age also increase. Mathematically,

Age = years of Experience  $\div d$  |  $d \rightarrow \text{constant, age when person started Job}$



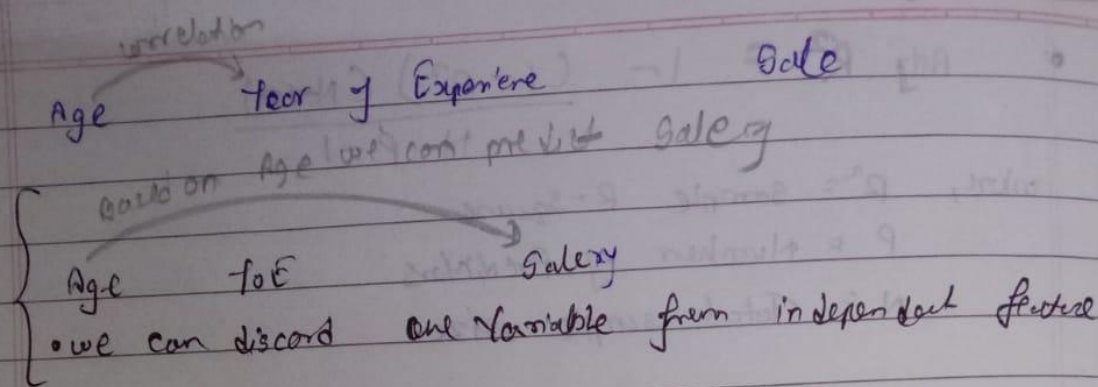
So, our LR Equation which was only supposed to have one linear relation, i.e.  $Y = mx_1 + mx_2 + c$  now has one more relation

→  $x_2 = ax_1 + d$



More than one Linear relation in an Equation, hence the name Multi-collinearity





$\rightarrow$  As the extent of the collinearity increases, there is a chance that we might produce an overfitted model. An overfitted model works with the test data but its accuracy fluctuates when exposed to other data sets.

- Variance Inflation Factor (VIF) = Regression of one X variable against other X variables.

$$VIF = \frac{1}{(1 - R^2)}$$

Correlation between a variable and other variables

### \* Regularization

low multicollinearity  $\rightarrow$  high multicollinearity

$\rightarrow$  When we use regression models to train some data, there is a good chance that the model will overfit the given training data set.

Regularization helps sort this overfitting problem by restricting the degrees of a given equation i.e. simply reducing the number of degrees of a polynomial function by reducing their corresponding weights.

In LR, we do not want huge weights/coefficient as a small change in weight can make a large difference for the dependent variable (Y). So, regularization constraints the weights of such features to avoid overfitting.

Note → In Linear Regression we are getting problem that is called overfitting.

→ we need to optimize this parameter  $m$  on  $C$  value.

overfitting means  $\Rightarrow$   $\begin{cases} \text{train acc is v.v. good} \\ \text{test acc is v.v. low.} \end{cases}$

we need to optimize our  $m$  on  $C$  parameter.

(11) Lasso (Least Absolute shrinkage and Selection Operator)

→ Lasso regression penalizes the model based on the sum of magnitude of the coefficients. The regularization term is given by regularization  $\Rightarrow \lambda * \sum | \beta_i |$

$$y = mx + c$$

$$WSS = (y - \hat{y})^2$$

$$m_{new} = m_{old} - \eta \frac{\partial L}{\partial m}$$

$$\boxed{(y - \hat{y})^2 + \lambda |m|}$$

$$m_{new} = m_{old} - \eta \frac{\partial [(y - \hat{y})^2 + \lambda |m|]}{\partial m}$$

$$m_{old} - \eta \frac{\partial [(y - (mx+c))\hat{y}] + \lambda |m|}{\partial m}$$

→  $\lambda$  is the shrinkage factor, and hence the formula for loss after regularization is

$$\boxed{\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j| = RSS + \lambda \sum_{j=1}^p |\beta_j|}$$

(12) Ridge

→ Ridge Regression penalizes the model based on the sum of squares of magnitude of the coefficient. The regularization term is given by

regularization  $= \lambda * \sum |\beta_j|^2$  ...  $\lambda$  is shrinkage factor

$$\boxed{\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p \beta_j^2 = RSS + \lambda \sum_{j=1}^p \beta_j^2}$$



## • Difference between Ridge and Lasso

→ Ridge regression shrinks the coefficient for those predictors which contribute very less in the model but huge weights, very close to zero. But it never makes them exactly zero. Thus, the final model will still contain all those predictors, though with less weights. This doesn't help in interpreting the model very well. This is where Lasso regression differs with Ridge regression. In Lasso, the L1 penalty does reduce some coefficient exactly to zero when we use a sufficiently large tuning ~~parameter~~ <sup>parameter</sup>  $\lambda$ . So, in addition to regularizing, Lasso also performs feature selection.

## Why use Regularization?

⇒ Regularization helps to reduce the variance of the models without a substantial increase in the bias. If there is variance in the model that means that the model won't fit well for dataset different than training data. The tuning parameter  $\lambda$  controls this bias and variance tradeoff. When the value of  $\lambda$  is increased up to a certain limit, it reduces the variance without losing any important properties in the data. But after a certain limit, the model will start losing some important properties which will increase the bias in the data. Thus, the selection of good value of  $\lambda$  is the key. The value of  $\lambda$  is selected using cross-validation methods. A set of  $\lambda$  is selected and cross-validation error is calculated for each value of  $\lambda$  and that value of  $\lambda$  is selected for which the cross-validation error is minimum.

Note:- The range of  $\lambda$  if we perform standardization in data then the minimum range is  $-3$  to  $+3$  (SD)  $\lambda$  is not more than  $m$  then then not less than  $m$

L1  
Lasso

L2  
Ridge

Path in  $\lambda$  loss  $+ \lambda |m_1 + m_2 + m_3|$   
Gradient function

$$m_{new} = m_{old} + \eta \frac{\partial l}{\partial m}$$

loss  $+ \lambda [ |m_1|^2 + |m_2|^2 + |m_3|^2 ]$   
Gradient function

$$m_{new} = m_{old} + \eta \frac{\partial l}{\partial m}$$