



Mo Tu We Th Fr Sa Su Date:

IA *ML*

MACHINE LEARNING ALGORITHMS

Introduction to ML (AI vs ML vs DL vs DS)

Machine learning is a subset or subfield of artificial intelligence (AI) that focuses on developing algorithms and statistical models that enables computers to learn and make predictions or decisions without being explicitly programmed to perform specific tasks. In essence it is a method of teaching computers to learn from data and improve their performance over time through experience.

What is AI?

AI, or Artificial Intelligence is a branch of Computer Science and Technology that focuses on creating intelligent machines and systems that can simulate human-like intelligence and perform tasks that typically require human intelligence.

What is an "AI" application?

AI application is able to do its own task without any human intervention.

Example: Netflix [Action movie Recommendation]

This Recommendation

is

called

Artificial Intelligence Application

Comedy movie Recommendation

: Amazon.in

↳ iPhone → Recommended →

Headphone

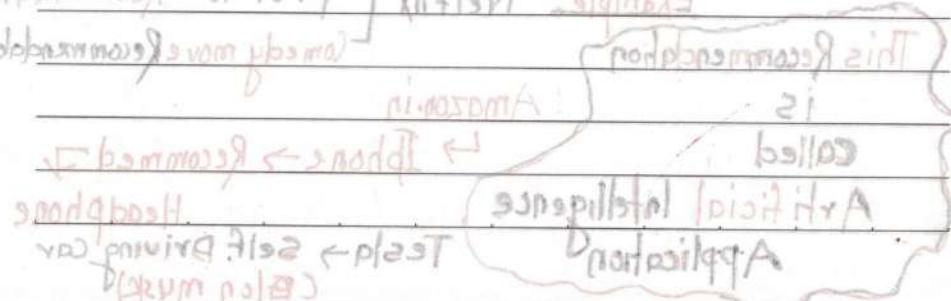
Tesla → Self Driving Car
(Elon Musk)

Multi Layered Neural Network

A Venn diagram with three overlapping circles. The top-left circle is labeled "ML". The top-right circle is labeled "DATA". The bottom circle is labeled "DL". The middle-left circle is labeled "MIMIC HUMAN BRAIN". The regions of overlap between ML and DATA, between DATA and DL, and between all three are shaded in red. The region between ML and DL is unshaded.

Role of ML \Rightarrow It provides stats tool to analyze the data
Visualize the data, to do prediction, forecasting

What is Deep Learning?
Deep learning is a subset or subfield of machine learning (ML) that focuses on training artificial neural networks with many layers (hence the term "deep") to perform tasks such as pattern recognition, image and speech recognition, natural language processing, and more.



MACHINE & Deep Learning

Supervised ML

- Regression
 - CLASSIFICATION

Unsupervised ML

- CLUSTERING
 - Dimensionality Reduction

Age	Weight	Output
21	45	basic
24	59	33.333333333333336
19	39	35
16	36	37

Take \rightarrow This model called (H_0)

\downarrow

Age \rightarrow Hypothesis \rightarrow Weight

\downarrow

Output

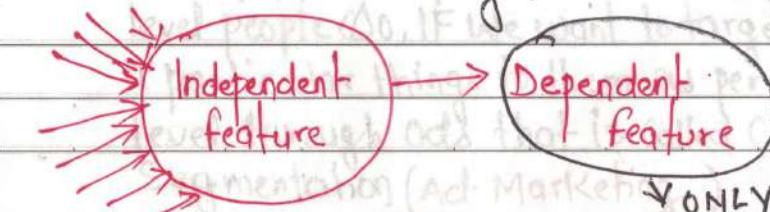
* Dependent feature
↓
Inweight (output)

* Independent feature

Age (Input)

Note: whenever we are solving Problems in the case of "Supervised ML" there, one Dependent feature

and there can be any number of independent feature



Regression Problem

Age	Weight → O/P
21	720
23	71
24	71.5

When output have a Continuous Variable then these became Regression Problem Statement

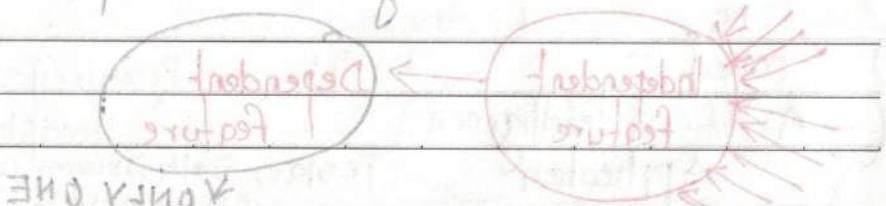
Classification Problem

No.of hours	No.of playhour	No.of sleep	(O/P)
INDEPENDENT FEATURE			Pass/Fail [Dependent]
*	*	*	P
*	*	*	F
*	*	*	P
*	*	(tugtoo)	F

Note:

Whenever we have Fixed Number of Category then that became a Classification Problem Statement

Suppose it has two output then it became Binary classification and if we have more than two different category it became Multiclass classification



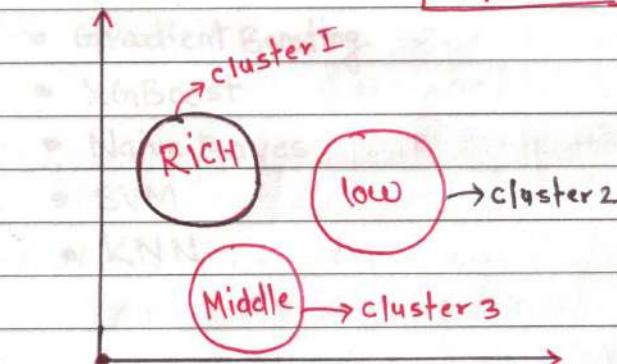
Unsupervised ML

Salary	Age	No Dependent Variable → No O/P feature
*	*	Number of autos → Input feature
*	*	Number of children → Input feature

In this scenario we use "Clustering"
"Clustering is the Grouping Algorithms" ↓
Customer Segmentation

Clustering basically means that based on data we have to find out similar group: eg: Group of People.

GROUPING

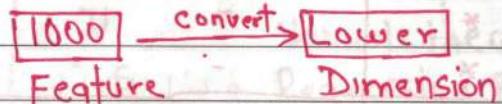


We are grouping them together

Eg: Suppose A Company launches a product "A" For Rich people and product "B" For Middle class people and product "C" For low income level people. If we want to target or send a particular things to them as per their level through ads that is called Customer Segmentation (Ad. Marketing)

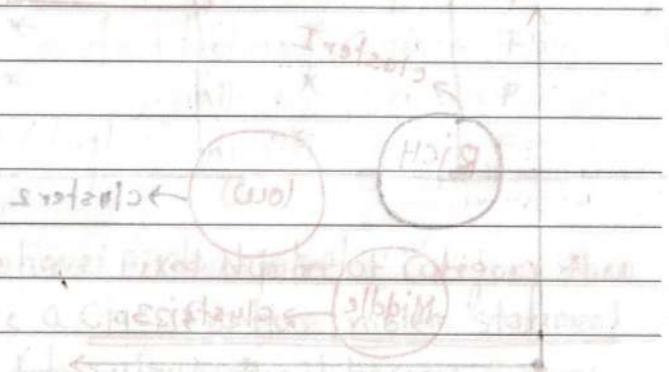
Dimensionality Reduction

Eg: Suppose if we have 1000 Feature, can we reduce this feature to lower dimension



Yes, it is possible with the help of Dimensionality Reduction Algorithms

PCA, LDA



"A" football or "B" football or "C" football :
"A" football or "B" football or "C" football or "D" football or "E" football or "F" football or "G" football or "H" football or "I" football or "J" football or "K" football or "L" football or "M" football or "N" football or "O" football or "P" football or "Q" football or "R" football or "S" football or "T" football or "U" football or "V" football or "W" football or "X" football or "Y" football or "Z" football

Machine Learning

Algorithms

Supervised

$$y = \theta_0 + \theta_1 x$$

Unsupervised

K-means

DBSCAN

Hierarchical clustering

K Nearest Neighbor clustering

Ada Boost

Random Forest

Gradient Boosting

XGBoost

Naive Bayes

SVM

KNN

What is linear Regression? = (x) or y

- Linear regression analysis is used to predict the value of a variable based on the value of another variable. The variable we want to predict is called the dependent variable. The variable we are using to predict the other variable's value is called the Independent variable.
- Linear regression is a fundamental machine-learning technique used for modeling the relationship between dependent variable (also called target or output) and one or more independent variable (also called features or predictors). It is commonly employed for tasks like predictive modeling and data analysis.
- In simple linear regression, there is only one independent variable, and the goal is to find a linear equation that best represents the relationship between the independent variable(s) and the dependent variable. Equation of Simple linear Regression

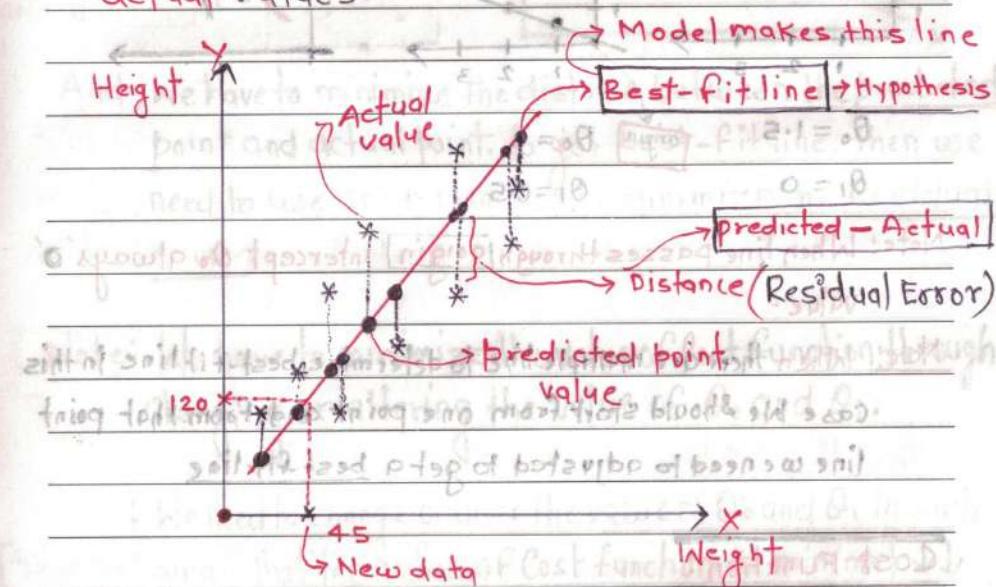
$$y = mx + c$$

Dependent Variable

Slope of line

Independent Variable

The objective of linear regression is to determine the values of m and c that minimize the difference between the predicted values ($mx+c$) and the actual values of the dependent variable 'Y' in the dataset. This is typically done using a method called Least squares, which finds the values of m and c that minimize the sum of squared differences between the predicted and actual values.



Residual Error = Distance Between Actual value point and predicted point value

Note: • If we sum of all distances the sum of Residual Error value must be very minimal to be a best-fit line

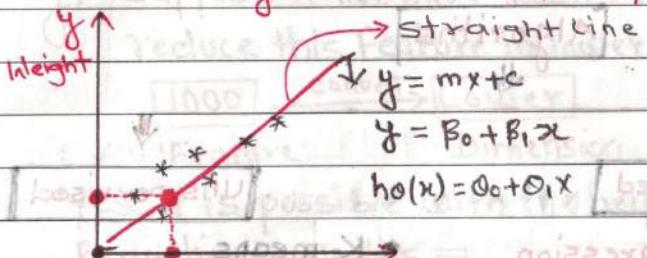
• Minimal Residual Error value will be Accepted even there are many lines.

To be best-fit line \Rightarrow Minimal Residual Error

Note: $[c \rightarrow]$ is the y-intercept (The point where the line intersects the y-axis)

LINEAR REGRESSION

Taking two feature X & Y



{Y is a linear function of X (Age)}

IN Linear regression we try to create a Model, with the help of above training dataset

TRAIN DATASET



Model



Hypothesis → OLP (Weight)

Take as input NEW AGE

In Linear Regression we are trying to find out a best-fit line which help to do prediction

LINEAR REGRESSION ⇒ It basically means we are going to create a linear line over there

Equation of straight line :

$$y = mx + c$$

$$y = \theta_0 + \theta_1 x^{(i)}$$

$$h_0(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

$$h_0(x) = \theta_0 + \theta_1 x$$

$$h_0(x) = \theta_0 + \theta_1 x$$

When, $x=0$

θ_0 = Intercept

θ_1 = Slope or Coefficient

$x^{(i)}$ = data points



↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

↓

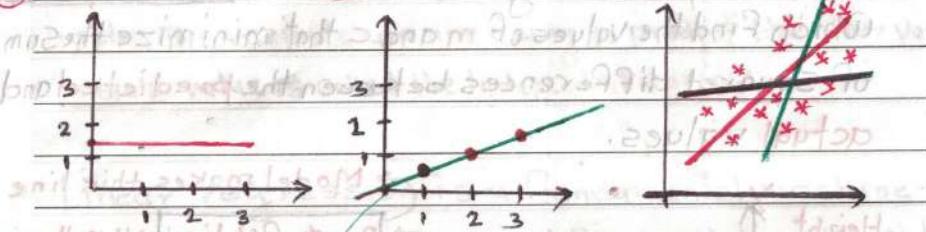
↓

Hypothesis $\{y \text{ is a linear function of } x\}$

$$h_0(x) = \theta_0 + \theta_1 x \Rightarrow \text{This is a Simple Hypothesis}$$

Which create or make the best-fit line Once over

- ① Top point
- ② Bottom point
- ③ Left-right



$$\theta_0 = 1.5$$

origin

$$\theta_0 = 0$$

$$\theta_1 = 0.5$$

Note: When line passes through origin intercept θ_0 always 0 value.

Note: When there are multiple line to determine best-fit line in this case we should start from one point and from that point line we need to adjusted to get a best-fit line.

Cost Function:

$$\sum_{i=1}^m \frac{1}{2m} [h_0(x^{(i)}) - y^{(i)}]^2 \rightarrow \text{This is Cost function and our main aim to minimize its value to get a best-fit line.}$$

OR

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}]^2$$

Squared Error Function

$m \Rightarrow$ No. of data point

Note: We have to find out sum of error Average that's why we divide by m .

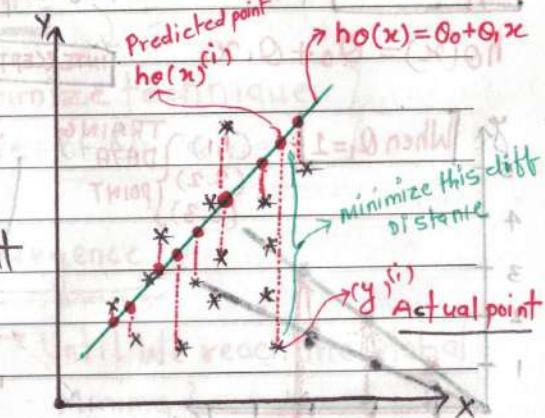
Explanation:

$$h_0(x) = \theta_0 + \theta_1 x^{(i)}$$

Hypothesis

$$h_0(x)^{(i)} \rightarrow \text{Predicted point}$$

$$y^{(i)} = \text{Actual point}$$



And, We have to minimize the distance between the predicted point and actual point. to get best-fit line. then we need to use Cost function to minimize the Residual error.

Note: We have to minimize the value of Cost function through changing or altering the value of θ_0 and θ_1 .

: We need to change or alter the value of θ_0 and θ_1 in such away that the value of Cost function is minimal.

: Differentiation is used to find out that point slope

$$0.0 = 1.8 \quad 2.5 \approx = (10)$$

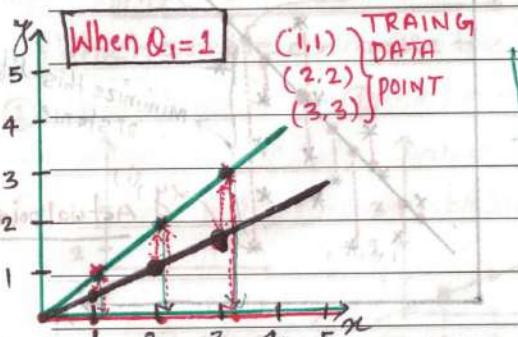
10. As usually add square yll be denoted the square, the square of sum of square difference is small value.

Hypothesis

Let's Consider, $\theta_0 = 0$

$$h_0(x) = \theta_0 + \theta_1 x_1$$

$$\text{INTERCEPT } J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}]^2$$



$$h_0(x) = 1 \times 1 \quad J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}]^2$$

$$\text{Input } h_0(x) = 1 \times 2 \quad = \frac{1}{2 \times 3} [(0)^2 + (0)^2 + (0)^2]$$

$$h_0(x) = 1 \times 3$$

$$[J(\theta_1) = 0] \quad \& \quad \theta_1 = 1$$

Now plot in above $J(\theta_1) = 0 \quad \& \quad \theta_1 = 1$

When, $\theta_1 = 0.5$

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}]^2$$

$$= \frac{1}{2 \times 3} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

$$[J(\theta_1) = \approx 0.58] \quad \& \quad \theta_1 = 0.5$$

When, $\theta_1 = 0.0$

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}]^2$$

$$= \frac{1}{6} [(0-1)^2 + (0-2)^2 + (0-3)^2]$$

$$[J(\theta_1) = \approx 2.3] \quad \& \quad \theta_1 = 0.0$$

Note: GRADIENT DESCENT automatically change the value of θ_1

Note: Unless and Until we have to change the θ_1 value when not reach at minimal point

Cost Function

Convergence Algorithm

It is a optimize technique

Optimize the Changes of θ_1 values

Repeat Until Convergence

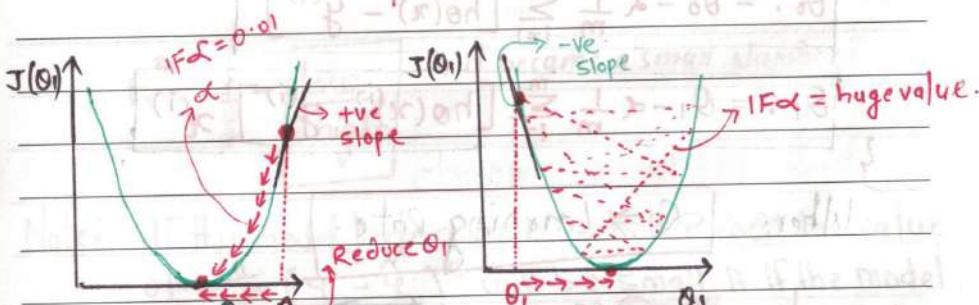
Until we reach the Global Minima }

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_j)}{\partial \theta_j}$$

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

Learning Rate

Note: This equation is actually help us to change our θ_1 value much more efficiently.



$$\theta_1 := \theta_1 - \alpha (+ve)$$

$$\theta_1 := \theta_1 - \alpha (-ve)$$

$$\theta_1 := \theta_1 + \alpha (-ve)$$

Note: By what speed we should be coming to near Global Minima $\rightarrow \alpha$
Learning Rate usually a small value,

Alpha

Find out derivative of J with respect to θ_0 & θ_1 , respectively. $(J(\theta_0, \theta_1))$

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m [(h_\theta(x^{(i)}) - y^{(i)})]^2$$

When, $J = 0 \& 1$

When, $J = 0 \Rightarrow \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$

When, $j = 1$

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

Note: Derivative of x^2 is $2x$.

Convergence theorem

Repeat Until Convergence

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}]$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}] x^{(i)}$$

Where, $\alpha \Rightarrow$ Learning Rate.

Derivative of $\frac{\partial}{\partial \theta_j} (\theta_0, x) = x$

$$\frac{\partial}{\partial \theta_0} = 1$$

$$\frac{\partial}{\partial \theta_1} = 1$$

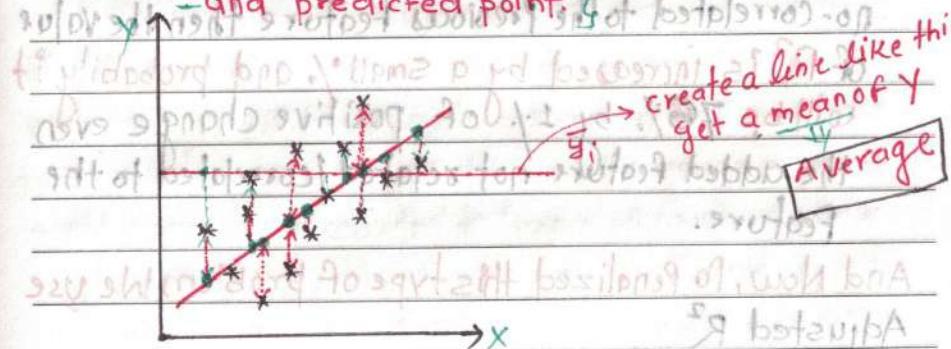
R squared

$$R^2 = \frac{1 - \frac{SS_{Res}}{SS_{Total}}}{SS_{Total}} \Rightarrow \begin{array}{l} \text{Sum of square Residual} \\ \text{Sum of square Total} \end{array}$$

$$R^2 = 1 - \frac{\sum_{i=1}^m (y_i - \hat{y}_i)^2}{\sum_{i=1}^m (y_i - \bar{y})^2}$$

$h_\theta(x)$
Predicted point
Mean of y

Note: Sum of Residual Error is difference between Actual point and predicted point.



$$R^2 = 1 - \frac{\text{low}}{\text{high}} \rightarrow \text{It became a small Number}$$

Note: If the model performance is better then the value of $\sum_{i=1}^m (y_i - \hat{y}_i)^2$ will be small. If the model

performance is bad it means that the value of $\sum_{i=1}^m (y_i - \bar{y})^2$ will be big.

Note: If $R^2 = 100\%$, then it is called Overfitting but Not Possible.

Adjusted R squared

Dataset

Let's Consider $R^2 = 71\%$.

Gender	Size of House	No. of Bedrooms	location	Price of House

Explanation

IF We have three feature our R^2 is 71% . If we add one feature more our R^2 is changed to 75% . But the feature we have added that is correlated to the previous feature that's why R^2 increased by 4% . But If we add an additional feature that is no-correlated to the previous feature then the value of R^2 is increased by a small %, and probably it will be 76% , by 1% of positive change even the added feature not related/correlated to the feature.

And Now, To Penalized this type of problem we use Adjusted R^2 .

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

Where, N = No. of datapoints

P = No. of independent Feature.

\downarrow No. of Predictors

Note: If any feature are added to our dataset that is not important and not correlated for that dataset and due to this unnecessary feature R^2 increased then after we need to calculate Adjusted R^2 square (R^2) value which minimize the value of R^2

Note: $\text{Adjusted } R^2 < R^2$

Ridge Regression

Ridge regression, also known as L2 regularization, is a technique used to prevent overfitting in linear regression, the objective is to minimize the sum of squared

loss function with a added a regularization term to the loss function.

which is $\frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \beta_j^2$

where y_i is the observed response value for the i th observation, \hat{y}_i is the predicted response value for the i th observation, β_j is the coefficient for the j th predictor variable, and λ is the regularization parameter.

The term $\lambda \sum_{j=1}^p \beta_j^2$ is called the L2 regularization term.

It is used to prevent overfitting by adding a penalty term to the loss function.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

As λ increases, the magnitude of the coefficients decreases.

Ridge Regression \rightarrow Ridge regression, also known as L2 Regularization, is a technique used to prevent overfitting in linear regression models. by adding a penalty term to the loss function. In traditional linear regression, the objective is to minimize the sum of squared differences between the predicted value and actual values (least squares). However, this can lead to overfitting when there are many input features or when some features are highly correlated.

Ridge Regression

adds a regularization term to the least squares loss function, which is proportional to the sum of the squared coefficients of the linear regression model. The regularization term encourages the model to keep the coefficients small, which in turn reduces the impact of individual features on the prediction. This helps in mitigating multicollinearity (High correlation between features) and improves the Generalization of the model.

The Ridge Regression loss function \rightarrow

$$\text{Loss} = \text{Least Squares Loss} + \lambda * (\text{Sum of Squared Coefficients})$$

Where,

λ (lambda) is a regularization parameter, which controls the strength of regularization. A higher λ value results in more Aggressive Regularization.

$$\text{Cost Function} = (h_0(x^{(i)}) - y^{(i)})^2 + \lambda (\text{slope})^2 \Rightarrow \text{Ridge Regression (L2 Norm)}$$

Purpose: Preventing Overfitting

Lasso Regression

\hookrightarrow L1 Regularization

Lasso Regression, or L1 regularization is a technique used to prevent overfitting and perform feature selection in linear regression models. It is similar to Ridge Reg. It adds a Penalty term to the loss function, but instead of using the sum of squared coefficients, it uses the sum of the absolute values of the coefficients.

The Lasso Regression loss function \rightarrow

$$\text{Loss} = \text{Least Squares Loss} + \lambda * (\text{Sum of Absolute Values of Coefficients})$$

Note: Like Ridge Regression, λ controls the strength of regularization. In Lasso, however, one key difference is that Lasso tends to drive the coefficients of less important features to exactly zero. This property makes Lasso useful for feature selection. Because it effectively eliminates irrelevant features from the model.

$$\text{Cost function} = (h_0(x^{(i)}) - y^{(i)})^2 + \lambda |\text{slope}| \Rightarrow \text{L1 Reg.}$$

Purpose: Prevent Overfitting. } mode
: Feature Selection. } slope

GRADIENT DESCENT

Gradient Descent is an algorithm that finds best fit line for given training data set.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - y_{predicted})^2$$

Mean Squared Error

$$\text{Cost function} \Rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - (m\chi_i + b))^2$$

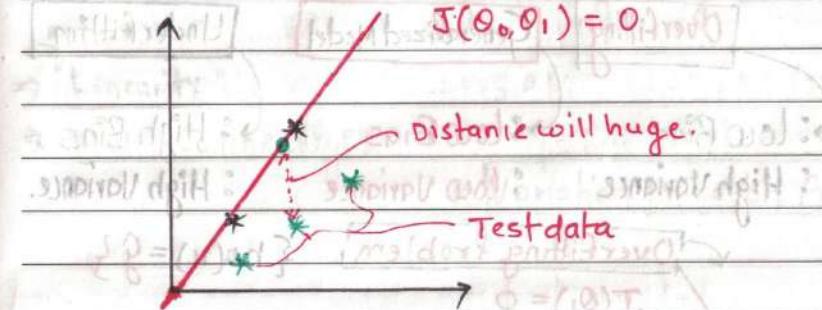
Note:

The Algorithm that are used to find out MSE is called Gradient Descent.

Cost function $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^n [h_\theta(x^{(i)}) - y^{(i)}]^2$

$$\theta_0 = 0$$

$$J(\theta_0, \theta_1) = 0$$



Note: Plotted two points are called Training data

Model is trained by Training Data

Overfitting \Rightarrow Model performs well \rightarrow Training data.

Fails to perform well \rightarrow Test data

High variance

Low Bias

Underfitting \Rightarrow Model Accuracy is bad with \rightarrow Training data.

Model Accuracy is also bad with \rightarrow Test data.

High variance

High Bias

Example**Model 1**

Training Accuracy = 90%. 92%.

Test Accuracy = 80%. 91%.

**Model 2****Model 3**

70%.

65%.

92%.

91%.

↓

Overfitting**Generalized Model****Underfitting**

: low Bias

: High Variance

: low Bias

: low Variance

: High Bias

: High Variance

Overfitting Problem{ $h_0(x) = g$ }

$J(\theta_1) = 0$

$= \frac{1}{2m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}]^2$

$= (\hat{y} - y^{(i)})^2 + \lambda (\text{slope})^2$

$= 0 + 1(2)^2$

$= 4$

Note: θ_0 is intercept of line.

$h_0(x) = \theta_0 + \theta_1 x$

When, $\theta_0 = 0$

then, $h_0(x) = \theta_1 x$ slope

$= (\hat{y}^{(i)} - y^{(i)})^2 + \lambda (\text{slope})^2$

$= (\text{small value}) + 1(1.5)^2$

$= \text{small value} + 2.25$

$\approx 3 \downarrow \downarrow$

Note: We have to reduce the value unless and until we get a line which is able to handle the **Generalized Model**.**Assumption of linear Regression**

If our features are in: Normal/Gaussian Distribution.
Or If our features follow these distribution then it's obviously good a model will get trained value.

Standardization: scaling data by using $Z\text{-score}$, $\mu=0$, $\sigma=1$.

Linearity: $y = \theta_0 + \theta_1 x + \epsilon$ no bound

Multicollinearity: x_1, x_2, \dots, x_n are highly correlated.

*** Variation Inflation Factor**

Logistic Regression

- Logistic Regression is a statistical and Machine learning model used for Binary classification tasks. It is a type of Regression analysis that predicts the probability of a binary outcome (1/0, YES/NO, True/False) based on one or more predictor variables. Despite its name, logistic regression is used for classification, not regression.
- Logistic regression is a type of Machine learning algorithm that are used for classification tasks, and models the probability that a sample belongs to a certain class using a logistic function.
- This is the first type of Algorithm in Classification

Example:

No. of Study(hr.)	No. of play(hr.)	P/F (Two category)
-	-	P
-	-	F

Binary (Y/N)

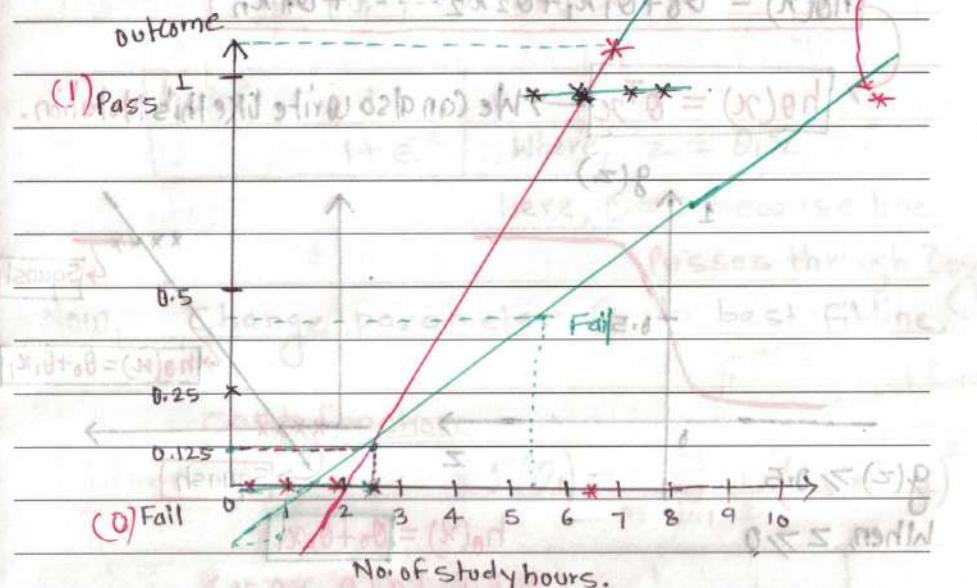
Note: Logistic Regression is very well with Binary classification

- We can solve multiclass classification with logistic regression.

Note: We have to reduce the value unless we get a line which is able to handle the Generalized Model.

Condition: $h_0(x) \leq 0.5 \rightarrow 0 \rightarrow \text{Fail}$

$h_0(x) \geq 0.5 \rightarrow 1 \rightarrow \text{PASS}$



Sigmoid Function

(Linear Regression can not be used to solve the above Problem)

Note: The above problem is related to Binary classification Problem that's why its value lies between 0 to 1

: Squash \Rightarrow to make straight

midrange of std $\times \times \times \times$ Squash

xxxx

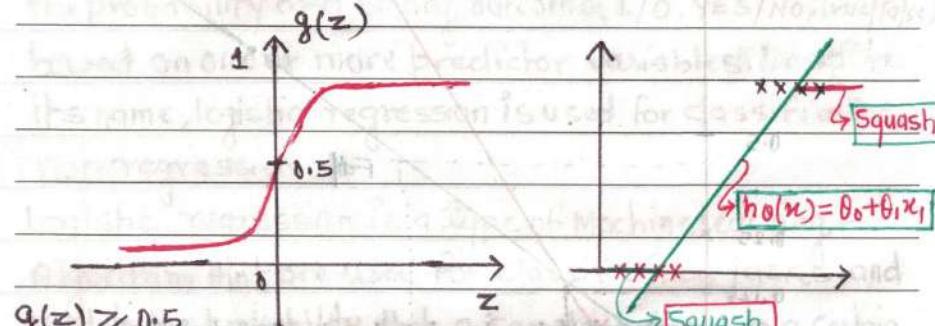
squash

Decision Boundary in case of Logistic Regression

Lets Consider,

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$h_0(x) = \theta^T x \rightarrow \text{We can also write like this Notation.}$$



$$h_0(x) = \theta_0 + \theta_1 x_1$$

$$h_0(x) = g(\theta_0 + \theta_1 x_1)$$

let,

$$z = \theta_0 + \theta_1 x_1$$

$$h_0(x) = \frac{1}{1 + e^{-z}}$$

Sigmoid or logistic function.

$$h_0(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}} \rightarrow \text{It squash the function}$$

This is Hypothesis

- With the help of Regression we are creating the straightline.
- With the help of the concept of Sigmoid We are able to Squashing

Training set $\{(x^1, y^1), (x^2, y^2), (x^3, y^3), \dots, (x^n, y^n)\}$

$y \in \{0, 1\} \rightarrow \text{Two Output}$

$$h_\theta(z) = \frac{1}{1 + e^{-z}}, \text{ where, } z = \theta^T x$$

Here, $\theta_0 = 0$, because line passes through origin

Now, Change parameter θ_1 to best fit line.

cost function

$$\text{Linear Regression: } J(\theta_1) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\text{logistic Regression} \rightarrow h_\theta(x) = \frac{1}{1 + e^{-(\theta^T x)}}$$

$$\text{cost function: } \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

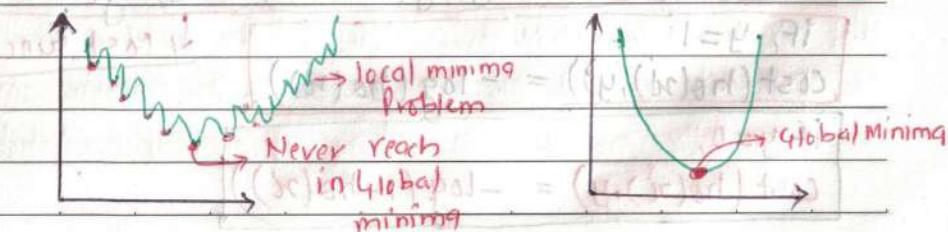
We can not use this cost function for logistic regression.

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

$$1 = p \neq 1 - (1 + e^{-(\theta^T x)}) \rightarrow \text{Non-Convex Function.}$$

Non Convex Function

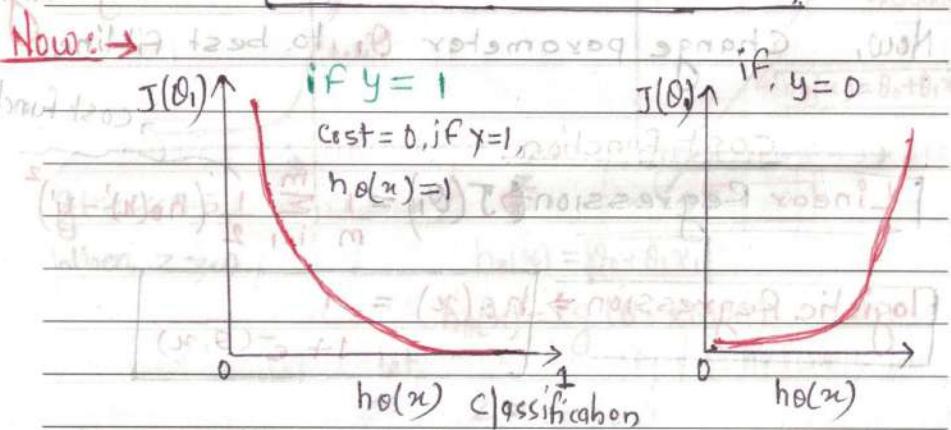
Gradient Descent



LOGISTIC REGRESSION Cost Function

$$J(\theta_0) = \begin{cases} -\log(h_0(x^i)) & y=1 \\ -\log(1-h_0(x^i)) & y=0 \end{cases}$$

Where, $h_0(x) = \frac{1}{1+e^{(\theta_0 x)}}$



$$\text{Cost}(h_0(x^i), y) = \begin{cases} -\log(h_0(x^i)) & \text{if } y=1 \\ -\log(1-h_0(x^i)) & \text{if } y=0 \end{cases}$$

$$\text{Cost}(h_0(x^i), y) = -y \log(h_0(x^i)) - (1-y) \log(1-h_0(x^i))$$

If, $y=1$ \downarrow Cost Function
 $\text{cost}(h_0(x^i), y) = -\log(h_0(x^i))$

If $y=0$,
 $\text{cost}(h_0(x^i), y) = -\log(1-h_0(x^i))$

$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(h_0(x^i)) + (1-y^i) \log(1-h_0(x^i))]$

$h_0(x^i) = \frac{1}{1+e^{(\theta^T x^i)}}$

(Repeat it Until Convergence)

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

Performance Metrics

- Performance metrics in machine learning are measures used to evaluate the performance of a machine learning model. These metrics help assess how well a model is making predictions or classifications on a given dataset.
- The choice of Performance metrics depends on the specific task and goals of the Machine learning Project.

⇒ Some of the Common Performance metrics used in Machine learning.

Accuracy:

This is the most straightforward metric and is simply the ratio of correctly predicted instances to the total instances. It is suitable for balanced datasets, where the classes are roughly equally distributed.

$$\text{Accuracy} = \frac{\text{Number of correct Predictions}}{\text{Total Number of Predictions}}$$

→ **Confusion Matrix** → A confusion matrix is a table used to evaluate the performance of a classification algorithm. It provides a more detailed breakdown of the model's predictions compared to a single accuracy score. The confusion matrix is especially useful when dealing with imbalanced datasets.

The confusion matrix is a matrix used to determine the performance of the classification models for a given set of test data. It can only be determined if the true values for test data are known. The matrix itself can be easily understood, but the related terminologies may be confusing. Since it shows the errors in the model performance in the form of a matrix, hence also known as an error matrix. : source → javatpoint.com Only This One

Components of a confusion matrix

Suppose we have a binary classification problem (two classes: positive and negative)

- True Positive (TP): Instances that are actually positive and are correctly classified as positive.
- True Negative (TN): Instances that are actually negative and are correctly classified as negative.
- False Positive (FP): Instances that are actually negative but are incorrectly classified as positive. (Type I error)
- False Negative (FN): Instances that are actually positive but are incorrectly classified as negative. (Type II error)

Confusion Matrix Presentation → 2x2 table

		Predicted Class	
		Negative	Positive
Actual Class	Negative	TN	FP
	Positive	FN	TP