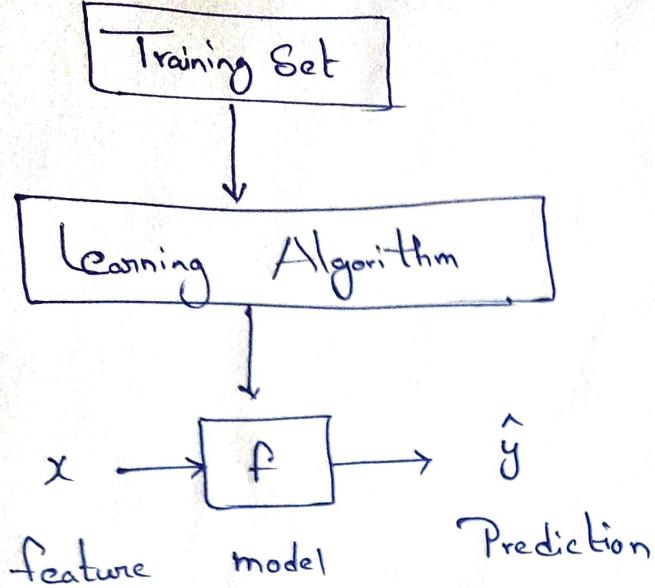


What is Machine Learning

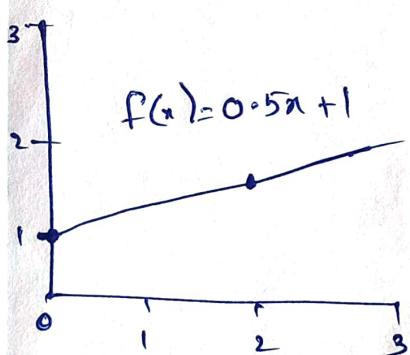
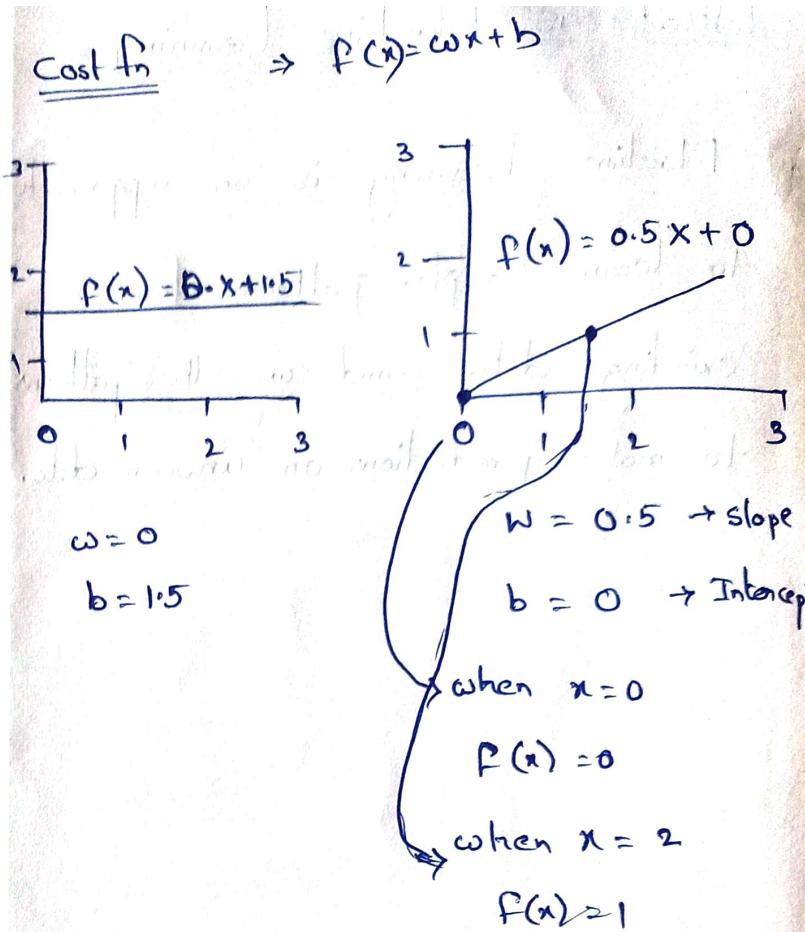
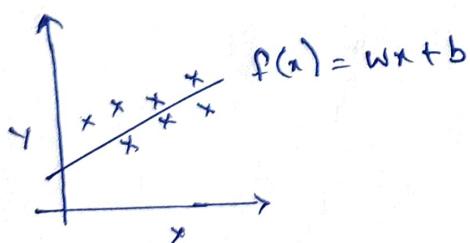
- * Machine Learning is an approach to learn complex patterns from existing data and use the patterns to make predictions on unseen data.



How to represent f ?

$$f_{w,b}(x) = wx + b \quad (\text{or} \quad f(x) = wx + b)$$

\hookrightarrow f takes the values of input ' x ' & depending on values of w & b it will output some value which is the predicted value (\hat{y})



$$w = 0.5$$

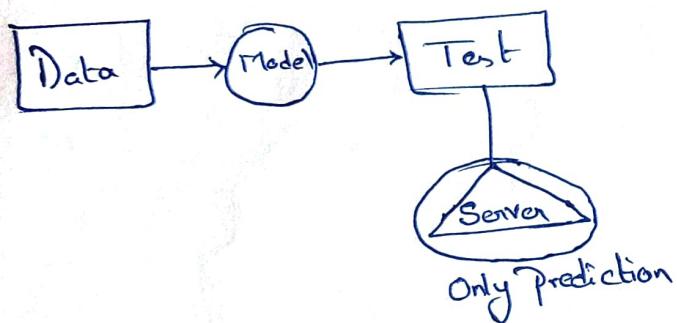
$$b = 1$$

Batch vs Online Machine Learning

Batch

↪ We train the model in a single shot using entire data. We do this offline.

↪ When the model is trained then we deploy it on Server



Con's

① We train the model once & then deploy it. There is no retraining the model on the new data. As a result overtime the model doesn't perform very well on new data.

Solution



This Cycle keeps on going (2wks, Week, 1month, ...)

② Lots of Data

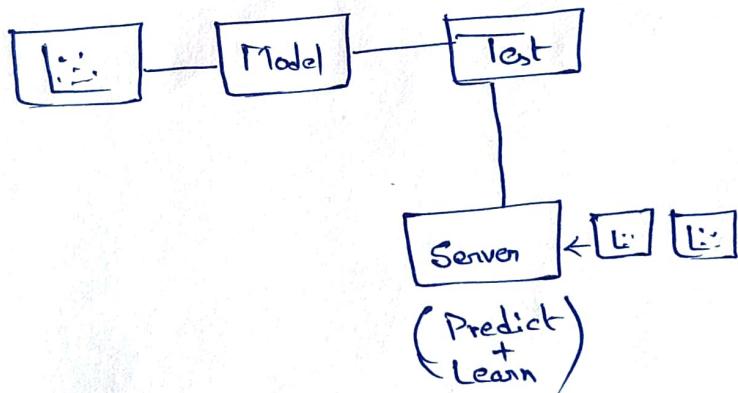
↪ Can't process with whole data

③ Hardware Limitation

④ Availability

Online → River Library

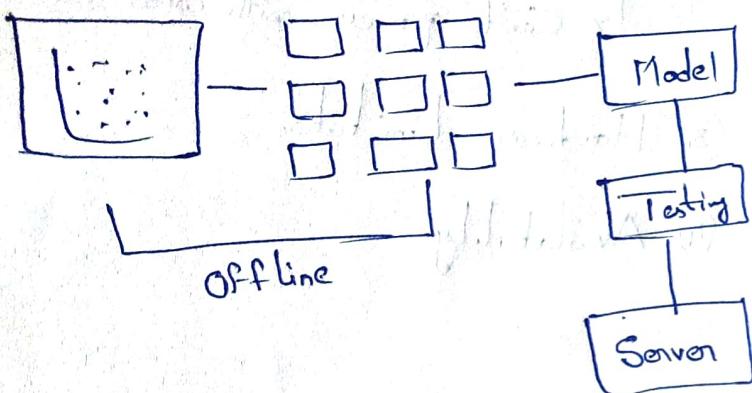
↪ It is done incrementally on the Server using mini batches



When to Use?

1. When there is fast concept/context Changer
2. Cost Effective
3. Faster Solution

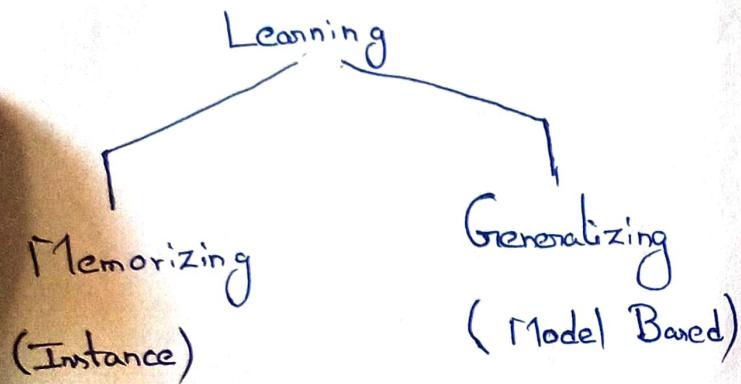
Out of Core Learning



Here we are using the technique of online learning in offline.

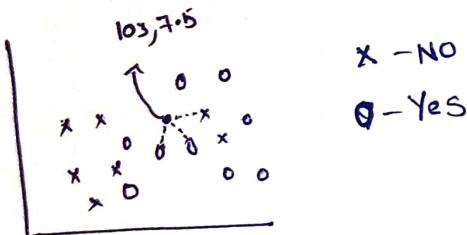
We are splitting the large data into small chunks & then we are training the model incrementally.

Instance Vs Model Based Learning



Instance Based

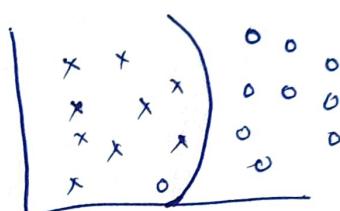
iq	CgPa	placement
80	8	Y
70	7	N



Here when even a new point comes we try to close it is to its three neighbors if two are yes and one neighbor is no. then the new point belongs to yes

Model Based

iq	cgpa	placement
80	8	Y
70	7	N



It tries to understand the behaviour of the points & based on that it generalizes.

Tensors

↳ A Container to store numbers

2 3
 ↓
Scalar → 0D-Tensor

$[1, 2, 3, 4] \rightarrow 1D\text{-Array}/\text{array}$

↳ Vector → 1D Tensor

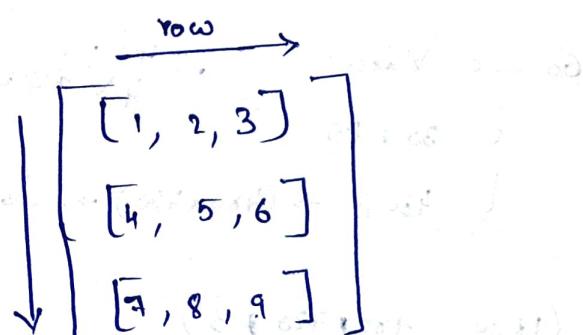
No. of axis = Rank (of Tensor)
= dimension

→ 1D Tensor

Vector → 4D

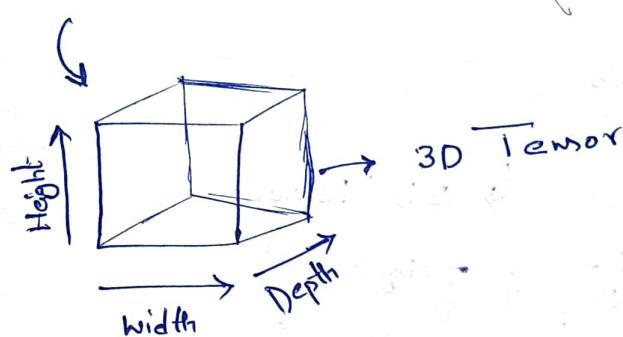
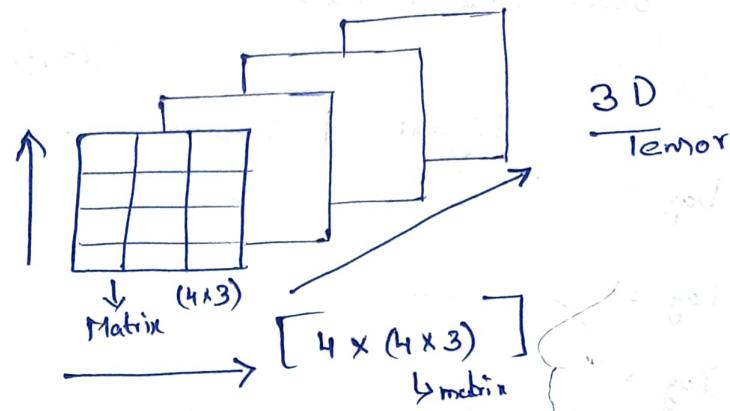
Ex- $[1, 2]$ → 1D Tensor

Vector → 2D



Matrix → 2D Tensor

ND Tensors



Ex-

NLP

Hi	Rahul	1	0	0	0	0
Hi	Ankit	0	1	0	0	0
Hi	Bala	0	0	1	0	0

2D Tensor
Hi Rahul → $\begin{bmatrix} [1, 0, 0, 0], [0, 1, 0, 0] \end{bmatrix}$

Hi Ankit → $\begin{bmatrix} [1, 0, 0, 0], [0, 0, 1, 0] \end{bmatrix}$

Hi Bala → $\begin{bmatrix} [1, 0, 0, 0], [0, 0, 0, 1] \end{bmatrix}$

3D Tensor
 $\begin{bmatrix} \begin{bmatrix} [1, 0, 0, 0], [0, 1, 0, 0] \end{bmatrix}, \begin{bmatrix} [1, 0, 0, 0], [0, 0, 1, 0] \end{bmatrix} \end{bmatrix}, \begin{bmatrix} [1, 0, 0, 0], [0, 0, 0, 1] \end{bmatrix}$

Ex 2

Time Series Data

Stock Market Data

(365, 2)

Highest | Lowest

Day 1

— | —

Day 2

— | —

Day 3

— | —

Day 365

— | —

↳ 2D Tensor

$$\left[\begin{array}{c} [-, -] \\ [-, -] \\ [-, -] \end{array} \right] \rightarrow \text{Day 1} \quad \left[\begin{array}{c} [-, -] \\ [-, -] \\ [-, -] \end{array} \right] \rightarrow \text{Day 365}$$

This is for One Year.

For 10 years

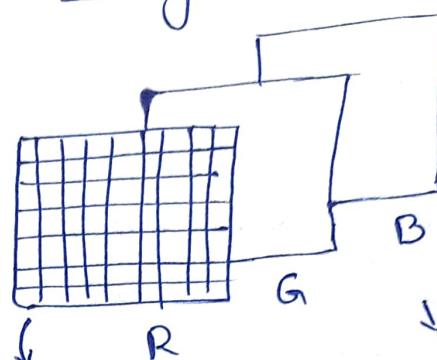
$$\Rightarrow (10, 365, 2)$$

↳ 3D Tensor

$$\left[\begin{array}{c} \left[\begin{array}{c} [-, -] \\ [-, -] \end{array} \right] \rightarrow \text{One Year} \\ ; \\ \left[\begin{array}{c} [-, -] \\ [-, -] \end{array} \right] \rightarrow 10^{\text{th}} \text{ year} \end{array} \right] \quad \left\{ \begin{array}{l} \text{3D} \\ \text{Tensor} \\ (10, 365, 2) \\ \downarrow \\ \text{Time Axis} \end{array} \right.$$

4D Tensors

↳ Images



(1200, 800)

↓ 1 Color Image
(3, 1200, 800)
↳ 3D Tensor

If we have 50 color Images

$$\Rightarrow (50, 3, 1200, 800)$$

↳ 4D Tensor

5D Tensors

↳ Videos

Videos are basically Collection of frames.

frames are nothing but images.

60 sec Video

→ 4D videos

↳ 30 FPS

↳ 480 P → (480 × 720) → 3 channels

$$\Rightarrow (1800, 480 \times 720 \times 3) \quad \text{↳ 4D Tensor}$$

$$\Rightarrow (4, (1800, 480 \times 720 \times 3)) \quad \text{↳ 5D Tensor}$$

Understanding Data

Asking Questions

① How big is the data?

↳ df.shape

② How does the data look like?

↳ df.head() / df.sample(5)

③ What is the data types of Col?

↳ df.info

④ Are there any missing values?

↳ df.isnull().sum()

⑤ How does the data look mathematically?

↳ df.describe()

⑥ Are there duplicate values?

↳ df.duplicated.sum()

⑦ How is the Correlation b/w col?

↳ df.corr()

⑧ What is the distribution of data?

Exploratory Data Analysis

Univariate Analysis

Categorical Data

a) Countplot

↳ sns.countplot(df['Survived'])

↳ df['Survived'].value_counts

b) PieChart

df['Survived'].value_counts.plot(kind='pie',
 autopct='%0.2f')

Numerical Data

a) Histogram

↳ import matplotlib.pyplot as plt
plt.hist(df['Age'])

b) Distplot

sns.distplot(df['Age'])

↳ Probability Density Function

c) Boxplot

↳ sns.boxplot(df['Fare'])

d) Min, Max, Mean

Bi-Variate Analysis

① Scatterplot (Numerical - Numerical cols)

↳ sns.scatterplot(df['total'], df['tip'])
↳ Bivariate

↳ sns.scatterplot(" ", " ", hue=df['sex'])

↳ Multivariate Analysis

② Bar Plot (Numerical - Categorical)

↳ sns.barplot

③ Box Plot (Numerical - Categorical)

④ Distplot (Numerical - Categorical)

⑤ HeatMap (Categorical - Categorical)

⑥ ClusterMap (Categorical - Categorical)

⑦ Pairplot

⑧ Lineplot (Numerical - Numerical)

Feature Engineering

① Feature Transformation

- Missing Values → One-Hot Encoding
- Handling Categorical Features
- Outlier Detection
- Feature Scaling

② Feature Construction

③ Feature Selection

④ Feature Extraction

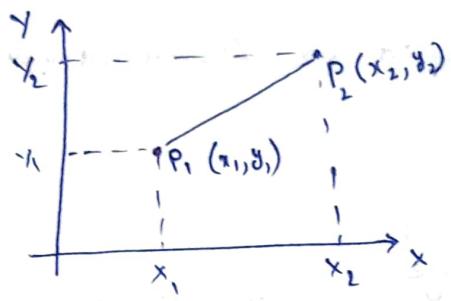
Feature Construction

↳ You Create a Completely new Column from the existing features.

SibSp	Parch

SibSp → ^{col}
 SibSp & Spouse
 → Parent & Child

Feature Scaling



Euclidean Distance b/w P_1 & P_2

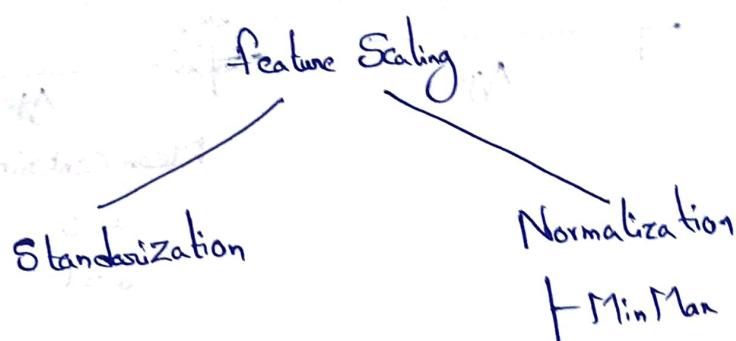
$$= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

x	y	
Age	salary	Purchase
50	83000	1
27	48000	0

$$(x_2 - x_1)^2 = 529$$

$$(y_2 - y_1)^2 = 1225000000$$

↳ This effects the euclidean distance.



Standardization

↳ Also called as Z-Score Normalization

f₁

Age → Age^{Transform}

27 → 0.08

$$Z = \frac{x_i - \mu}{\sigma}$$

$\mu \rightarrow$ Mean
 $\sigma \rightarrow$ Standard deviation
 $G \Rightarrow$ Sigmoid

$$\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n}}$$

↳ no. of data points

29

30

31

32

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

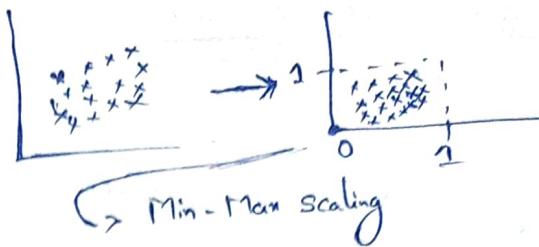
:

:

Feature Engineering

- ↳ Feature Engineering is the process of using domain knowledge to extract features from raw data.

These features can be used to improve the performance of ML algorithms.



Max Abs Scaling

- ↳ Used when we have Sparse data

Sparse data — Where 0's are more in data

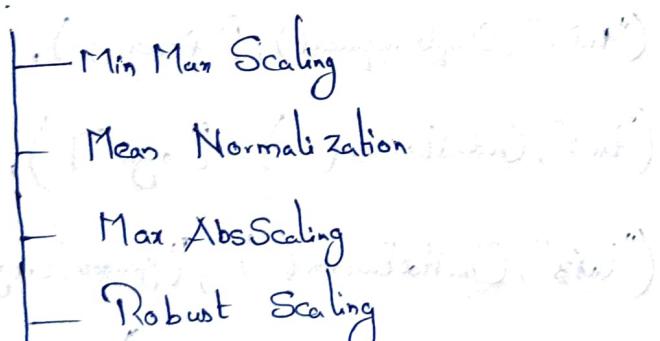
Robust Scaling

$$x_i = \frac{x_i - X_{\text{median}}}{\text{IQR}}$$

- ↳ Robust to outliers

Normalization Vs Standardization

1. Is feature Scaling important?
2. Most of the problems will be solved using standardization
- 3.



Min - Max Scaling

* mostly used for numerical data

$$x_i = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}}$$

(and they are having no outliers)

Feature Transformation

Encoding Categorical Data

Nominal Categorical Data

- ↳ No relationship b/w few values.
- (or) no hierarchy b/w data.

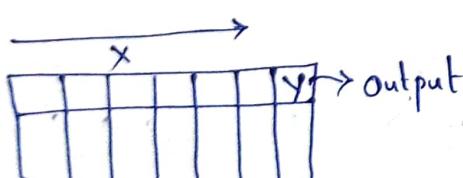
Ex:- States, Names, brand ...

Ordinal Categorical Data

- ↳ Hierarchy of the data

Ex:- Family Heirarch, Military Ranks,

Ordinal Encoding



- If the inputs are categorical data then we go for "Ordinal Encoding".

- If the Output column is categorical then we go for "Label Encoder".

↳ sklearn.preprocessing import

OrdinalEncoder

↳ For input Columns

↳ sklearn.preprocessing import LabelEncoder

↳ For Output Column

One-Hot Encoding

↳ sklearn.preprocessing import OneHotEncoder

Column Transformer

↳ from sklearn.compose import ColumnTransformer

Transformer = ColumnTransformer(transformers=[

("tf₁", SimpleImputer(), ['Pclass'])],

("tf₂", OrdinalEncoder(), ['Sex']),

("tf₃", OneHotEncoder(), ['Parch'; 'SibSp'])]

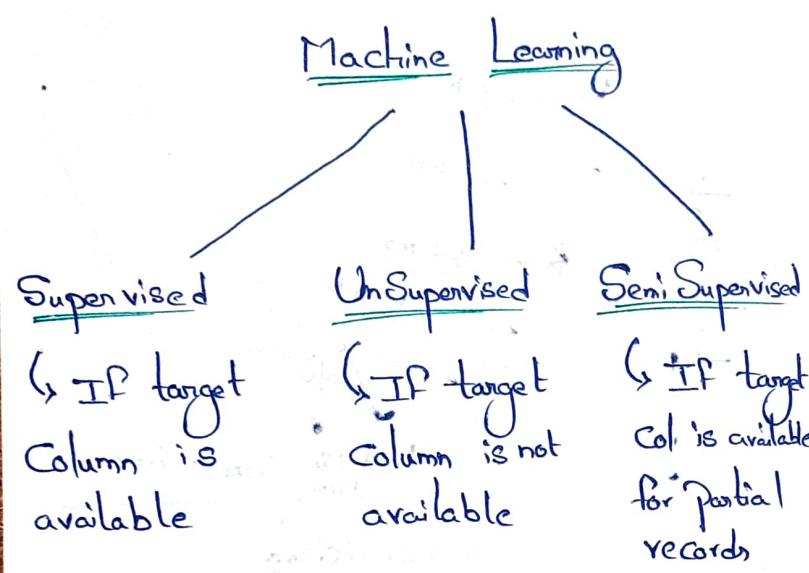
], remainder = "passthrough")

transformer.fit_transform(x_train)

transformer.transform(x_test)

①

Intro to Machine Learning



feature 1 ↓ Height (x)	feature 2 (target col) ↓ Weight (y)
5.5	57
6	64
4.7	49
5.3	61
5.0	?

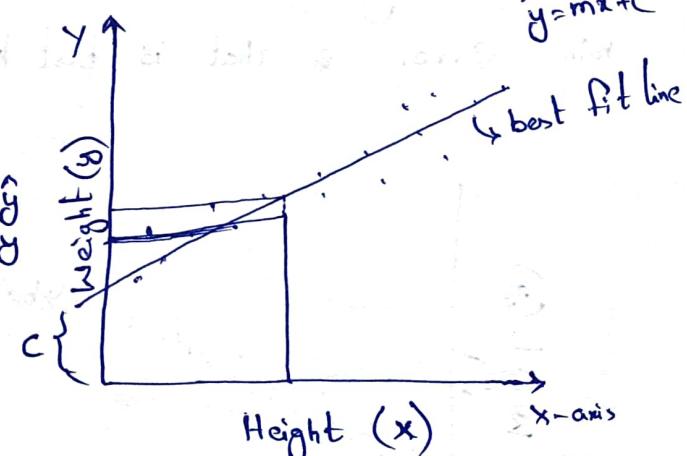
We are given some heights & their corresponding weights, we need to find the weight for a particular height based on the given data.

Supervised Learning

Classification	Regression
Target Column Has discrete values	Target Column Has Continuous Values
Binary / Multiple Values	Ex - 1.2, 1.4, 1.6, 2.3, 2.5, 2.7 ...
Ex - 0/1 (or) 1, 2, 3, 4	

How Linear Regression Model Works ?

First the model takes all the data points & plots them on a Co-ordinate System.



Then it finds a best fit line where error is minimum.

* It deals with regression based problems

Let's consider an example & understand how it works.

Gradient Descent Algorithm

$$\text{Error} = \sum_{i=1}^n |\hat{y} - y|$$

↳ Absolute Error

\hat{y} → Predicted Value

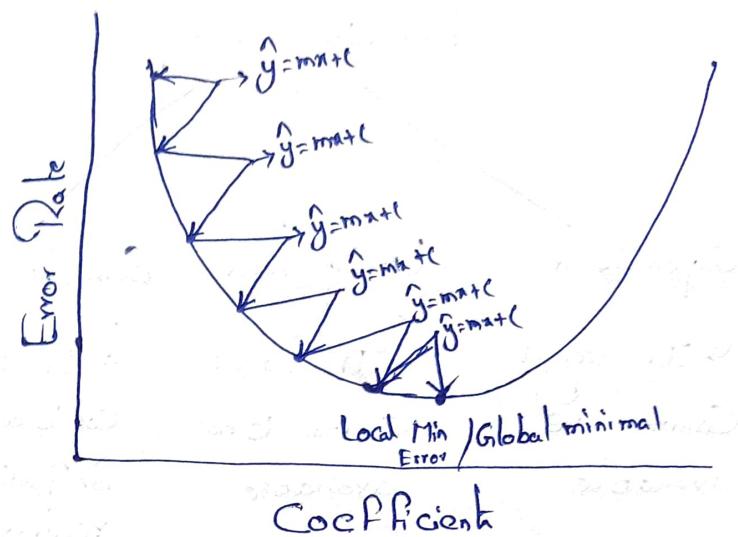
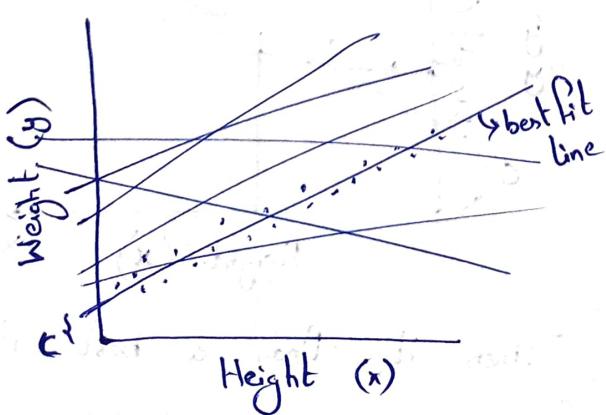
y → Actual Value

$$\hat{y} = mx + c \rightarrow \text{Straight line eq}$$

↓ ↓
slope intercept

How is it going to find best fit line?

When we start training the model it randomly draws lines & checks if the error is close to zero. Then it keeps doing it till it finds the min error & that is best fit line.



This algorithm is used to reduce the error.

At first a random initialization of m & c is done & the error is checked. If the error is not close to zero then m & c are updated and the process is repeated till the global minima is reached & the m & c values are noted.

(2)

Updates of m & c

↳ Convergence Algorithm

How is new m & c values evaluated for every iteration

Let's Consider MSE

$$\text{Error} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2$$

$$\Sigma \rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - (mx + c))^2$$

$$m_{\text{new}} = m_{\text{old}} - \eta \times \frac{\delta E}{\delta m} \quad \begin{matrix} \text{rate of change} \\ \text{of error w.r.t} \\ \text{change in slope} \end{matrix}$$

$$c_{\text{new}} = c_{\text{old}} - \eta \times \frac{\delta E}{\delta c} \quad \begin{matrix} \text{rate of change of} \\ \text{error w.r.t change} \\ \text{in intercept} \end{matrix}$$

$$\Sigma \rightarrow \frac{1}{n} \sum_{i=1}^n ((y_i)^2 + (mx + c)^2 - 2y_i(mx + c))$$

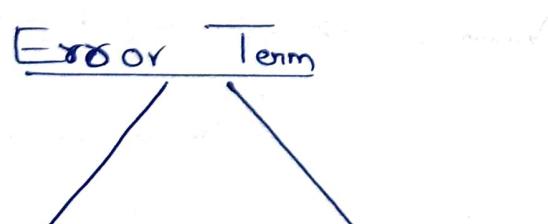
$$\frac{\delta \Sigma}{\delta m} \Rightarrow \frac{1}{n} \sum_{i=1}^n 2x(mx + c - y_i)$$

$$\Sigma \rightarrow \frac{1}{n} \sum_{i=1}^n 2x(\hat{y} - y_i)$$

$$\frac{\delta \Sigma}{\delta c} \Rightarrow \frac{2}{n} \sum_{i=1}^n (\hat{y} - y_i) \cdot n$$

$$\frac{\delta \Sigma}{\delta c} \Rightarrow \frac{2}{n} \sum_{i=1}^n (\hat{y} - y_i)$$

$\left. \begin{matrix} \frac{\delta \Sigma}{\delta m} \\ \frac{\delta \Sigma}{\delta c} \end{matrix} \right\} \rightarrow$ Extent to which the how much reduction happened



Mean Squared Error

Mean Absolute Error

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (\hat{y} - y_i)^2$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n |\hat{y} - y_i|$$

Through Hyperparameter Tuning we can adjust the learning rate

Usually Learning Rate would be

$$\eta = 0.01$$

(or)

$$0.001$$

(or)

$$1$$

Simple Linear Regression

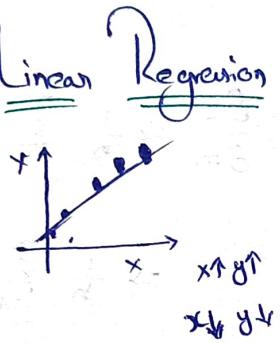
$$\hookrightarrow \hat{y} = mx + c$$

Multiple Linear Regression

$$\hookrightarrow \hat{y} = m_1x_1 + m_2x_2 + m_3x_3 + \dots + m_nx_n + c$$

Assumptions of Linear Regression

1.) Linearity:



$$(x_1, x_2, x_3, \dots, x_n) \rightarrow y$$

\hookrightarrow The inputs are in linear to the output

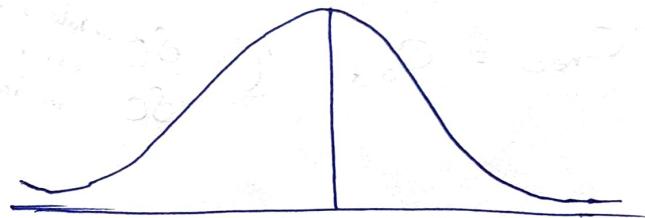
2.) Independence

$$(x_1, x_2, x_3, x_4, \dots, x_n)$$

\hookrightarrow Each feature is independent of each other

3.) Normality

If we can plot all the residual values on the graph then they follow normal / Gaussian Distribution



Normal / Gaussian Distribution

mean = median = mode

4.) Homoscedasticity

The Variance of the residuals should be constant.

Variance \rightarrow spread / dispersion

\hookrightarrow should be constant

Normalization

Linear Regression - OLS

↳ Scikit Learn uses this by default.

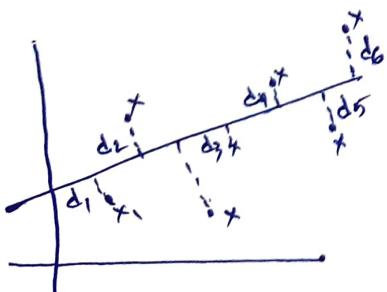
↳ Preferred when there are fewer dimensions in data.

$$Y = mx + b$$

$$\Rightarrow b = \bar{Y} - m\bar{x}$$

$$m = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\begin{matrix} \bar{x} \\ \bar{y} \end{matrix} \rightarrow \left. \right\} \text{Mean}$$



$$E = d_1 + d_2 + d_3 + d_4 + \dots + d_n$$

$$\hookrightarrow E = d_1^2 + d_2^2 + d_3^2 + d_4^2 + \dots + d_n^2$$

$$E = \sum_{i=1}^n d_i^2 \rightarrow \text{Error Function}$$

why are we squaring? why not mod?

Squaring

Reason 1 → Converts negative distance to positive

Reason 2 → We are penalising Outliers

Reason 3 → It is differentiable.

Mod

↳ It is not differentiable

$$d_i = (y_i - \hat{y}_i)$$

$$\Rightarrow E = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\hat{y}_i = mx + b$$

$$E(m, b) = \sum_{i=1}^n (y_i - mx_i - b)^2$$

$$\frac{\delta E}{\delta b} = \frac{\delta}{\delta b} \sum_{i=1}^n (y_i - mx_i - b)^2 = 0$$

$$= \sum \frac{\delta}{\delta b} (y_i - mx_i - b)^2 = 0$$

$$\Rightarrow \sum (y_i - mx_i - b) = 0$$

$$\Rightarrow \frac{\sum y_i}{n} - \frac{\sum mx_i}{n} - \frac{\sum b}{n} = 0$$

$$\downarrow \frac{\bar{y}}{n} - m\bar{x} - \frac{\bar{b}}{n} = 0$$

$$\Rightarrow \bar{y} - m\bar{x} - b = 0$$

$$\Rightarrow b = \bar{y} - m\bar{x}$$

$$E = \sum (y_i - mx_i - \bar{y} + m\bar{x})^2 = 0$$

$$\frac{\delta E}{\delta m} = \sum \frac{\delta}{\delta m} (y_i - mx_i - \bar{y} + m\bar{x})^2 = 0$$

$$\Rightarrow \sum 2(y_i - mx_i - \bar{y} + m\bar{x})(-x_i + \bar{x}) = 0$$

$$\Rightarrow \sum -2(y_i - mx_i - \bar{y} + m\bar{x})(\hat{x}_i - \bar{x}) = 0$$

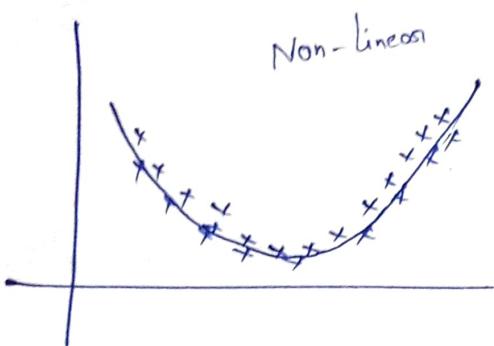
$$\Rightarrow \sum [(y_i - \bar{y}) - m(x_i - \bar{x})](x_i - \bar{x}) = 0$$

$$\Rightarrow \sum [(y_i - \bar{y})(x_i - \bar{x})] - m(x_i - \bar{x})^2 = 0$$

$$\Rightarrow \sum (y_i - \bar{y})(x_i - \bar{x}) = m(x_i - \bar{x})^2$$

$$\Rightarrow m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Polynomial Linear Regression



Simple PR

degree = 0

Polynomial degree = 0 $\hat{y} = c \cdot x^0 \rightarrow$ Constant value

" = 1 $\Rightarrow \hat{y} = c \cdot x^0 + m \cdot x^1 \rightarrow$ Simple
Linear
Regression

" = 2 $\Rightarrow \hat{y} = c \cdot x^0 + m_1 \cdot x^1 + m_2 \cdot x^2$

" = n $\Rightarrow \hat{y} = c \cdot x^0 + m_1 \cdot x^1 + m_2 \cdot x^2 + \dots + m_n \cdot x^n$

Exploratory Data Analysis

Steps

(i) Load the dataset

`df = pd.read_csv(file path)`

`df.head()` → Top 5 features

`df.shape` → no. of rows & columns

(ii) `df.info()` check for null values

`df.isnull().sum()` → Shows null value for each feature

null value → missing value / NAN

Replacing null values is called as

Data Imputation

Measures of Central Tendency

KNN Algorithm

(iii) Check for Distribution

Symmetric

Skewed

Mean

Median

da DF.describe()

↳ gives the descriptive stats of the data

(iv) Check if it is in SNF

SNF - Standard Normal Form

↳ $\mu = 0$ → Mean

$\sigma = 1$ → Standard Deviation

Why Convert to SNF?

↳ To make sure the model is not biased to any one feature.

Range → -1 to 1

↳ when we scale all the features

will be in this range

(iv) Finding Correlations

* Convert Object → Categorical data type

(v) One Hot Encoding

↳ Convert Categorical data to

numerical data (Ex: 0, 1, ...)

Feature	Red	Green
Red	1	0
Green	0	1
Yellow	0	0
Green	0	1

dataframe = pd.get_dummies(df)

dataframe.head(30)

dataframe = pd.get_dummies(df, drop_first=True)

↳ Less no. of
Features (viii)

fit_transform

↳ 1.

helps to get the value of

$\mu \approx 6$

Model Preparation (Training & Testing)

(vi) Co-efficient Correlation

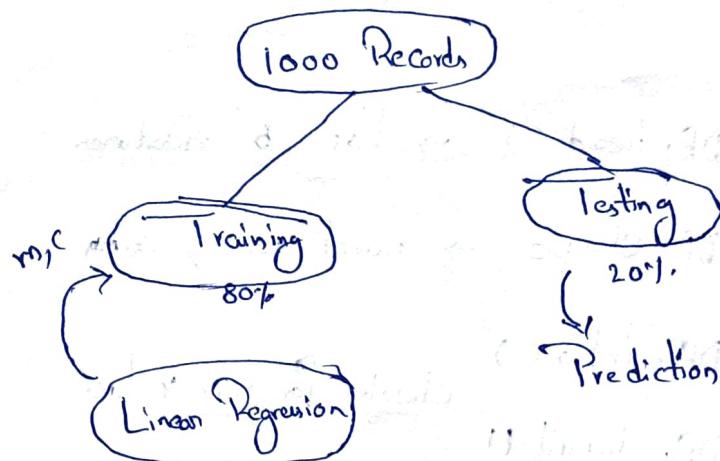
dataframe.Corr()

(vii) Scaling of Numeric Features

from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()

dataframe[...] = dataframe.fit_transform()



Model Splitting

from sklearn.model_selection import train_test_split

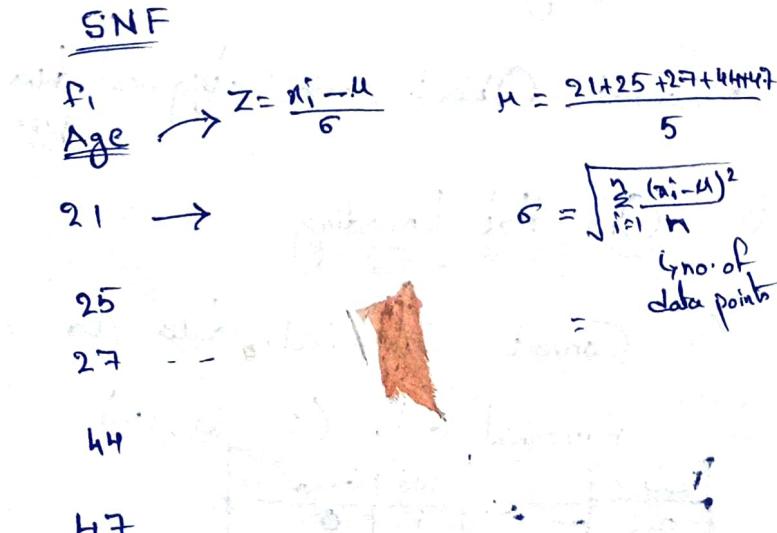
Model Training

from sklearn import LinearRegression

model = LinearRegression()

model.fit(x_train, y_train)

y_pred = model.predict(x_test)



Model Evaluation

$$\hookrightarrow \text{(i) } \text{MSE} = \frac{1}{n} (\hat{y} - y)^2$$

$$\text{(ii) } \text{RMSE} = \sqrt{\text{MSE}} \rightarrow \downarrow \downarrow$$

(iii) R^2 — closer to 1

(iv) Adjusted R^2

- * As the no. of features increase, R^2 value also keeps increasing irrespective of the feature.

So, to overcome this drawback we have adjusted R^2 .

Adjusted R^2

$$\hookrightarrow 1 - \frac{(1-R^2)(n-1)}{n-p-1}$$

n - no. of records
p - no. of features

- * As the no. of features increase, adjusted R^2 stays constant (or) decreases based on the feature.

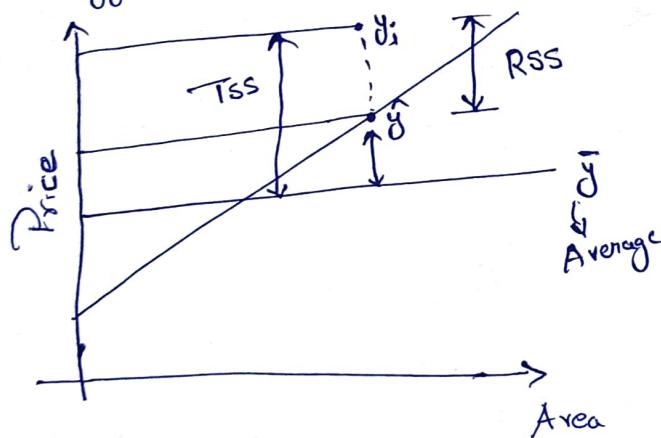
Performance Metrics

R^2

- * $1 - \frac{\text{RSS}}{\text{TSS}}$ \rightarrow Residual Sum of Square

\rightarrow Total Sum of Square

1 — $\frac{\text{Smaller number}}{\text{Bigger number}}$

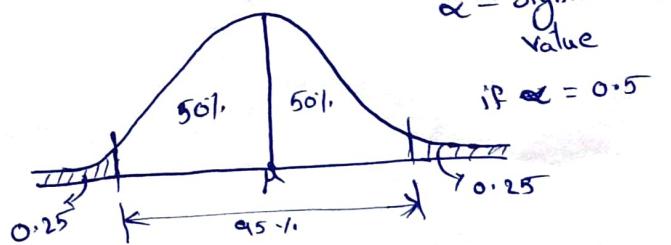


$$\text{RSS} = \sum (y_i - \hat{y})^2$$

$$\text{TSS} = \sum (y_i - \bar{y})^2$$

- * If R^2 is close to 1 then it is a good accuracy.

P-Value \rightarrow Statistics



f_1 - area $\rightarrow 0.0000$ (p-value)

f_2 - bedrooms $\rightarrow 0.0000$

f_3 - sex $\rightarrow 0.18000 \rightarrow 100 - 18 = 82$
Passing only times to failing

If 'P' value is higher then the feature is irrelevant and it should be removed.

Real Time Use Case

HR (Employee Performance Evaluation)

Finance

Z-Score

When there are not too many outliers

TQR

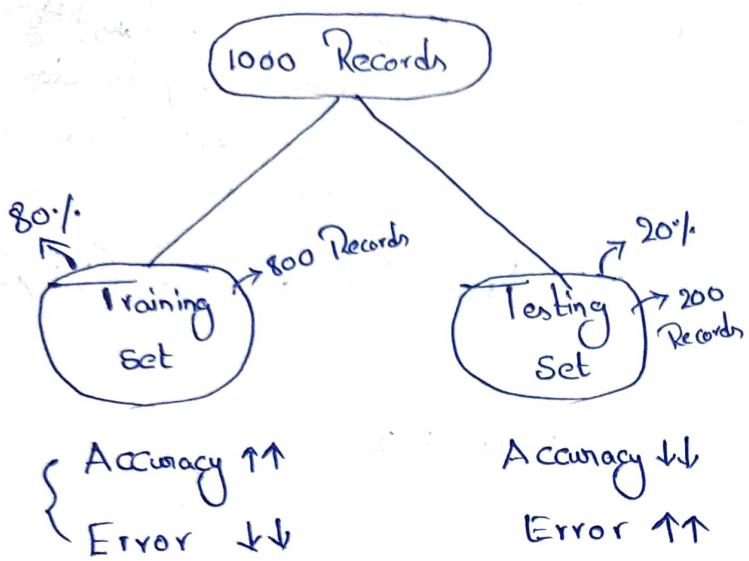
When the Outliers are too much more than 50%, then we go for this method

Ex:-

ECommerce Data

Real Estate Data

Overfitting & Underfitting



↳ Overfitting

↳ This is the main problem in ML

(Low Bias & High Variance)



Bias & Variance \Rightarrow Error

Q) Why two variables are used for error?

Data Distribution

i) During training we might have data with same distribution but during testing the data might not have same distribution.

So to represent the errors during training & testing differently we have 2 Errors.

(ii) Complex / noisy data points

During training the error might be low because of the noisy data points

(o) Outliers, this affects the testing phase.

so '2' variables to estimate error during training & testing.

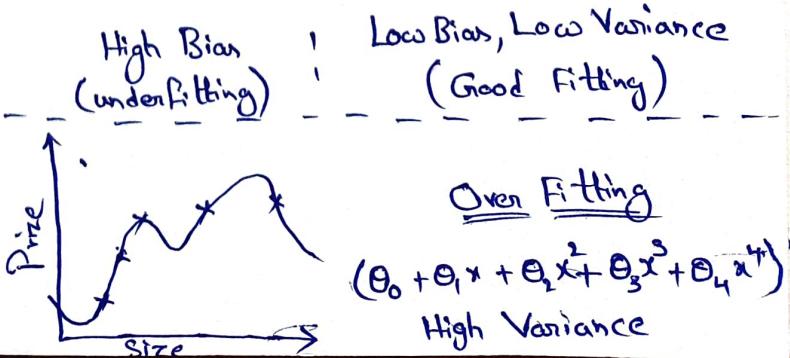
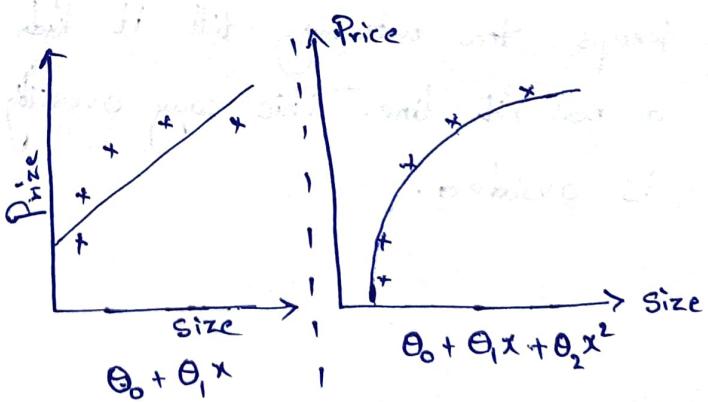
↳ Generalization happens because of the Outliers.

≡ Generalizing about the drivers of Uber/da.

Regularization Regression

↳ To avoid/handle overfitting problem

we want
↳ Low Bias & Low Variance



Regression

↳ L1 Regularization (LASSO)

↳ L2 Regularization (RIDGE)

↳ L3 Regularization (L1 + L2)

LASSO → Very widely Used
no. of features

$$\text{MSE} + \lambda \sum_{i=1}^K |w_i|$$

↳ Hyper Parameter
Error Function

Ridge → Also used for Feature Selection

$$\text{MSE} + \lambda \sum_{i=1}^K (w_i)^2$$

when slope/weightage is close to zero, that feature is removed.

Elastic Net Regularization (L1 & L2)

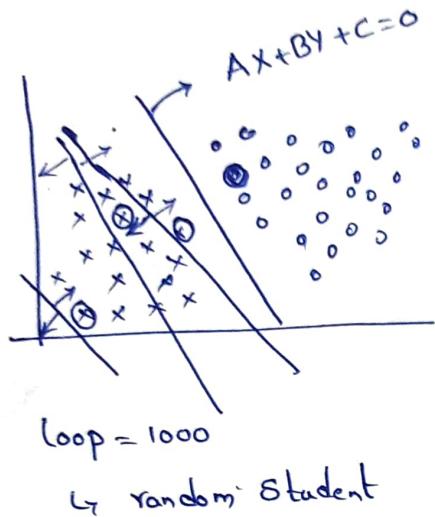
↳ This helps in Overcoming the problem of overfitting.

↳ when MSE/Error becomes zero

(i) close to zero, we are adding extra value, so that the model keeps retraining till it finds a best fit line. This way overfitting is avoided.

Logistic Regression

Perception Trick



In Logistic Regression we will replace step function with Sigmoid

$$\hat{y} = w_1 x_{i1} + w_2 x_{i2} + w_0 \xrightarrow{\downarrow z} \text{Sigmoid}$$

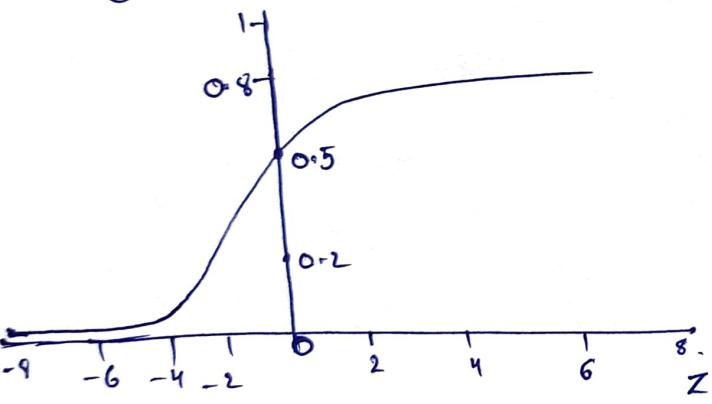
$$z \Rightarrow \sum w_i x_i \xrightarrow{\downarrow \text{sigmoid}} \sigma(z)$$

if $\underline{z} > 0$.

if z is positive $\Rightarrow \geq 0.5$

z is negative $\Rightarrow \leq 0.5$

Sigmoid Function



$-\infty < z < \infty$

$0 < y < 1$

This fn output will always be in the range of 0 to 1.

what we used to do

$$\hat{y} = w_1 x_{i1} + w_2 x_{i2} + w_0 \xrightarrow{\downarrow z} \begin{cases} z \geq 0 \rightarrow 1 \\ z < 0 \rightarrow 0 \end{cases}$$

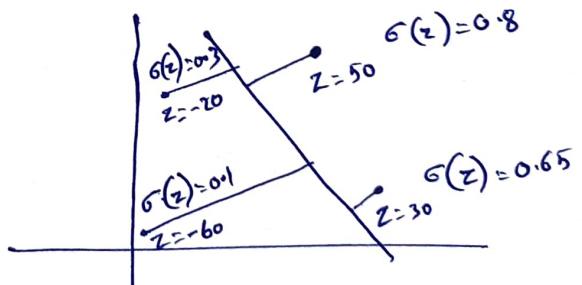
$\xrightarrow{\text{step function}}$

Impact of Sigmoid

$$w_n = w_0 + \eta(y_i - \hat{y}_i)x_i$$

$$y_i = \sigma(z)$$

where $z = \sum w_i x_i$



y_i	\hat{y}_i	$y_i - \hat{y}_i$
1	0.8	0.2
0	0.65	-0.65
1	0.3	0.7
0	0.15	-0.15

Logistic Regression

→ Deals with Binary Classification

Study Hours

Pas/Fail

Fail

2

3

4

5

6

8

12

Pas

"

Pas

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

"

Why we can't use Linear Regression

$$h_{\theta}(x) = \frac{1}{1+e^{-z}} \rightarrow \text{Sigmoid fn}$$

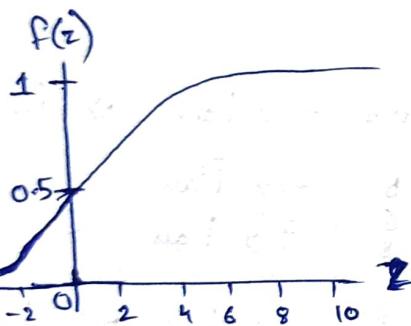
for Logistic Regression?

1.) Best fit line (outlier)

↳ changes

2.) >1 & <0 {square line}
↳ Linear Regression
Can't do it.

Sigmoid Function



if $z > 0$

$$f(z) \geq 0.5$$

else

$$f(z) < 0.5$$

$$h_{\theta}(x) = \theta_1 x + \theta_0 \rightarrow \text{Linear Regression}$$

$$h_{\theta}(x) = \sigma(\theta_1 x + \theta_0) \quad z = \theta_1 x + \theta_0$$

↳ Activation Fn (sigmoid)

$$= \sigma(z)$$

Linear Regression Cost function

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

$$h_{\theta}(x) = \theta_1 x + \theta_0$$

↳ Convex

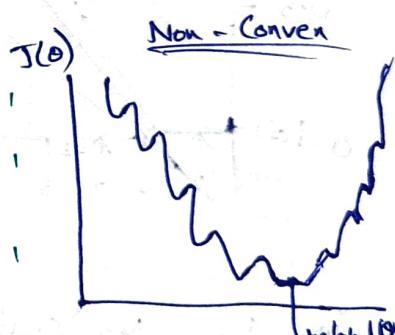
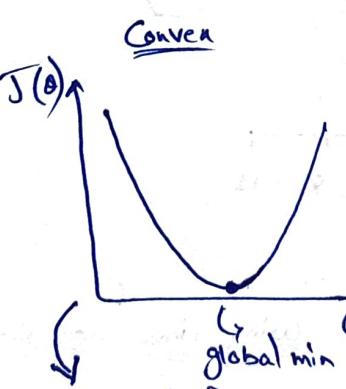
Logistic Regression Cost function

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

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

↳ Sigmoid activation fn

↳ Non Convex



Gradient Descent

↳ Linear Regression

↳ Logistic Regression

In Order to replace this non-Convex Curve with Convex Curve we are using Sigmoid / Logloss Cost function.

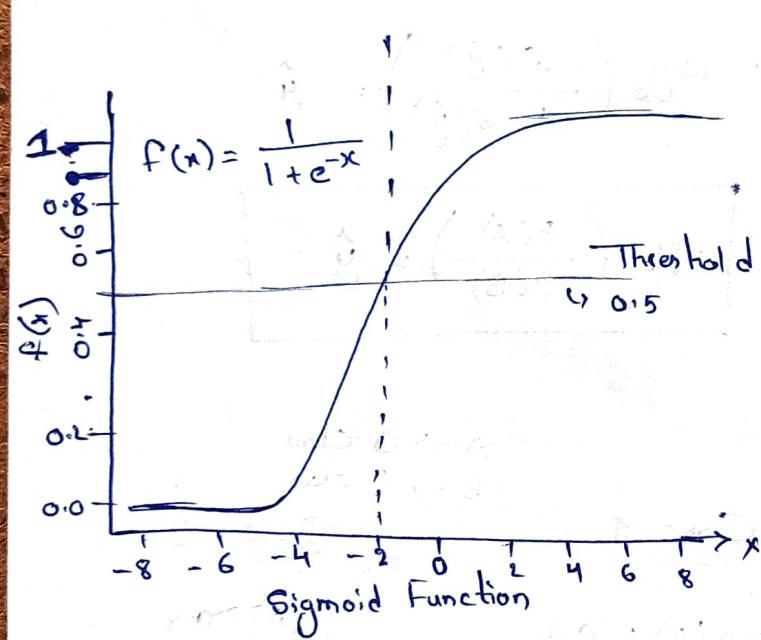
Logistic Regression

↳ Mostly used for Binary Classification

↳ It gives us the probability

Sigmoid \rightarrow Probability
 $0 \rightarrow 0.6$
 $1 \rightarrow 0.4$

Activation
function



* Sigmoid Function gives us the probability based on the threshold value

IF the threshold $> 0.5 \rightarrow 1$

threshold $< 0.5 \rightarrow 0$

In Linear Regression we want to find a straight line/best fit line

$$\hat{y} = mx + c$$

but in logistic Regression we want to find the probability of it being a straight line/best fit line

$$\therefore \sigma(\hat{y}) = \frac{1}{1 + e^{-\hat{y}}}$$

so we apply sigmoid fn on the

line $y = mx + c$

$$\therefore \sigma(\hat{y}) = \frac{1}{1 + e^{-\hat{y}}}$$

$$= \frac{1}{1 + e^{-(mx+c)}}$$

$$= \frac{1}{1 + \frac{1}{e^{mx+c}}}$$

$$= \frac{1}{e^{mx+c} + 1}$$

$$\sigma(\hat{y}) = \frac{e^{mx+c}}{e^{mx+c} + 1}$$

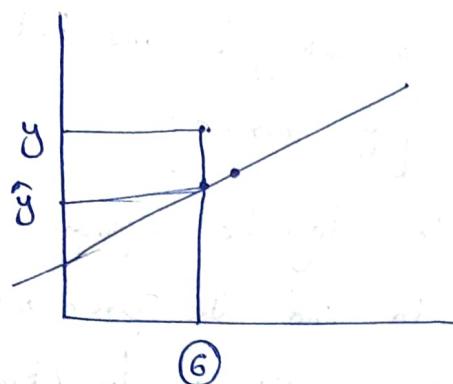
$$P(A) + P(B) = 1$$

$$P(A) = 0.6$$

$$P(B) = 0.4$$

$$\frac{\sigma(\hat{y})}{1-\sigma(\hat{y})} = e^{mx+c}$$

$$\log_e \left(\frac{\sigma(\hat{y})}{1-\sigma(\hat{y})} \right) = \log_e e^{mx+c}$$



In linear Regression we try to find the difference b/w the Predicted value & actual value. But in Logistic Regression we try to find what is the probability of height weight being 58 for the given height 6.

$$\Rightarrow \log_e \left(\frac{\sigma(\hat{y})}{1-\sigma(\hat{y})} \right) = mx+c$$

$$\Rightarrow \log_e \left(\frac{\sigma(\hat{y})}{1-\sigma(\hat{y})} \right) = \hat{y}$$

$$\boxed{\log_e \left(\frac{P(A)}{P(B)} \right) = \hat{y}}$$

$$A \rightarrow 0 \text{ class}$$

$$B \rightarrow 1 \text{ class}$$

Case 1

$$P(A) = 0.2$$

$$P(B) = 0.8$$

Case 2

$$P(A) = 0.8$$

$$P(B) = 0.2$$

Sigmoid fn

$$\hat{y} = \log_e \left(\frac{0.2}{0.8} \right)$$

Sigmoid fn

$$\hat{y} = \log_e \left(\frac{0.8}{0.2} \right)$$

$$= -1.38 \quad \text{class 0} \quad = 1.38 \quad \text{class 1}$$

0 class \downarrow class 0

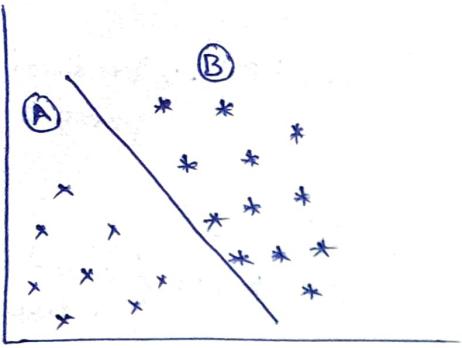
0 class

1 class \downarrow

$$\boxed{P(B) = \frac{1}{e^{mx+c} + 1}}$$

Ans (3)

Ans (3)



<u>Proper Classification</u>	<u>Misclassification</u>
$\hat{y} = 1, y = 0$	$\hat{y} = 0, y = 0$
$= -\log_e(1-\hat{y})$	$= -\log_e(1-0)$
$= -\log_e(1-1)$	$= -\log_e(1)$
$= \infty$	$= 0$
\hookrightarrow Error is very high	\hookrightarrow Error is zero

Cost Function

Linear Regression

$$\hookrightarrow \text{Cost function} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

Logistic Regression

$$\hookrightarrow \text{Cost function} = \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1+e^{(m_i + b)}} - y_i \right)^2$$

To get a smooth function of J

Convex Curve

$$\text{Cost}(\hat{y}, y) = \begin{cases} -\log_e \hat{y} & ; y = 1 \\ -\log_e (1-\hat{y}) & ; y = 0 \end{cases}$$

\hookrightarrow log loss

Misclassification \Rightarrow case 1

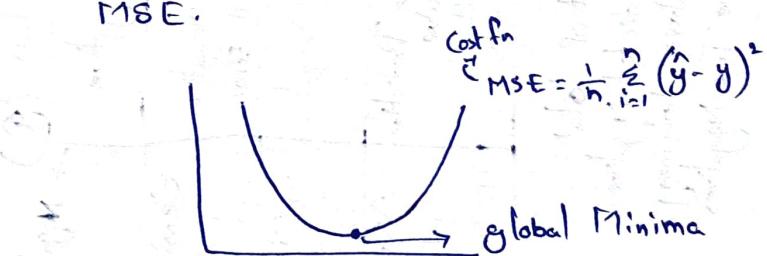
<u>Actual</u>	<u>Predicted</u>	<u>Proper Classification</u>
$\hat{y} = 0$	$y = 1$	$\hat{y} = 1, y = 1$
$\hat{y} = 1$	$y = 0$	$= -\log_e 1$
$\hat{y} = 0$	$y = 0$	$= 0$
$\hat{y} = 1$	$y = 0$	$\Rightarrow 0$
$\Rightarrow \infty$		

- * When there is a misclassification the error is very high & the training keeps going till the error is min.

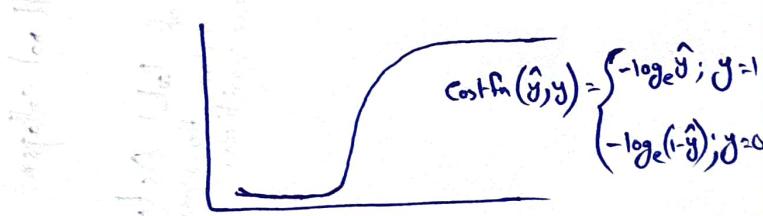
Why different Cost function?

In Linear Regression we deal with

Convex Curve, so the Cost fn is MSE.



But in Logistic Regression we don't deal with Convex Curve. So we need a different Cost fn. (sigmoid fn)



$$\text{Cost Fn} = \sum_{i=1}^N -y_i \log(\hat{y}) - (1-y_i) \log(1-\hat{y})$$

↳ Another way of writing Cost fn

$$\text{Precision} = \frac{TP}{(TP+FP)} \rightarrow \text{Predicted Positive}$$

$$\text{Recall} = \frac{TP}{(TP+FN)} \rightarrow \text{Actual Positive}$$

↳ Very serious problem ↴ ↴

f₁-Score

$$\hookrightarrow \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

Q1 When Precision is important & when Recall is important ?

↳ Recall should be used when we want to reduce False Negative.

Precision should be used when we want to reduce False Positive.

		Building will collapse		Building will not collapse	
		Hypothesis		Negating hypothesis	
		Predicted class (\hat{y})	Predicted class (\hat{y})	Building will not collapse	
Building will collapse (y)	Positive	True Positive (TP)	False Positive Type I Error	Sensitivity $\frac{TP}{TP+FN}$	Out of total Actual Positive values how many are True Positives.
	Negative	False Negative Type II Error	True Negative	Specificity $\frac{TN}{TN+FP}$	
Precision		Negative Predictive Value		Accuracy $\frac{TP+TN}{TP+TN+FP+FN}$	
		Out of total Predicted Positive values, how many are True Positives			

Accuracy

↳ It is not a good measure to consider when there is imbalance in the data.

$$\begin{array}{l} 0 \rightarrow 700 \text{ records} \\ 1 \rightarrow 300 \text{ records} \end{array}$$

(↳ 30%)

↳ Biased towards zero data.

F-Beta Score

$$\hookrightarrow (1+\beta^2) \frac{\text{Precision} \cdot \text{Recall}}{(\beta^2 \cdot \text{Precision}) + \text{Recall}}$$

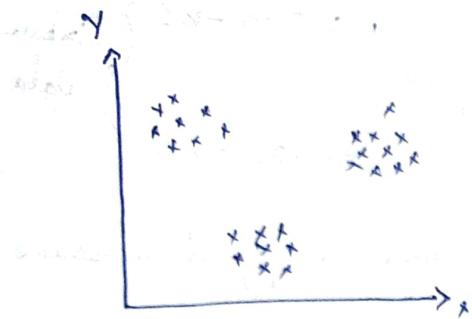
$\beta > 1$: Recall is more important than Precision (e.g., 2)

$\beta < 1$: Precision is more important than Recall (e.g., 0.5)

$\beta = 1$: If both Precision & Recall are important (F₁ Score)

$$\begin{aligned} F_1 \text{ Score} &= (1+1)^2 \frac{\text{Precision} \cdot \text{Recall}}{(1)^2 \text{Precision} + \text{Recall}} \\ &= 2 \cdot \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \end{aligned}$$

Logistic Regression → Multiclass

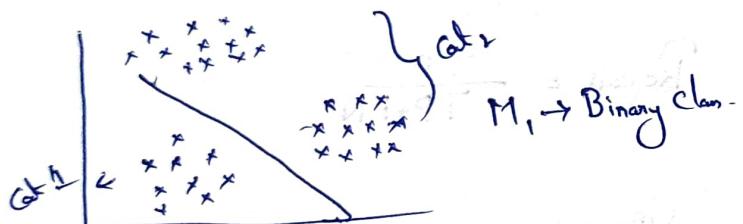


	f_1	f_2	f_3	O/P	O_1	O_2	O_3	One Hot Encoding
-	-	-	-	O_1	1	0	0	
-	-	-	-	O_2	0	1	0	
-	-	-	-	O_3	0	0	1	
-	-	-	-	O_1	1	0	0	
-	-	-	-	O_3	0	0	1	
-	-	-	-	O_2	0	1	0	

One Versus Rest

f_1	f_2	O/P
-	-	O_1
-	-	O_2
-	-	O_3
-	-	O_2
-	-	O_3

→ All are equal & value addition
→ value addition
→ value addition



$$M_1 \rightarrow I/P \{f_1, f_2, f_3\}, O/P \{O_1\}$$

$$M_2 \rightarrow I/P \{f_1, f_2, f_3\}, O/P \{O_2\}$$

$$M_3 \rightarrow I/P \{f_1, f_2, f_3\}, O/P \{O_3\}$$

New Test Data $\rightarrow M_1 \rightarrow 0.25 \rightarrow$ Probability

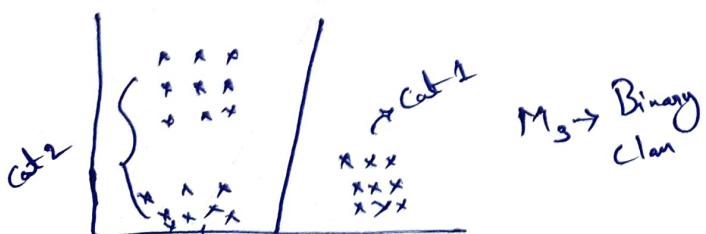
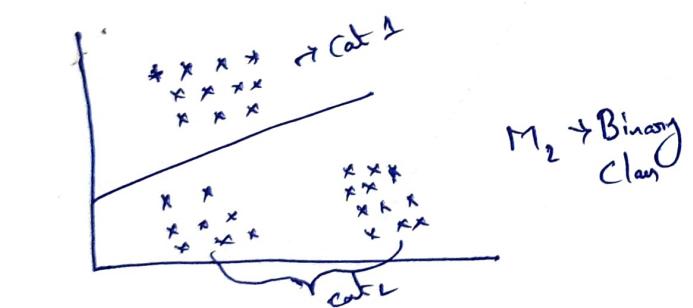
$$\hookrightarrow M_2 \rightarrow 0.20$$

$$\hookrightarrow M_3 \rightarrow 0.55$$

$$[0.25, 0.20, 0.55]$$

$\uparrow \quad \uparrow \quad \uparrow$
 $M_1 \quad M_2 \quad M_3$

\therefore The output for the new test data is O_3 , since M_3 used O_3 as output.



Decision Trees

↳ Can be used for both Regression & Classification based problems

Impurity Measure:-

* Model tries to predict which feature has lesser impurity

↳ Lesser impurity (or) Higher Information Gain

① Gini Index

② Entropy

Gini Index

↳ Using this we can decide which node will be a root node/parent node

$f_1 \quad f_2$

Gender	Class	Stay in Hostel
M	9	Yes
F	10	No
F	8	Yes
F	8	No
M	9	Yes
M	10	No
F	11	Yes
M	11	Yes
F	8	Yes
M	9	No
M	11	No
M	11	Yes
F	10	No
M	10	Yes

$$\text{Gini} \Rightarrow 1 - \sum_{i=1}^c p_i^2$$

Gini (Class)

Class	Stays in Hostel	Probability	
		P(Y)	P(N)
8	Y=2 N=1	2/3	1/3
9	Y=2 N=1	2/3	1/3
10	Y=1 N=3	1/4	3/4
11	Y=3 N=1	3/4	1/4

$$\text{Gini}(\text{class}) = 1 - (2/3)^2 - (1/3)^2 = 4/9$$

$$\text{Gini}(\text{class}=9) = 1 - (2/3)^2 - (1/3)^2 = 4/9$$

$$\text{Gini}(\text{class}=10) = 1 - (1/4)^2 - (3/4)^2 = 3/8$$

$$\text{Gini}(\text{class}=11) = 1 - (3/4)^2 - (1/4)^2 = 3/8$$

Weighted Average for Gini (class)

$$\begin{aligned}
 &= \frac{\text{No. of records}}{\text{Total no. of records}} * \frac{\text{Gini of class}}{\text{Gini of class}} \\
 &= \frac{3/4}{14} * \frac{4}{9} + \frac{3/4}{14} * \frac{4}{9} + \frac{4/14}{14} * \frac{3}{8} \\
 &\quad + \frac{4/14}{14} * \frac{3}{8} \\
 &= 0.404
 \end{aligned}$$

Gender	Stay in Hostel	P(Y)	P(N)
M	Y=5 N=3	5/8	3/8
F	Y=3 N=3	3/6	3/6

$$\text{Gini}(M) = 1 - (5/8)^2 - (3/8)^2 = \frac{15}{32}$$

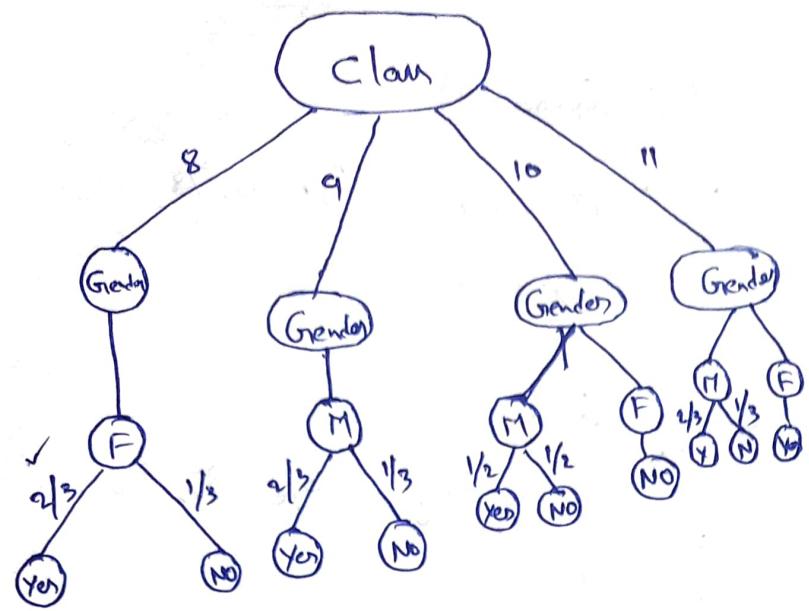
$$\text{Gini}(F) = 1 - (3/6)^2 - (3/6)^2 = 0.5$$

Weighted Average for Gini (Gender)

$$\begin{aligned}
 \text{Gini}(\text{Gender}) &= \frac{8}{14} * \frac{15}{32} + \frac{6}{14} * 0.5 \\
 &= 0.482
 \end{aligned}$$

$$\text{Gini}(\text{class}) < \text{Gini}(\text{Gender})$$

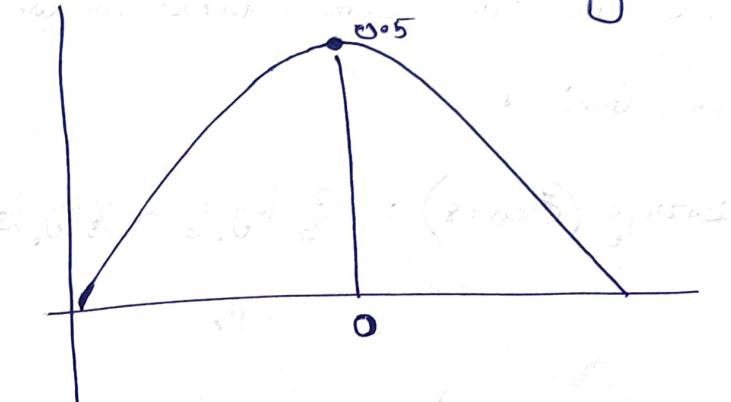
$\therefore \text{Gini}(\text{class})$ will be the parent node.



Pure Sample

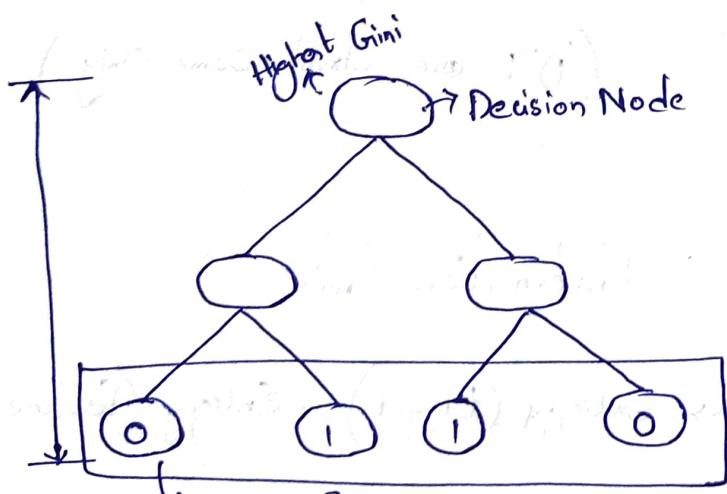
$$\left\{ \begin{array}{l} Y = 2 \\ N = 2 \end{array} \right. \Rightarrow 1 - \left(\frac{1}{2} \right)^2 - \left(\frac{1}{2} \right)^2 = 0.5$$

Gini Value $\rightarrow \{0 \text{ to } 0.5\}$
 ↓
 Lowest Highest



Pure Sample $\rightarrow \left\{ \begin{array}{l} Y = 3 \\ N = 0 \end{array} \right. \Rightarrow 1 - \sum_{i=1}^3 P_i^2 = 1 - \left(\frac{1}{3} \right)^2 = 0$

$$\left\{ \begin{array}{l} Y = 0 \\ N = 3 \end{array} \right. \Rightarrow 1 - \left(\frac{3}{3} \right)^2 = 0$$



As we are approaching
the lower levels the
Gini Value decreased.

Entropy

It is a measure of Randomness

$$\text{Entropy} \Rightarrow - \sum_{i=1}^c P_i \log_2 P_i$$

Computation is very expensive

We use the same table we used for Gini:

$$\text{Entropy} (\text{class}=8) = -\frac{2}{8} \log_2 \frac{2}{8} - \frac{1}{3} \log_2 \frac{1}{3}$$

$$= 0.918$$

$$\text{Entropy} (\text{class}=9) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3}$$

$$= 0.918$$

$$\text{Entropy} (\text{class}=10) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4}$$

$$= 0.811$$

$$\text{Entropy} (\text{class}=11) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4}$$

$$= 0.811$$

Entropy (class)

$$\Rightarrow \frac{3}{14} * 0.914 + \frac{3}{14} * 0.918 + \frac{4}{14} * 0.811$$

$$+ \frac{4}{14} * 0.811$$

$$\Rightarrow 0.864$$

Entropy (Grenden=F)

$$= -\frac{5}{14} \log_2 \frac{5}{14} - \frac{3}{14} \log_2 \frac{3}{14}$$

$$= 0.954$$

Entropy (Grenden)

$$= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$= 1$$

Entropy (Grenden)

$$\Rightarrow \frac{8}{14} * 0.954 + \frac{6}{14} * 1$$

$$\Rightarrow 0.974$$

Gini / Entropy

↳ Results are ^{always} similar

(1) T are almost same Only

Information Gain

↳ Entropy (target) - Entropy (feature)

$$\text{Entropy} (\text{target}) = -P_Y \log_2 P_Y - P_N \log_2 P_N$$

$$= -\frac{8}{14} \log_2 \frac{8}{14} - \frac{6}{14} \log_2 \frac{6}{14}$$

$$= 0.98522$$

$$\begin{aligned}\text{Information Gain of Class} &\Rightarrow 0.9852 - 0.856 \\ &= 0.1278\end{aligned}$$

$$\begin{aligned}\text{Information Gain of Gender} &\Rightarrow 0.9852 - 0.974 \\ &= 0.011\end{aligned}$$

\therefore Class is providing more information in comparison to gender.

Naive Bayes

↳ This algorithm works on the principle of Bayes' theorem

↓
Conditional Probability

→ Probability of A given B

$$P(A|B) = \frac{P(AnB)}{P(B)}$$

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

$$P(B|A) = \frac{P(BnA)}{P(A)}$$

$$P(AnB) = P(B|A) * P(A)$$

Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	False	No
Overcast	Hot	High	False	Yes
Rainy	Hot	High	True	No
Sunny	Mild	High	False	Yes
Sunny	Cold	Normal	False	Yes
" "	"	"	True	No Yes
Overcast	"	"	"	Yes
Rainy	"	"	False	Yes
Rainy	Mild	High	False	No
Sunny	"	Normal	"	Yes
Rainy	"	"	True	Yes
Overcast	"	High	True	Yes
" "	Hot	Normal	False	Yes
Sunny	Mild	High	True	No

$P(Y | \text{outlook})$

$P(Y | \text{sunny, Hot, Normal, False})$

=

$P(N | \text{Sunny, Hot, Normal, False})$

=

Features
P(B|A) * P(A)

$$P(\text{Sunny} | \text{Yes}) * P(\text{hot} | \text{Yes}) * P(\text{Normal} | \text{Yes}) * P(\text{False} | \text{Yes}) * P(\text{Yes})$$

$$P(\text{Sunny}) * P(\text{hot}) * P(\text{Normal})$$

$$* P(\text{False})$$

$$\hookrightarrow P(B)$$

$$P(\text{Sunny} | \text{No}) * P(\text{hot} | \text{No}) * P(\text{Normal} | \text{No})$$

$$* P(\text{False} | \text{No}) * P(\text{No})$$

$$P(\text{Sunny}) * P(\text{hot}) * P(\text{Normal})$$

$$* P(\text{False})$$

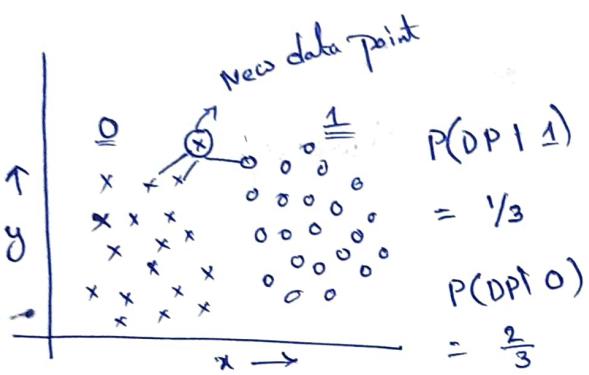
$$\therefore P(Y | \text{sunny, Hot, Normal, False}) = 0.6$$

$$P(N | \text{Sunny, Hot, Normal, False}) = 0.4$$

$$\therefore \text{Target} = \text{Yes}$$

KNN Algorithm

↳ Works for both Regression & classifications



Final prediction = 0 class

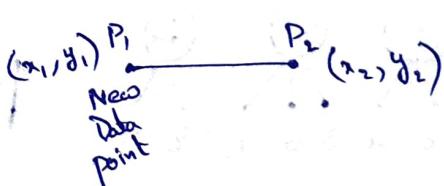
Classification Problem

$$\text{Regression} \rightarrow \arg \left(\frac{1cr + 2cr + 3cr}{3} \right) = 2cr$$

Distance Metric

$k \rightarrow$ Odd number

↳ To avoid bias towards one class

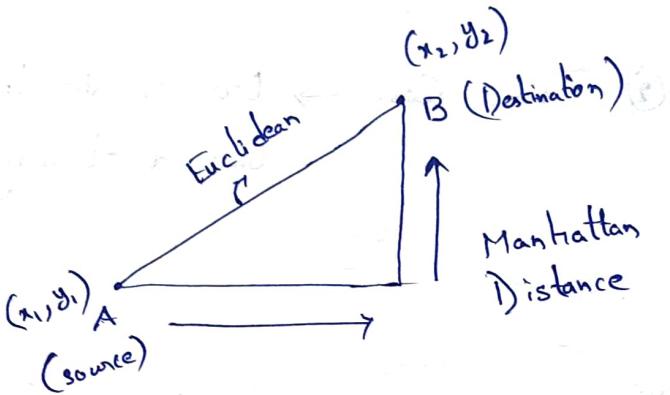


① Euclidean Distance $\rightarrow P=2$ (Default)

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

② Manhattan Distance

$$|x_2 - x_1| + |y_2 - y_1|$$



Categorical (or)
Binary

Hamming Distance

"abc" \rightarrow 1 (distance)
"dbc" mismatch

$$\begin{array}{r} \text{Binary} \rightarrow 1 0 0 1 0 0 1 \\ 0 1 0 1 0 1 0 \\ \hline 1 + 1 + 0 + 0 + 0 + 1 + 1 \end{array}$$

$$\text{Distance} = 4$$

Pros & Cons of KNN

Pros

- ① KNN doesn't make any assumptions about the distribution of data.
↳ It works for any distribution of data.
- ② Time Cost → Less at the time of training

Con's

- It does all the calculations during prediction.
- ① It is a Lazy Learner
↳ At time of Prediction it will take too much time.
 - ② Finding optimal value of 'K'
- * When we have huge amount of data, KNN should be avoided.

Applications

- ↳ Smaller data Points
- ↳ Data Imputation → null values
- ↳ SMOTE → OverSampling
↳ uses KNN Algorithm

Optimization

↳ Reduce Time & Space Complexity

KD Tree

Binary Search Tree

Ball Tree

Clustering

