

feature engineering

↳ feature transformation.

↳ { match - smooth
mathematical transformation }

1). dog - transform.

2). Reciprocal

3). power (sq / sqrt)

1). Box - Cox.

2). Yeo - Johnson

↳ { normal distribution }

{ to convert in normal distribution }

function transformer:-

sk - learn

use

↳ sns.distplot,
function
transform

pd → pd.skew()
(-, +) skewed
data

Power
transform

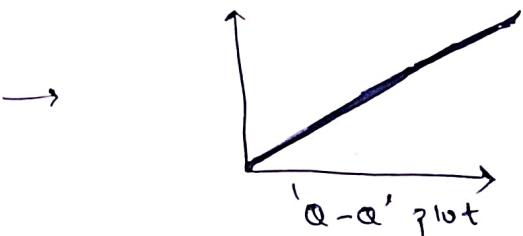
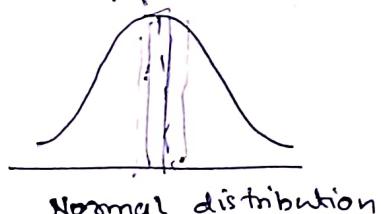
Quantile
transform.

+ log
+ Reciprocal
+ sq / sqrt.
+ custom.

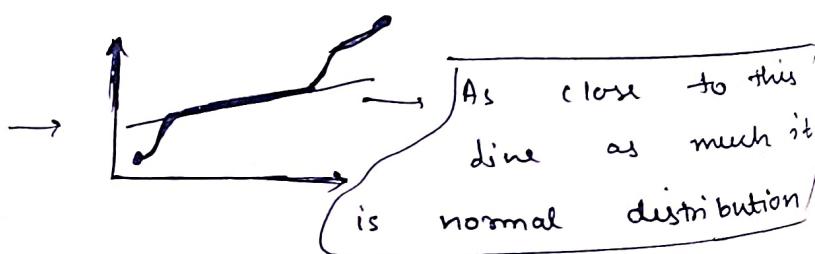
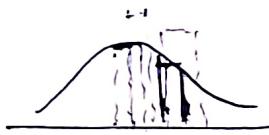
+ Box - Cox
+ yeo - johnson.

'Q-Q plot'

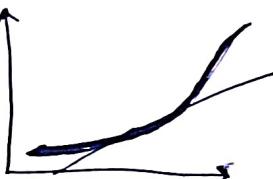
↳
stats
scipy
use.

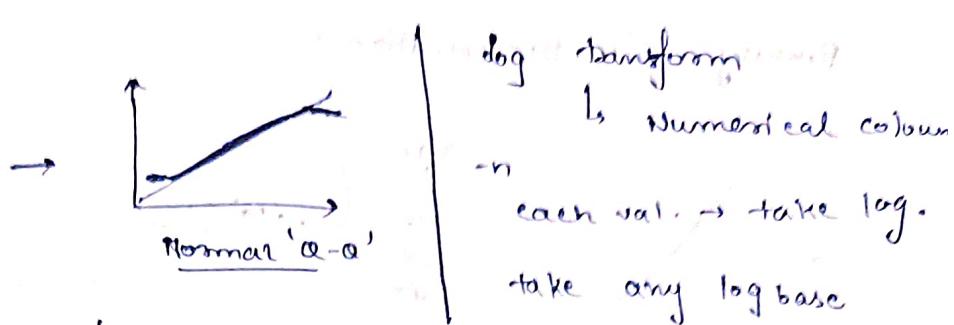
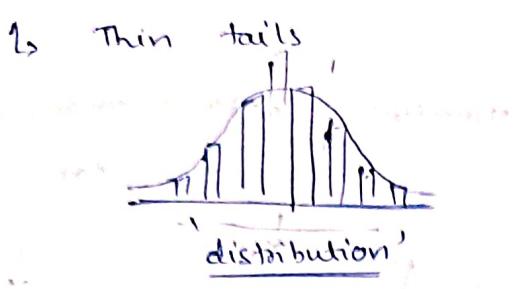


↳ Data peaked in middle.



↳ skewed - data.





- ⇒ dog transform used on 'Right-skewed' data to normalize.
- get equivalent in the scale.

Reciprocal $\Rightarrow (1/n)$	\sqrt{n}	Square. (n^2) \Rightarrow 'left' skewed data'
Image \leftrightarrow Text. audio \leftrightarrow Text	Mild skew	Make val non-negative, highlight outliers

Powers transformers :-

- ⇒ Box-Cox transformer :-

given distrib. $u_i^{(\lambda)} = \begin{cases} \frac{u_i^{(\lambda)} - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \ln(u_i) & \text{if } \lambda = 0 \end{cases}$

(convert 'Normal distribution') {try to calculate ' λ ' val.} range $[-5 \leftrightarrow +5]$

only applicable on
condition $n \geq 0$

$$-5 < \lambda < 5$$

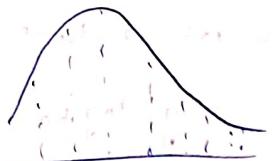
- ⇒ Yeo-Johnson :-

Applicable on every value, we will take each λ , then raise it with the power of ' n '.

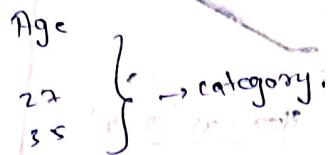
$$(Applicable on) u_i^{(\lambda)} = \begin{cases} [(u_i + 1)^{\lambda} - 1] / \lambda, & \lambda \neq 0, u_i \geq 0 \\ \ln(u_i) + 1, & \lambda = 0, u_i > 0 \\ - [(-u_i + 1)^{2-\lambda} - 1] / (2-\lambda), & \lambda \neq 2, u_i \leq 0 \\ - \ln(-u_i + 1), & \lambda = 2, u_i \leq 0 \end{cases}$$

Binning & Binarization :-

and to handle outliers



Skewed d-set.



Numerical data \rightarrow categorical data.

Binning

(Descretization)

Binning :- A process to transform continuous variable into

discrete variable by creating

range of values of it.

outliers detected remove from data.

To handle outliers

To improve the value spread.

uniform the

data - spread.

Age

23, 35, 57, 81, 85

{ 0-10, 10-20, 20-30 }



bining apply

Binning

unsupervised

Supervised

equal width
(uniform)

Decision
tree binning

Equal frequency

Select the no. of bin

K-Mean binning.

like 1st, 2nd, ..., 10th percentile

use K-Means

based on distribution

'bisect' at mid-way

OR

Calc distance of each
point from each centroid

Suppose

bins = centroids = 5

value - spread even

random created

\rightarrow K-Bin Discretizer

\rightarrow encoding \rightarrow ordinal / one

\rightarrow strategy

Handling missing values :-

→ If we have some missing values in our dataset then we have to handle them.

Missing value.

→ We can handle them by removing item.

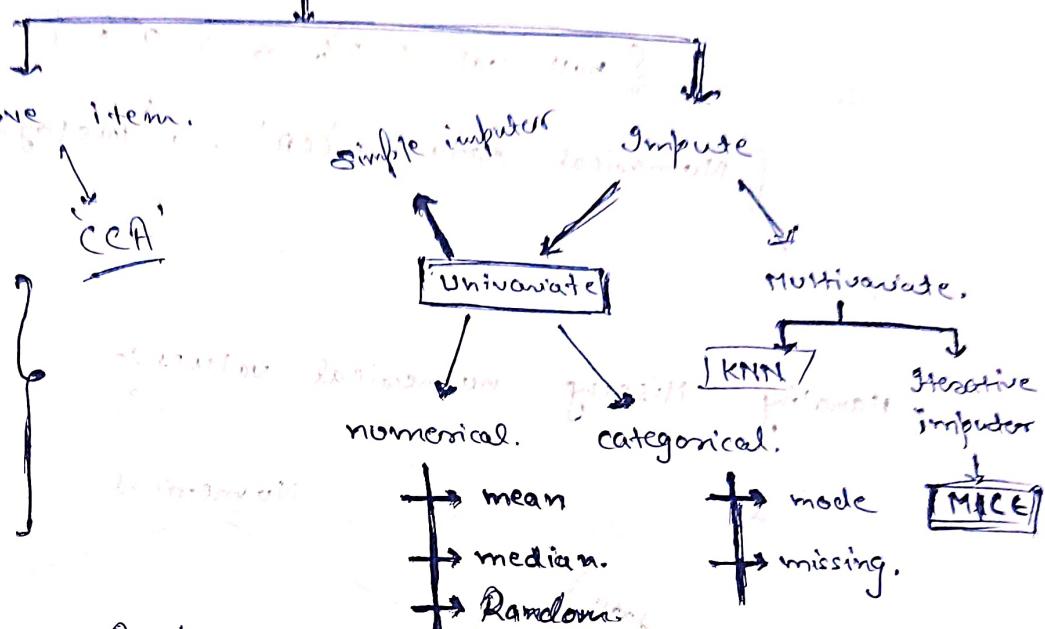
→ Remove value.

→ simple Imputer

→ KNN imputer

→ Iterative

→ Missing indicators.



'CCF' - complete Case Analysis:-

dist-wise deletion: Consists in discarding

observations if in any of variable are missing.

Advantage

Disadvantage

- easy to implement as no data manipulation required.
- Preserves variable distribution.

Disadvantage :-

- It can exclude a large fraction of original dataset.
- Excluded observations could be informative for the analysis.
- When using our models in production → it will not handle missing values.

[MCAR → 5% < CCF]

CCA conditions:-

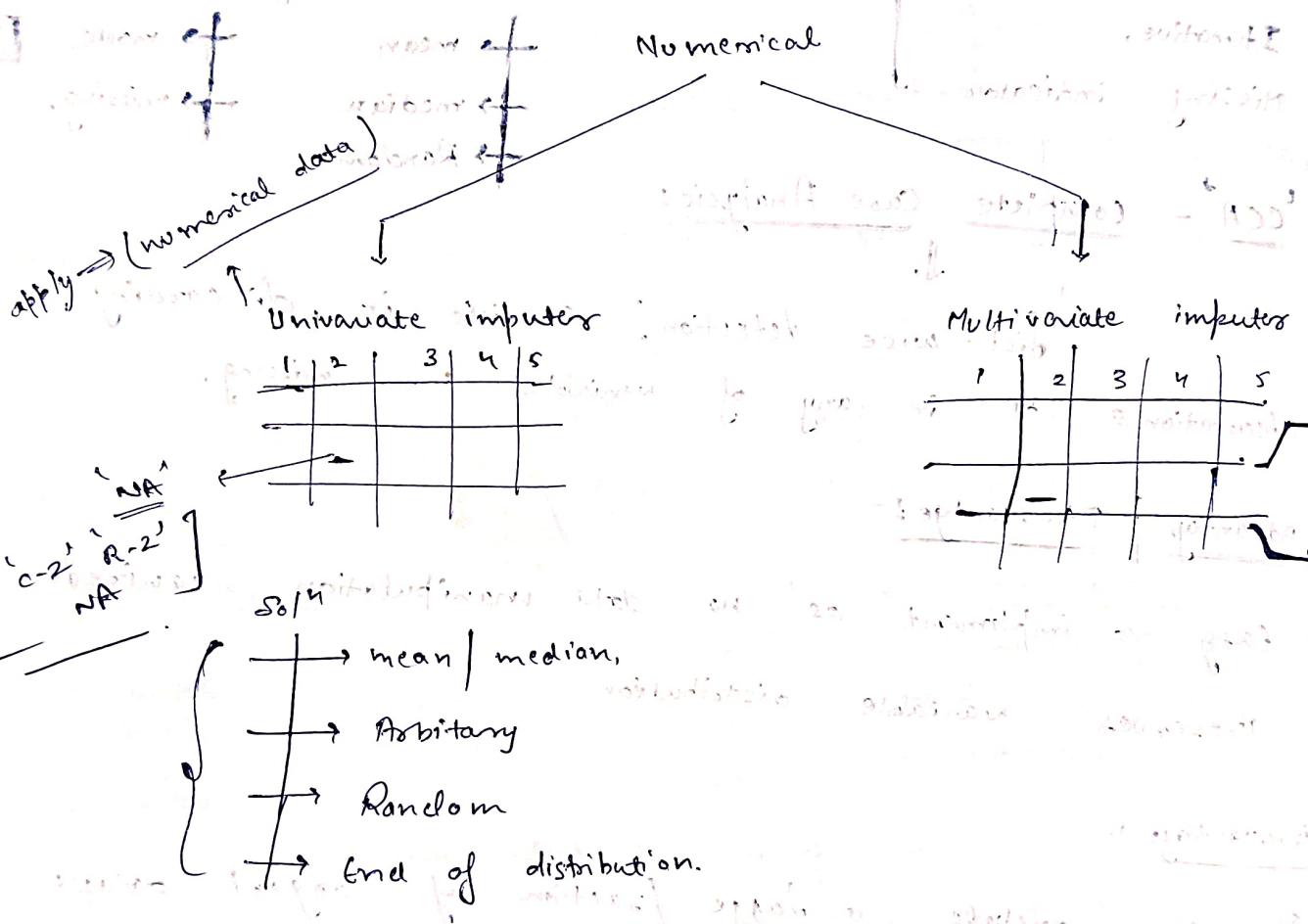
[column → apply → data missing < 5%]

↓ { null val < 5% & > 0% } [condition]

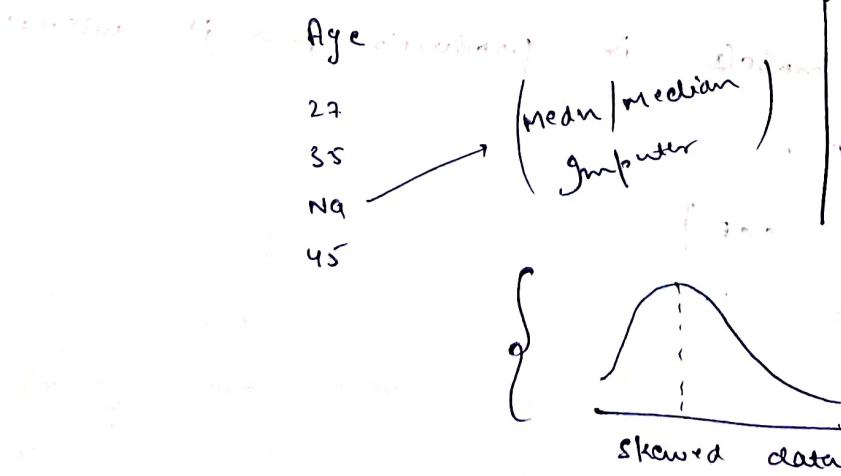
[Numerical data → 'cca' → 'histogram' plot before & after]

Comparable

Handling Missing numerical values :-



Mean | Median :-



Normal distribution

mean → apply

median → apply

Benefits :-

- 1). Simplicity (server deploy).
- 2). $\leq 5\%$ less serialize.

When to use :-

- 1). semi data missing $< 5\%$.
- 2). Missing completely at random.

Disadvantage :-

- 1). Distribution shape change.
- 2). Outliers formed.
- 3). Correlation / Covariance changes

mean others column so result vary.

→ Variance shrinks

Arbitrary Value Imputation :-

Application Categorical → 'NA' → 'Missing' → fill val

{ Benefit :- easy to apply }

⇒ PDF graph distribute

⇒ covariance also changes

⇒ variance

Numerical data

↳ 'Arbitrary No. use'

(Used when data isn't missing at Random)

End of distribution

Imputer :-

{ Arbitrary val → missing replace }

PDF changes

Variance varies.

(IQR - inter quartile range)



↳ Normally distributed

↳ (Mean + 3σ)
or
(Mean - 3σ)

missingness.

↳ skewed data (IQR proximity)

Q1 - 1.5 IQR
Q1 + 1.5 IQR

{ Benefit }
easy

Data handling

Missing

Categorical

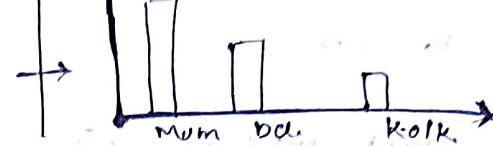
Data :-

replace that value by mode
condition :- missing completely at random

mode

missing

It change distribution of data.



mode \rightarrow more time in data

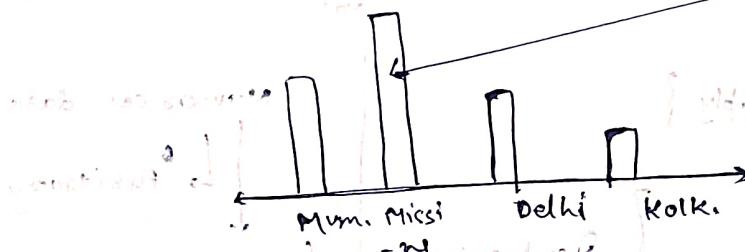
Missing category

Imputation :-

missing val. $> 10\%$ or more

replace with mode (x)

or create a new category Name \rightarrow Missing.



Random Imputation :-

Age

2

9

8

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

-

</div

Missing

Indicator :-

Age	Fare	Age - NA'
27	29	f] → for all present value.
41	32	f] → for all present value.
NA	31	(+) → missing val.
52	41	

Automatically select value for imputation:

grid search CV

grid search cross validation

give best parameters

automates the search with reliable evaluation.

step ↴

→ Define model

→ Define a parameter

→ gridSearchCV tries All combination

→ Cross-val.

KNN Imputer :-

	1	2	3	4	5
1	-	-	-	-	-
2	-	-	-	-	-
3	-	-	-	-	-

KNN Imputer
Iterative Imputer (Multivariate case)
 $(-) = \text{NA}$

Calculate 'the distance'

to find $\text{feature } 2, 4$ → calculate distance from each row

'least value' → 'feature 1 val. use replace'

KNN

'k = no. of neighbours'

find the value (take val. from neighbor)

calculate mean

general distance

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

This case val. is missing more than

non-euclidean formula

$$\text{dist}(u, y) = \sqrt{(\text{weight}) \times (\text{distance from pr. coordinates})^2}$$

calculate = 83

Apply

SN	F-1	F-2	F-3	F-4
1	33	—	67	21
2	—	45	68	12
3	23	51	71	18

$$\text{weight} \times (s1-45)^2 + (f1-68)^2 + (t1-12)^2$$

$(3/3)$

$$\text{weight} = \frac{\text{total no. of co-ordinates}}{\text{total no. of pr. coordinates}}$$

Advantages/Disadvantages :-

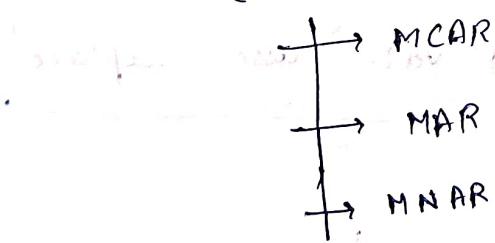
1). More accurate

2). More no. of calculation.

3). Production → train set display on server

Multivariate Imputation by chained equations :-

(MICE)



Adv | Disadv:-

Slow

→ server training
memory

R&D	Admin	M.Span
8	15	30
-	5	20
15	10	41
12	15	26
2	15	M.NAN

→ NAN replace by mean of each column.

step-II. move left to right in column

replace missing value with NAN, only first one

R&D	Admin	M.Span
8	15	30
NAN	5	20
15	NAN	41
12	15	M.NAN
2	15	M.NAN

→ training data (input)

R&D	Admin	M.Span
8	15	30
NAN	5	20
15	mean 15	41
12	15	mean

→ prediction

(check result)

+ next step

- 22

⇒ Apply linear regression / ML model to train & predict.

likewise same for again repeat each step

Then \rightarrow ((Iteration - 1) - (Iteration - 0))

loop till the difference ≈ 0

positive

negative

zero

loop

loop

loop

Converged

Converged

Converged

Converged

Converged

Converged

Converged

Converged

Converged

Outliers - Detection:

Outlier :- A datapoint out of bound 8 or different behavior

linear algebra

AdaBoost

logistic regression

all

affected by outliers

Outlier's treatment

Trimming

- thin (data)
- fast

Capping

(Apply limit)

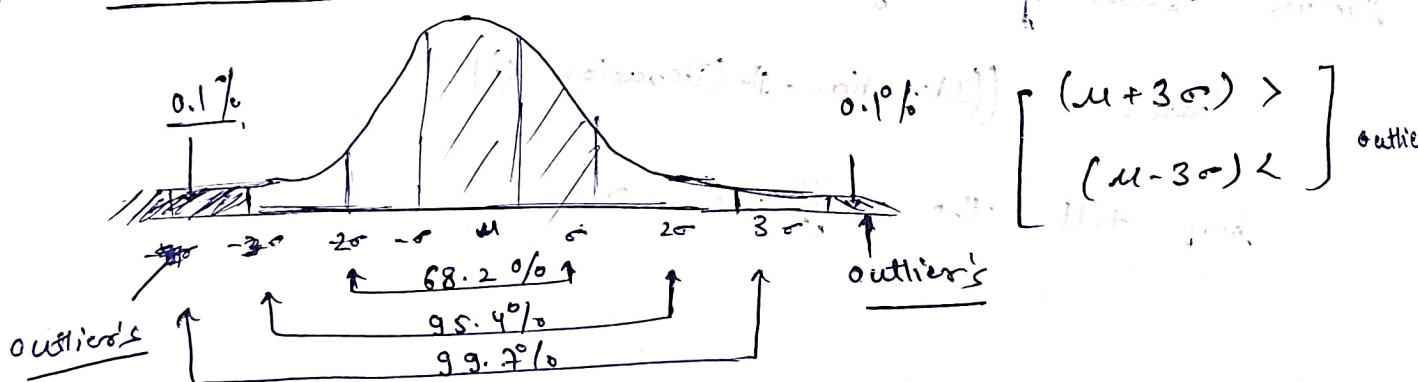
nesting value

discretization

for numerical values

How to detect Outliers:-

1). Normal distribution:-



2). Skewed Distribution:-

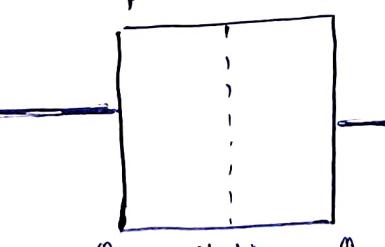
IQR (Inter quartile range)

outlier

↓
0 0 0

minimum

$(Q_1 - 1.5 \times IQR)$



(25%ile)

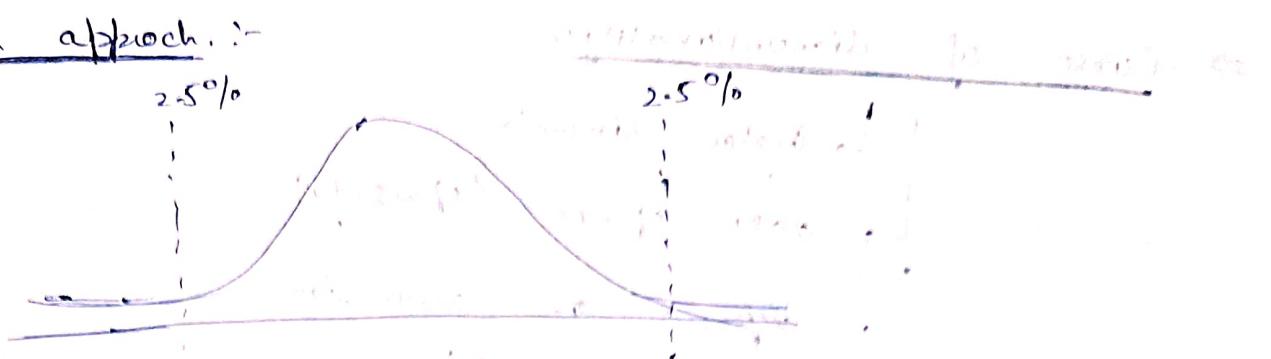
(75%iles)

outlier
100

Max.

$(Q_3 + 1.5 \times IQR)$

3). Percentile-based approach:



We select a particular range & apart from it everything is an outlier.

Technique for outlier Detection and Removal.



(a) z-score treatment

(z-score treatment)

IQR-based filtering

Percentile

Winsorization
(Percentile method-capping)

Outlier's removal using Z-score treatment

$$Z\text{-Score} = \left(\frac{x_i - \mu}{\sigma} \right) \quad \text{consider range } [-3, 3].$$

where μ is mean & σ is standard deviation.

↳ condition (T/F)

if true remove

outlier treatment

(a) Stepwise

⑤ Outlier

remove row

Trimming & Capping

↳ Trimming or Capping (apply cap to a particular)

Outlier distribution using IQR :- Suppose range

$$\left[\min(Q_1 - 1.5 \text{ IQR}), \max(Q_3 + 1.5 \text{ IQR}) \right] \quad (\text{IQR proximity rule})$$

↳ Outliers

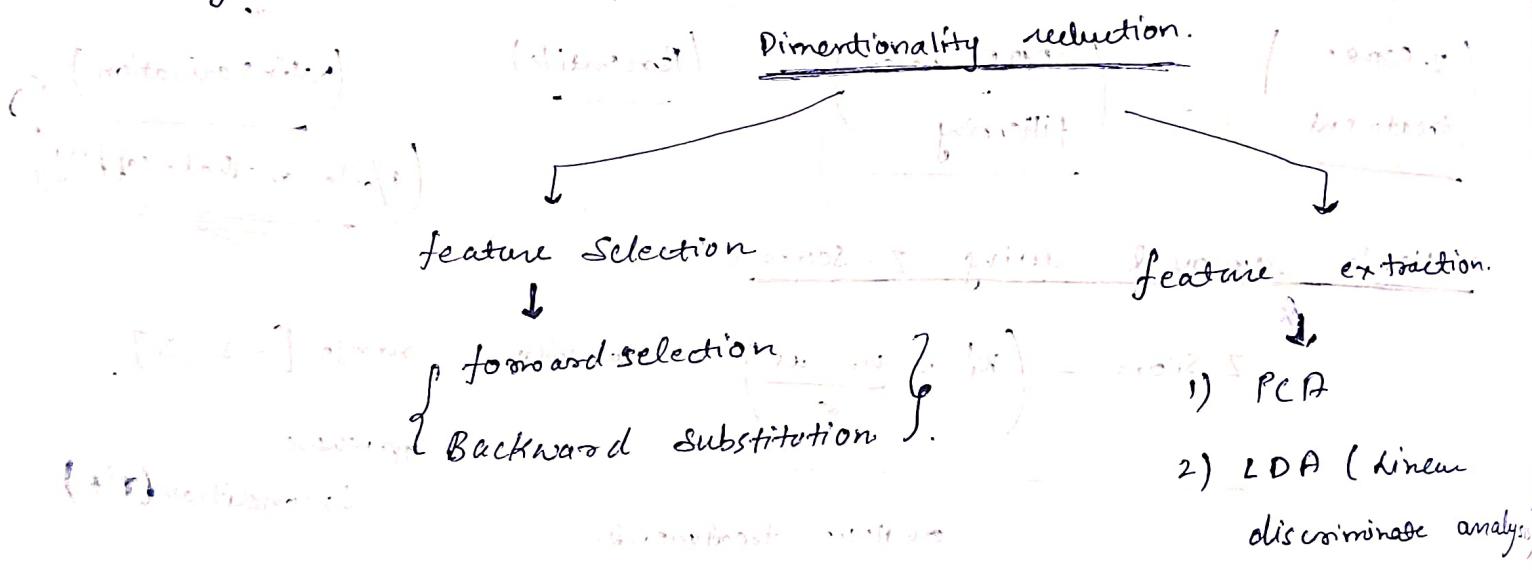
↳ remove

Curse of dimensionality:-

- ↳ higher dimension
- ↳ data space (sparsity).
- ↳ depend on dimension
- ↳ (1D, 2D, 3D)
- ↳ dimension (\uparrow) finding particular point (difficult).

Dimensionality reduction \rightarrow soln.

- \Rightarrow Performance decreases (\downarrow)
- \Rightarrow higher computation.



Principle Component Analysis (PCA):-

- ↳ unsupervised ML problem.
- ↳ only input no output.
- ↳ old technique.

feature extraction technique to reduce curse of dimensionality

Benefit of PCA:-

[even reduce '10-D' in '3-D']

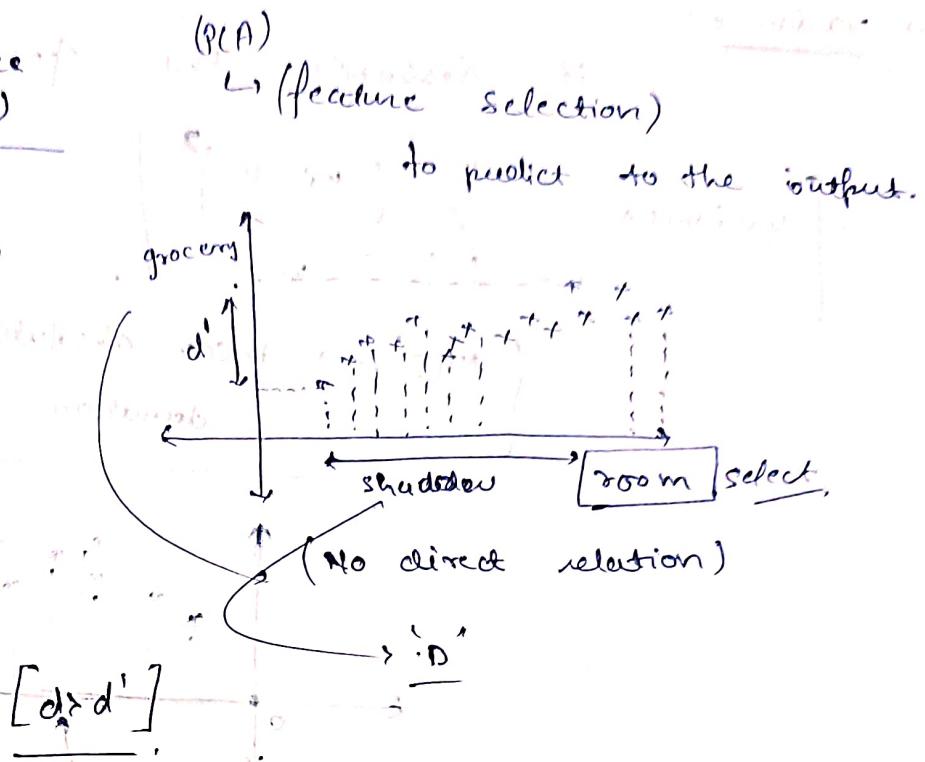
- { 1) visualization.
- { 2) faster

geometric intuition :-

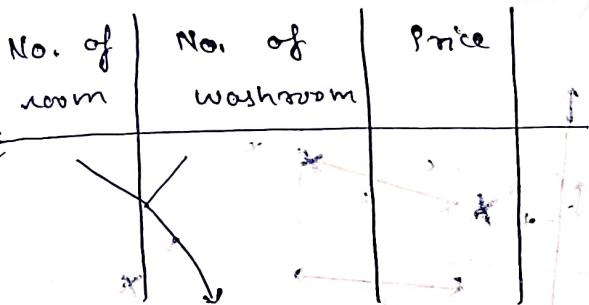
No. of rooms	No. of grocery shops	price (L)
3	2	60
4	0	130
5	6	170
2	10	90

↑
(remove)

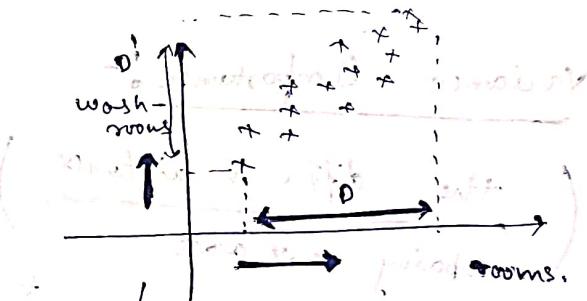
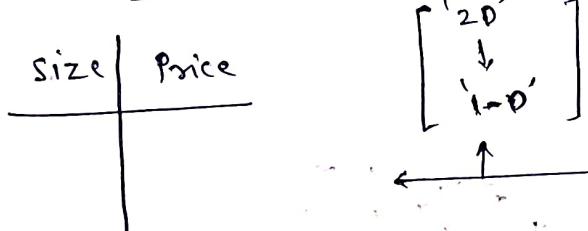
$$\left\{ \begin{array}{l} D = \text{shadow} \\ d' = \text{shadow of } x\text{-axis} \end{array} \right\}$$



changes in 'd-set' (feature combining)



'combine & get the size'

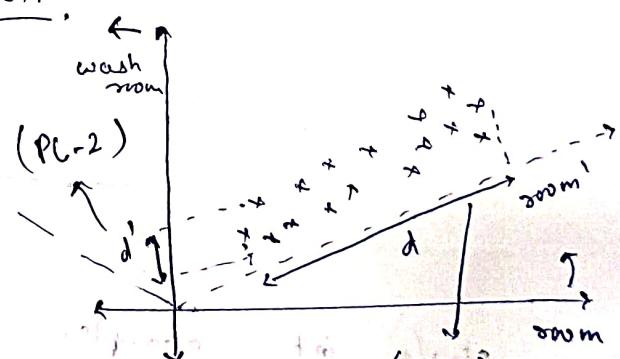


(feature extraction works here)

$(d > d')$

No. of principle component $\leq n$

(total no. of original feature)
in data.

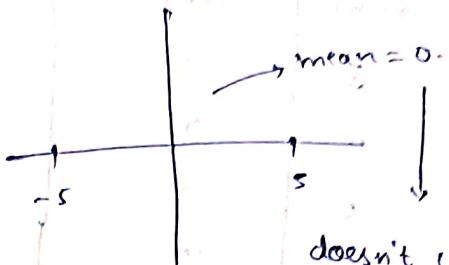


⇒ Variance :-

it describes the spread of data.

$$\text{Variance} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}$$

{'MAD' → mean absolute deviation}



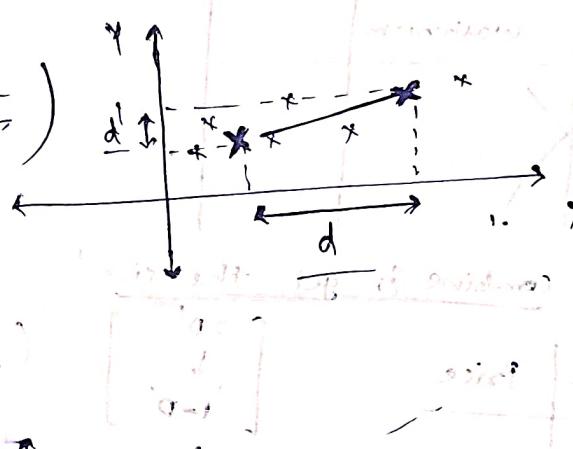
doesn't classify diff. datasets



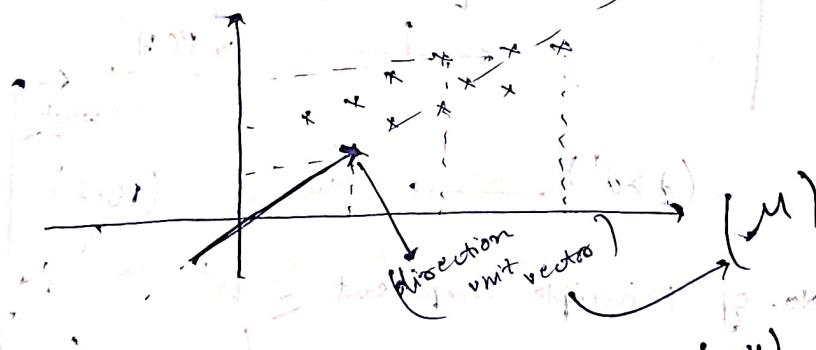
(variance calculate)

Variance Importance :-

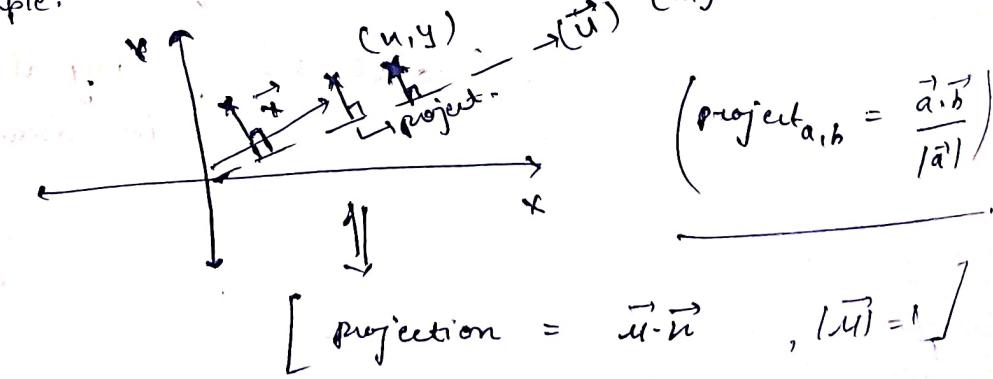
(the diff. in distance
compared to x,y axis)

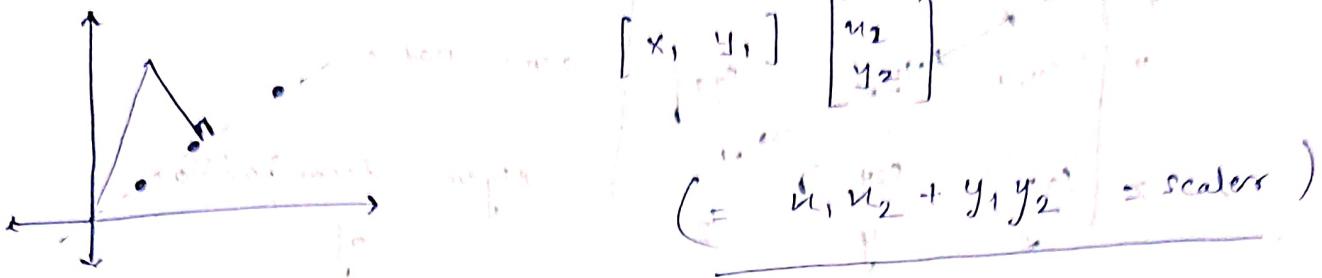


⇒ problem formulation :-



single-point example.





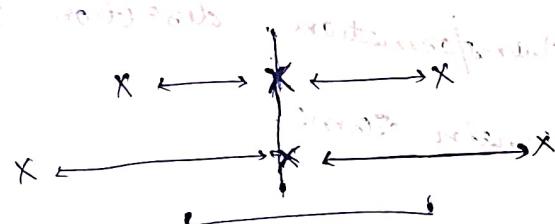
variance allocation

variance allocate $= 1/n \sum_{i=1}^n (u_i - \bar{u})^2$

$$\left[\frac{\sum_{i=1}^n (u^T u_i - u^T \bar{u})^2}{n} \right] = (\bar{u} - \bar{u}) \text{ variance}$$

optimization

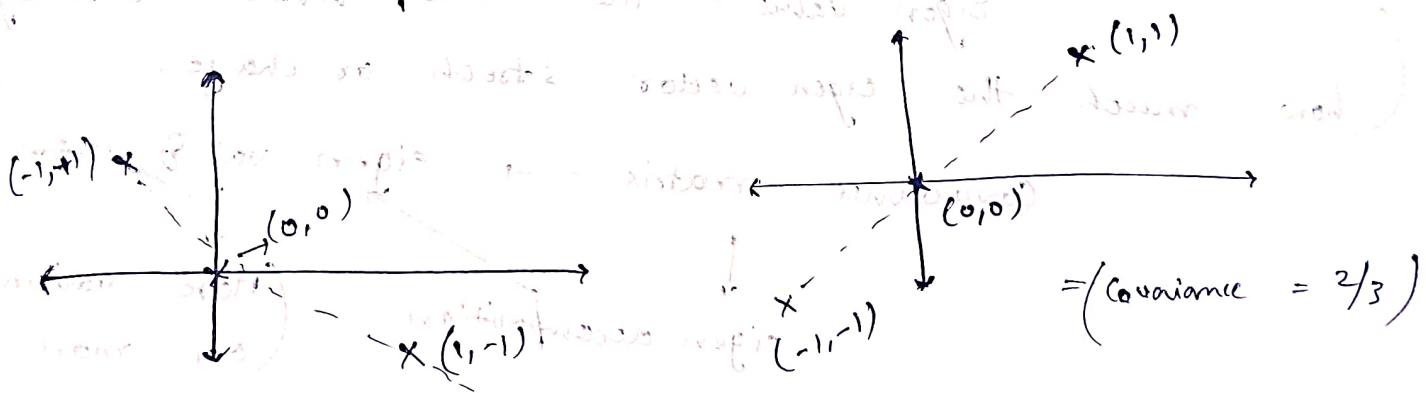
Covariance & Covariance Matrix



Variance problem
with

[diff. on x,y axis]

mean problem



(Covariance \rightarrow relation b/w x, y axis)

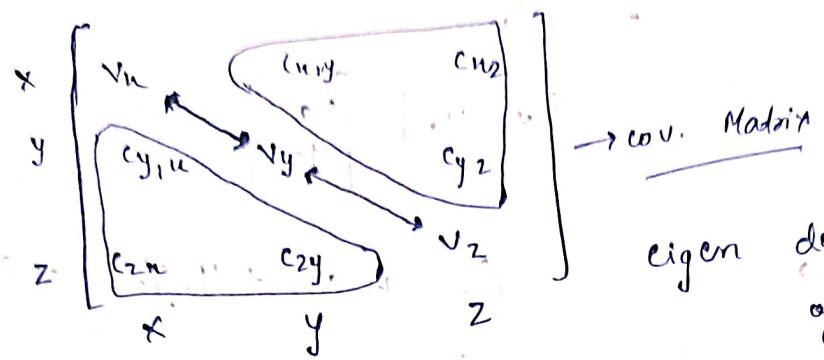
$$\text{Covariance} = (-2/3)$$

(diff. co-relation).

$$\begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} \text{ cov} = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} \text{cov}(u_1, u_1) & \text{cov}(u_2, u_1) \\ \text{cov}(u_1, u_2) & \text{cov}(u_2, u_2) \end{bmatrix} \xrightarrow{\text{matrix}} \begin{bmatrix} -1 & 1 \end{bmatrix} \text{ var.}$$

$$\text{cov}(u_1 b) = \text{cov}(b, u_1)$$

$$= \begin{bmatrix} u_1 & u_2 \end{bmatrix} \begin{bmatrix} \text{var}(u_1) & \text{cov}(u_2, u_1) \\ \text{cov}(u_1, u_2) & \text{var}(u_2) \end{bmatrix} \xrightarrow{\text{matrix}} \begin{bmatrix} \text{var}(u_1) & \text{cov}(u_2, u_1) \\ \text{cov}(u_1, u_2) & \text{var}(u_2) \end{bmatrix} \quad (= \text{square matrix})$$



Linear Transformation, Eigen Vector and Eigen Value :-

(geogebra.org)

Matrices are linear transformation.

Eigen vector: special vector which remains same after transformation, direction don't change. Span remain same.

Eigen value: the value which usually describes how much the eigen vector stretch or change.

Covariance matrix \rightarrow eigen val & vector.

eigen decomposition \downarrow More variance more stretching.

- Step-by-step soln :-

1). mean centre

2). find covariance matrix.

3). find the eigen val. & vector for covariance matrix.

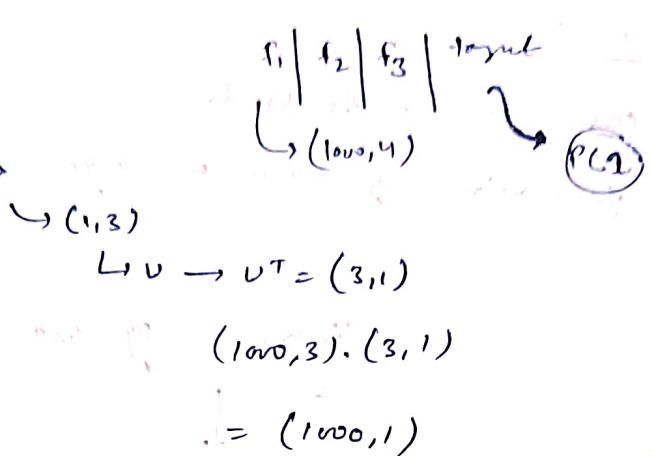
3D \rightarrow 3 vectors
PC1 PC2 PC3

\rightarrow convert in 2D, 1D

How to transform point :-



shape (1000, 3)



vertical input

vertical output

$\left[\begin{array}{c} f_1 \\ f_2 \\ f_3 \end{array} \right]$

horizontal input

horizontal output

vertical input

vertical output



vertical input

horizontal output

$\left(\begin{array}{c} f_1 \\ f_2 \\ f_3 \end{array} \right)$

horizontal input

vertical output

\downarrow

vertical input

horizontal output

Machine learning Algorithms:-

linear regression :-

→ easy to understand.

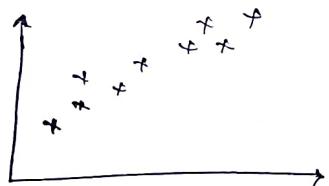
→ foundation of every ml algo.

($y = mx + c$) → linear

linear regression

Simple linear
Regression.

1 target var.
1 data var. (input)



Multiple linear
regression

More than 1
input var.

Polynomial
regression

Poly. equation
 $[y = an^2 + bn + c]$

$y = mx + c$, ↳ sort of linear data.

(m, c)

closed form

direct
method.

solve eqn

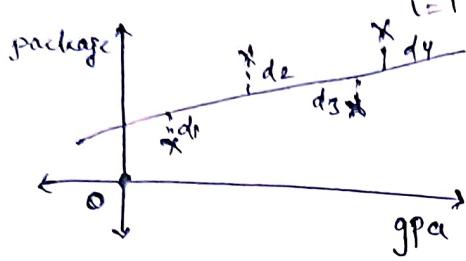
or

Non-closed
form.

↓
gradient
descent.

$$B = \bar{y} - mx$$

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



$$\epsilon = d_1 + d_2 + d_3 + \dots + d_n$$

$$E = d_1^2 + d_2^2 + d_3^2 + \dots + d_n^2 \text{ is 'tve'}$$

$|d_1| + |d_2| + \dots + |d_n| (x)$
not differentiable

$$[\epsilon = \sum_{i=1}^n d_i]$$

$$\text{error func.}$$

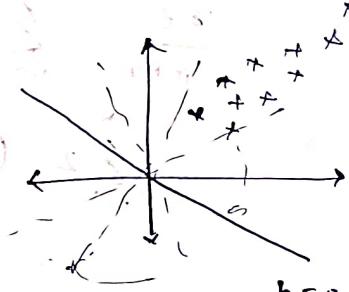
$$d_i = (y_i - \hat{y}_i)$$

$$[\epsilon = \sum_{i=1}^n (y_i - \hat{y}_i)^2]$$

$$e = \sum_{i=1}^n (y_i - \hat{y}_i) = E(m, b)$$

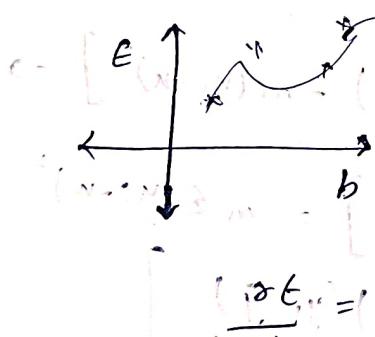
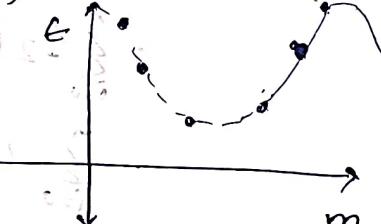
$$\left\{ \begin{array}{l} y = f(x) \\ E(m, b) \end{array} \right.$$

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



$$b = 0, E(m) = \sum_{i=1}^n (y_i - m x_i)^2$$

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



local maxima & minima concept :-

$$\left[\frac{\partial E}{\partial m} = 0, \quad \frac{\partial E}{\partial b} = 0 \right] \Rightarrow (m, b)$$

$$\frac{\partial E}{\partial b} = \frac{\partial}{\partial b} \sum_{i=1}^n (y_i - m x_i - b)^2 = 0$$

$$= \sum_{i=1}^n \frac{\partial}{\partial b} (y_i - m x_i - b)^2 = 0$$

$$\sum_{i=1}^n 2(y_i - m x_i - b) = 0$$

$$= \sum_{i=0}^n (y_i - m x_i - b) = 0$$

$$= \sum_{i=0}^n y_i - \sum_{i=0}^n m u_i - \sum_{i=0}^n b = 0$$

divide by n .

$$= \frac{1}{n} \sum_{i=0}^n y_i - \left[\frac{\sum_{i=0}^n m u_i}{n} + \frac{\sum_{i=0}^n b}{n} \right] = \frac{0}{n}$$

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

$$\boxed{(\bar{y} - m \bar{u}) - b = 0}$$

$$b = \bar{y} - m \bar{u}$$

$$\Sigma = \sum (y_i - m u_i - \bar{y} + m \bar{u})^2$$

$$\frac{\partial \Sigma}{\partial m} = \sum_{i=0}^n \frac{\partial}{\partial m} (y_i - m u_i - \bar{y} + m \bar{u})^2 = 0$$

$$= \sum_{i=0}^n 2 (y_i - m u_i - \bar{y} + m \bar{u})$$

$$= \sum_{i=0}^n 2 (y_i - m u_i - \bar{y} + m \bar{u}) - (m u_i + \bar{u}) = 0$$

$$= \sum_{i=0}^n 2 (y_i - m u_i - \bar{y} + m \bar{u}) (u_i - \bar{u}) = 0$$

$$= \sum_{i=0}^n (y_i - m u_i - \bar{y} + m \bar{u}) (u_i - \bar{u}) = 0$$

$$= \sum_{i=0}^n [(y_i - \bar{y}) - m(u_i - \bar{u})] (u_i - \bar{u}) = 0$$

$$= \sum_{i=0}^n [(y_i - \bar{y})(u_i - \bar{u})] - m \sum_{i=0}^n (u_i - \bar{u})^2 = 0$$

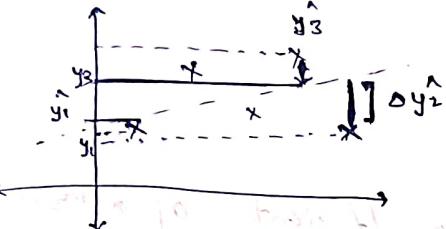
$$= \sum_{i=0}^n [(y_i - \bar{y})(u_i - \bar{u})] = m \sum (u_i - \bar{u})^2$$

$$\therefore \boxed{m = \sum_{i=1}^n \frac{(u_i - \bar{u})(y_i - \bar{y})}{\sum_{i=1}^n (u_i - \bar{u})^2}}$$

Regression Metrics :-

- 1) MAE
- 2) MSE
- 3) RMSE
- 4) R2 score
- 5) Adjusted R-2 score.

\Rightarrow MAE :- Mean absolute Error :-



$$\text{mae} = \frac{\sum |y_i - \hat{y}_i|}{n}$$

This advantage -

\Rightarrow Mod. funcn is not differentiable at '0'.

$$\text{Absolute error} = |y_1 - \hat{y}_1| + |y_2 - \hat{y}_2| + |y_3 - \hat{y}_3| \dots$$

$=$ ~~total error~~ \rightarrow total

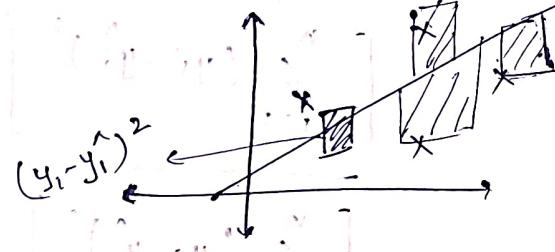
$$\text{mean absolute error} = \frac{|y_1 - \hat{y}_1| + |y_2 - \hat{y}_2| + |y_3 - \hat{y}_3| \dots}{n}$$

$$\left[\sum_{i=1}^n |y_i - \hat{y}_i| \right] \Rightarrow \begin{array}{l} \text{Advantages} \\ \downarrow \\ \text{direct ob.} \end{array}$$

\Rightarrow Same unit (gPA vs LPA)

\Rightarrow robust to errors
(detect & less diff. --).

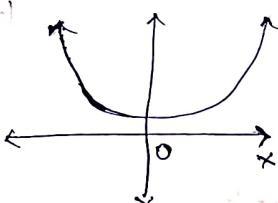
\Rightarrow MSE :- Mean Squared Error :-



$$mse = \frac{\sum (y_i - \hat{y}_i)^2}{n}$$

Advantage :-
use as loss function.

differentiable at every point.



$$\begin{cases} Y = LPA \\ mse = (LPA)^2 \end{cases}$$

Disadvantage :- Not robust to outliers.

RMSE :- Root mean square error.

$$= \sqrt{MSE}$$

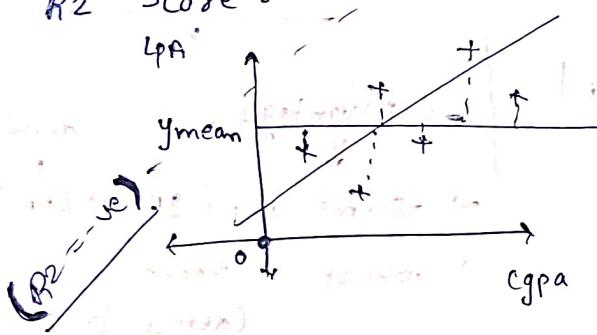
$$= \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$$

$MSE = LPA$ benefit

$$y = LPA$$

disadvantage :- Robust - collinear.

R² score :- coefficient of determination



R^2 or goodness of fit.

$$R^2 = 1 - \frac{SSR}{SSM}$$

$SSR \rightarrow$ sum of squared errors in regression line.

$SSM \rightarrow$ sum of sq. error in mean line.

$$R^2 = 1 - \frac{\left[\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right]_{reg.}}{\left[\sum_{i=1}^n (y_i - \bar{y})^2 \right]_{mean}}$$

Adjusted R² Score:-

$(R^2 - score)$
↳ R² adjusted

$$R^2_{adj.} = 1 - \left[\frac{(1-R^2)(n-1)}{(n-1-k)} \right]$$

$R^2 =$
 $n = \text{no. of variables}$
 $k = \text{independent}$

$k = 1, 2, 3, \dots$