

IEFOR EN-F8F .OM

Sr. No.	Title	Page No.	Sign./Remarks
			override 100% ✓

PAGE NO.:

⇒ what is multicollinearity b/w ?

$$Y = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

→ In this eq we assume that w_i is denote how y increase or decrease with respect to x_i while all other x_j (ie when $j \neq i$) is constant

→ Coefficient w_i gives relation b/w y and x_i while all others x_j is constant

Is multicollinearity b/w

No, If we want prediction only using linear regression
yes, If we need to find relation b/w y & x by using w
because here w is not important it has no meaning If there is multicollinearity

TYPES

Structured

by creating new feature
Ex: one hot encoding
city A | B | C
A | B | C need to remove one column

Data base
Data itself contain multicollinearity

How to remove it → increase data

→ remove one correlated column

→ use Lasso, Ridge regression

→ use PCA Principle Component Regression

⇒ multicollinearity more affect on parametric ML algo Ex: Linear, logistic regression, Naive bayes, KNN algo where we take assumption that all input columns are independent

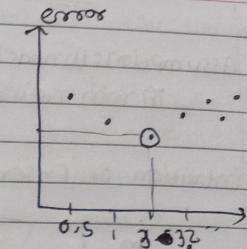
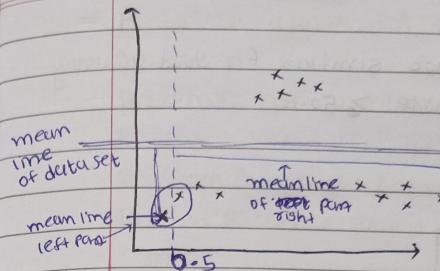
while non parametric. Ex Tree base algo and algo where no assumption that input cols are independent

⇒ Fit(), transform(), Fit_transform()

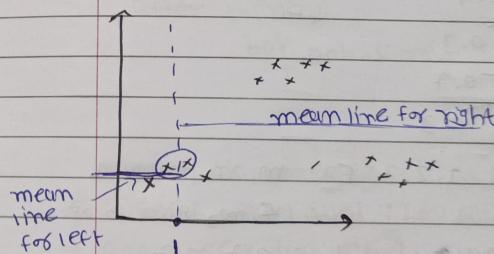
Transformers: Standard Scaler, Min-Max Scaler, PCA, Imputer etc
fit() is only compute Ex: compute mean, standard deviation, compute none values using median
transform() transform data Ex: convert each data point to $Z = \frac{x - \mu}{\sigma}$
Imputer put values to na values

→ fit() compute parameters and weights

Regression Trees



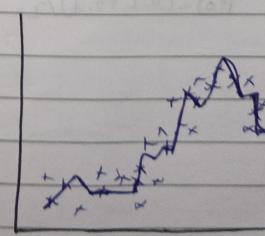
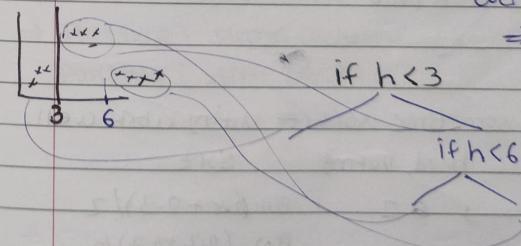
⇒ In this method every time we take two successive point and draw line at mid of both and for every part (left and right) we take min SSE (some square error) and Plot graph for every two successive point and consider those vertical line for which SSE is min.



⇒ we do this for every two successive and take that line which has min error

⇒ now for every part left of line and right of line we do same step

⇒ when point less than threshold we stop process for that part

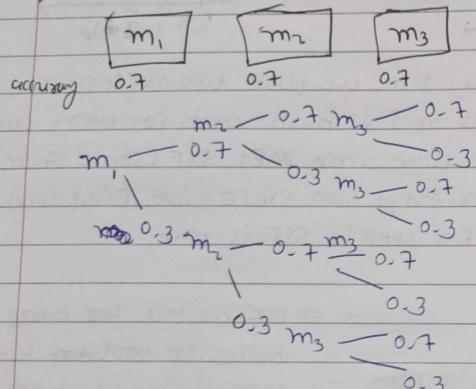


Regression Curve

Voting Ensemble

Assumption: i) model should be not similar for good accuracy
ii) min accuracy of every model $\geq 50\%$.

intuition of Ensemble techniques



OP is mux OP from model 1, 2, 3 Ex: m_1, m_2 gives 1

& m_3 gives 0 OP

then we choose those Puth which has 0.7, 0.7
and find mux prob

Hard & Soft Voting

⇒ Some time hard and sometime soft voting work well.

hard voting	Soft
m_1	$P(1) = 0.8$
	$P(0) = 0.2$
m_2	$P(1) = 0.3$
	$P(0) = 0.7$

$$\rightarrow 1 \quad P(1) = (0.8 + 0.3)/2$$

$$P(0) = (0.2 + 0.7)/2$$

→ Voting & boosting ~~not~~ can have diff weight while Bugging has sume.

Bugging

⇒ Bootstrap aggregation : for sume model diff data training

Bugging: Row sample with replacement \rightarrow bootstrap = True

Pusting: " " without " " bootstrap = False

Random Subspaces: Column sampling with or without replacement

Random Patches: Row and column sampling with or without

OOB (out of Box): there are some ~~some~~ data exist which is not feed into any of model

→ So we can use that data in testing arr also

oob-score = True inside bug= $\text{bug} = \text{True}$ inside $\text{Bugging}(\text{list of } \dots)$

after training

bug.oob-score give accuracy of that OOB data

Bugging Vs Boosting

Bugging: Type of model used: Low bias high var, High bias low variance

⇒ Bugging = Parallel

Boosting = sequential

⇒ Bugging = same weightage

Boosting = diff weights

Random forest

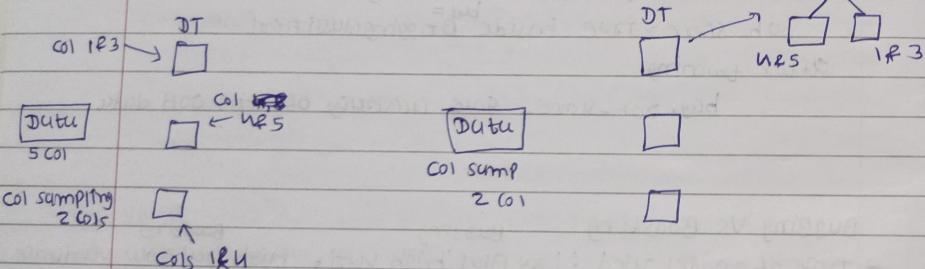
⇒ In Random forest noisy data distribute among all the DT and not feed in one DT so impact of noisy data distribute over trees not stored in one DT

⇒ Difference between Random Forest and Bagging:

Bagging can have \rightarrow DT only also
↓
CNN
SVM

Random forest has DT only

✓ Bagging



⇒ Bagging: Tree level sampling

before every tree columns are fixed

⇒ RF: node level sampling

in every tree, every node has select random columns.

Feature Importance using Random Forest and Decision Trees

→ For DT:
 (sum of importance of every node with column C₁)
 gives feature importance of C₁
 (sum of importance of all cols)

node importance = IG = information gain of
that node

$$\Sigma \text{ Feature impo} = 1$$

→ impurity based feature importance can be misleading for high cardinality features (many unique values)
 In this kind of case use sklearn.inspection.permutation_importance

Adaboost

→ different weights

+	+
+	-
+	-
+	-

①

②

③

+ve	-ve
+	+
+	-
+	-

upsampling wrong classified point

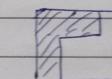
+	+
+	-
+	-
+	-

PAGE NO.:
1

+ve	-ve
+	+
+	-
+	-

+	+
+	-
+	-
+	-

+	+
+	-
+	-
+	-



decision stump

(decision tree
with depth = 1)

can use any model

Eg logistic regression, SVM

(KNN not work ∵ simple base, ...)

$$h(x) = \text{sign}(d_1 h_1(x) + d_2 h_2(x) + d_3 h_3(x))$$

x ₁	x ₂	y	y _{pred}	weight
3	7	1	1	0.2
2	9	0	0	0.2
wrong pred	1	1	0	0.2
9	8	0	0	0.2
3	7	0	0	0.2

* $x_1 > 5, x_2 < 10, x_1 > 2, x_2 < 8$

choose DS which has high IG

let that DT has model m,
 ⇒ now find weight of that DS

weight = d_i

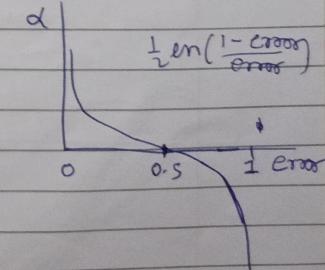
depends on error rate

$$\text{error for } m_i = 0.2 + 0.2$$

$$\text{error} = 0.4$$

$$d_i = \frac{1}{2} \ln \left(\frac{1 - \text{error}}{\text{error}} \right)$$

$$d_i = 0.2$$



→ update weights

$$\text{misclassified point} = w e^{\alpha_1}$$

$$\text{correctly classified point} = w e^{-\alpha_1}$$

x_1	x_2	y_{true}	weight	updated weight	range
3	7	1	0.2	0.16	0.166
2	9	0	0.2	0.24	0.25
1	4	1	0.2	0.24	0.25
9	8	0	0.2	0.16	0.166
3	7	0	0.2	0.16	0.166
				0.96	1

If model has good then d of that model is also high and weight change is high
→ $\alpha < 0$ mean model gives more than 50% wrong Pred so let's say 80% wrong so we take $\sim \alpha$ of it and get 80% correct

5 random numbers between 0-1

let random numbers are 0.13, 0.43, 0.62, 0.3, 0.81

so we take 1, 3, 3, 2, 0.12

now for next model and sume

step as * to till now.

Final op of any input is $a_1 h_1(x) + a_2 h_2(x) + a_3 h_3(x) + \dots$

⇒ we can adjust Learning Rate

$$LR \cdot \alpha$$

So change in weight is small

so update of amplitude will be small and reduce overfitting.

Gradient boosting

C6 salary pred1 pred2

8	3	4.8	1.8
7	4	4.8	
6	8	4.8	
9	6	4.8	
5	3	4.8	

i _q	CG	base					
		lpcu	Pred1	res1	Pred2	res2	Pred3
90	8	3	4.8	-1.8	-1.8	-1.62	-1.62
100	7	4	4.8	-0.8	-0.8	-0.72	-0.72
110	6	8	4.8	3.2	3.2	2.88	2.88
120	9	6	4.8	1.2	1.2	1.08	1.08
80	5	3	4.8	-1.8	-1.8	-1.62	-1.62

here overfitting
hence constant ex * ignore PAGE NO.

Gradient Boosting

i _q	CG	lpcu	Pred1	res1	Pred2	res2	Pred3
90	8	3	4.8	-1.8	-1.8	-1.62	-1.62
100	7	4	4.8	-0.8	-0.8	-0.72	-0.72
110	6	8	4.8	3.2	3.2	2.88	2.88
120	9	6	4.8	1.2	1.2	1.08	1.08
80	5	3	4.8	-1.8	-1.8	-1.62	-1.62

$$\text{Pred}_1 = \text{lpcu} - \alpha \text{Pred}_2$$

⇒ Pred 1 of base model is avg of O/P column lpcu

and then we train model 1 with

$$\text{IIP} = \text{CG}, i_9 \quad \text{and model gives OP}$$

so actual OP from Pred2 can be predicted by

$$\text{OP of base model} + \text{OP of model 1(Pred2)}$$

(Pred1)

⇒ To reduce overfitting we introduce learning rate to

$$\text{Pred of OP} = \text{OP of model 1} + \alpha \text{OP from model 2}$$

$$\text{So actual output} = m_0 + \alpha m_1$$

⇒ after that find res2 = actual OP - actual pred

lpcu

For $\alpha = 0.1$ we calculated res2 is in above table

$$3 - [4.8 + 0.1(-1.8)] = -1.62$$

$m_1 \Rightarrow$ for model 2

Train on IIP = i₉, i₆ O/P = res2 and model 2 gives OP of Pred 3

$$\text{Calcual OP} = m_0 + \alpha_1 m_1 + \alpha_2 m_2$$

$$\text{Pred from } m_2 \quad \text{OP of } m_1 \quad \text{OP of } m_2$$

$$\boxed{\begin{aligned} F(x) &= h_0(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots \\ \text{Final OP} & \end{aligned}}$$

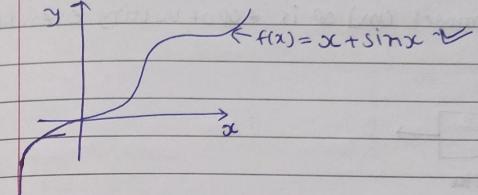
hyper para

hy pen

→ In Boosting we use DT us function and each fun is act like small part of final OP function (Additive modeling)

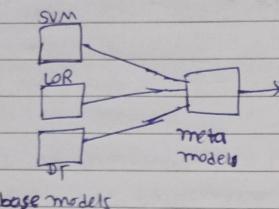
Additive Modelling

Additive modelling function has sum of small functions



Stacking

→ Same as Voting but we add new ML model which take input as OP of other model and OP is OP of Voting ensemble

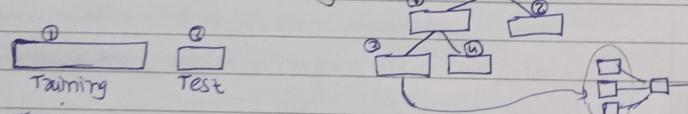


⇒ Stacking we have problem:

If base model is like DT then its give OP which contain overfitting and meta model also become overfitted model

→ Solution:

i) Hold out method: Bagging



② We also divide Training into part and that part is test for metamodel

ii) Cross-validation & Bootstrap

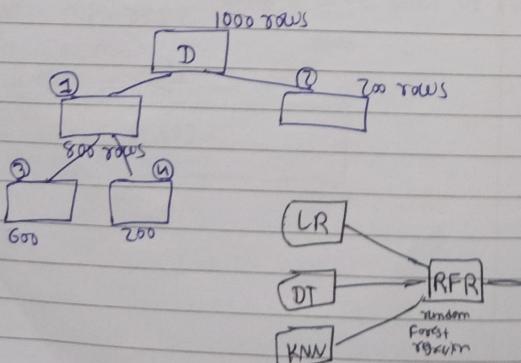
CG is package

7 70 4

9 101 5

8 120 6

1000
rows



Step 1: Train 3 base models on ③

Step 2: Using this 3 model and pred on ④

package	LR-Pred	DT-Pred	KNN-Pred
---------	---------	---------	----------

dataset →

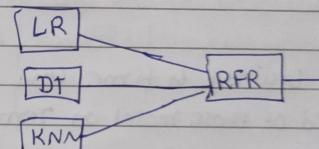
shape = (200, n)
(rows cols)

Step 3: Train metamodel on Step 2's dataset

IP = Pred of 3 model
OP = package (given)

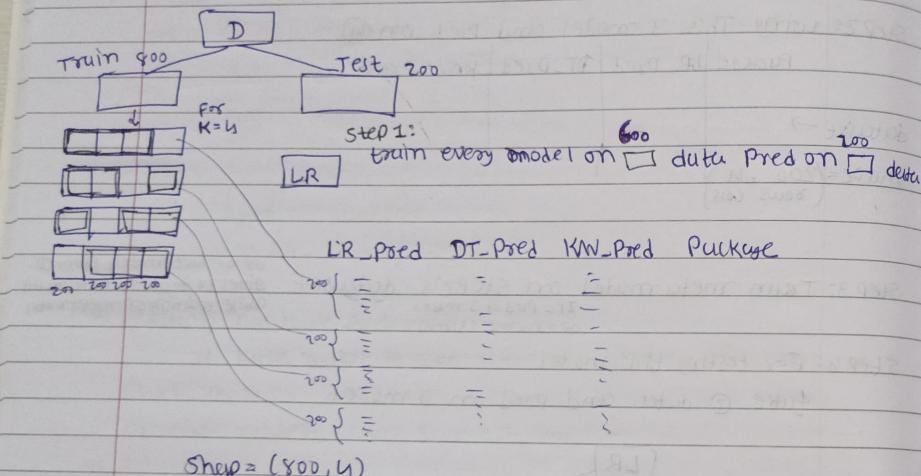
Step 4: For testing this model

take ② data and pred on ④ model



⇒ by this method we can overcome overfitting (while training of metamodel base models does not have training on ④)

ii) K-Fold approach - Stacking



⇒ Here every model trained 4 times (and above m_{test} dataset contain Pred of those model on remaining 200 data in each time)

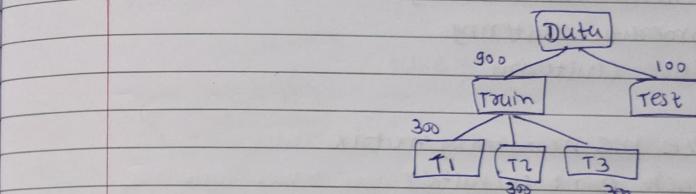
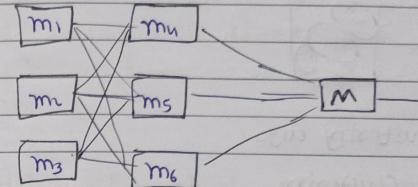
Step 2: using above data train meta model

Step 3: For base model

Train base models on Train data (800) and you get
your base models

Step 4: now train this model on Test (200) data

★ Multi-layer stacking



Step 1:

Train m_1, m_2, m_3 on T_1 and Pred on T_2 and we get
 m_1 Pred | m_2 Pred | m_3 Pred | package

shape $(300, 4)$

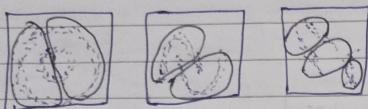
↑
Step 2: Train m_4, m_5, m_6 on This dataset and pred on T_3
 m_4 Pred | m_5 Pred | m_6 Pred | package

shape $(300, 4)$

↑
Step 3: now train meta model M on This data (300, 4)

Hierarchical clustering intuition

\Rightarrow K-means clustering not gives good result on



so we need other clustering algo

Types of Hierarchical clustering

- 1) Agglomerative clustering
- 2) Divisive clustering

Algo:

- 1) initialize the proximity matrix

2) make each point a cluster

3) inside a loop

a) merge the 2 closest clusters

b) update the proximity matrix

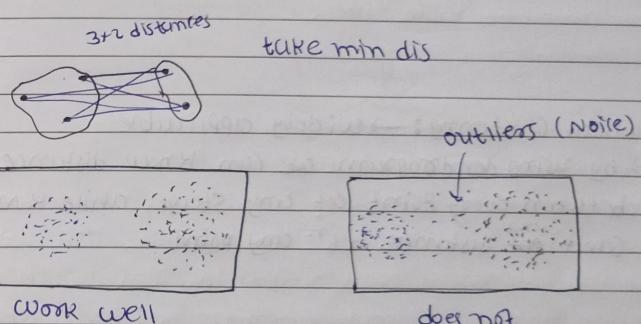
4) until only one cluster is left

	P ₁	P ₂	P ₃	P ₄	P ₅	
P ₁	0	-	-	-	-	P ₁ , P ₂ , C, P ₅
P ₂	-	0	2	-	-	P ₂ , 0
P ₃	-	2	0	-	-	C, 2, 2, 0, 2
P ₄	-	-	-	0	-	P ₄ , -
P ₅	-	-	-	-	0	

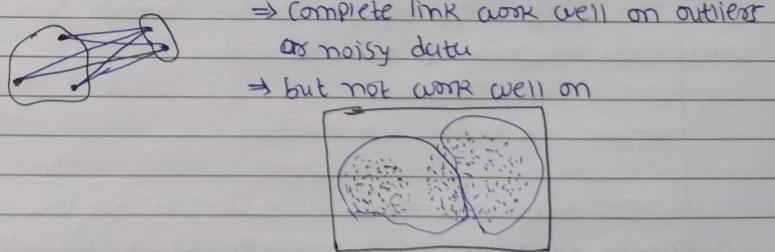
Type of Agglomerative clustering

- 1) Min (single-link)
- 2) Max (complete-link)
- 3) Average
- 4) Ward

1) Single Link (min)

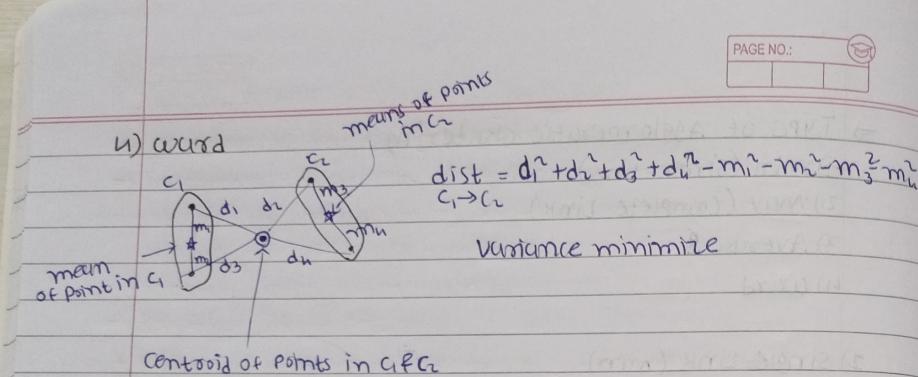


2) Complete link (max)



3) Group average

$$\frac{\sum d_i}{6} \quad i=1, 2, \dots, 6$$



Adv of Hierarchical Clustering: → widely applicable
 → by using dendrogram we can know distance between every point at any stage while K-means can't know distance b/w any point.

disAdv: → can't use on big datasets because distance matrix takes ^{large} memory

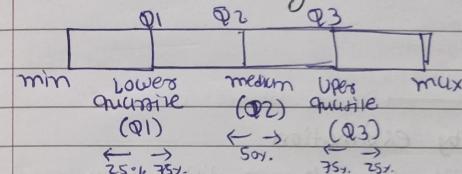
Statistics ← Descriptive
Inferential

→ Coefficient of variation (CV) : represented amount of variability in a dataset

$$CV = \frac{\sigma}{\bar{x}} \times 100$$

Can compare spread of 2 diff column we use CV

* 5 number summary in statistics



* Correlation and causation

→ The phrase "Correlation does not imply causation" means that two variables are associated with each other, it does not mean that one causes the other.

→ E.g. more ice cream sales more murder reason was ~~weather~~

* Random Variable

→ Algebraic Variable

x is unknown value $x + 5 = 10$

$$x = 5$$

→ Random Variable: a random variable is a set of possible values from a random experiment.

$$X = \{H, T\}$$

↑ sample space values that can have RV.

★ Parametric density estimation

→ is method to estimate prob density function (PDF) of random variable by assuming that underlying distribution belongs to specific parametric family of prob distri such as normal, expo, Poisson distri.

⇒ we have sample $\sigma \in \Omega$

by using assumption of particular distribution find Probability of ~~at~~ point and make graph

★ Non-parametric density estimation

→ If distri not match with any known distri then use -

→ It is statistical technique used to estimate PDF without making any assumption about underlying distribution.

→ It is called non-parametric density estimator because it doesn't require any assumption or predefine distri function.

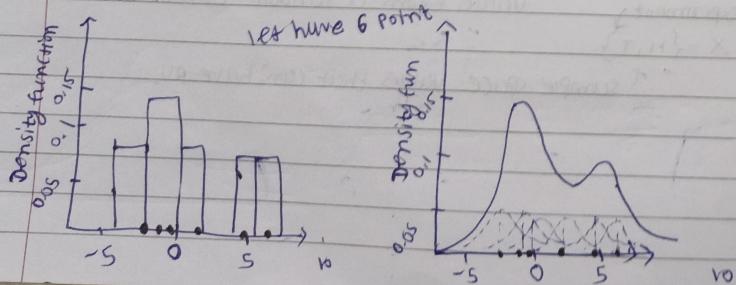
→ This technique involves construct estimate of PDF using the available data. This done by creating kernel density estimate.

adv: can find PDF of any kind of data

disadv: computationally expensive

not give good result with less data.

KDE: Kernel density fun to smooth out the data and create a continuous estimate of the underlying density function.



→ For every point consider its mean and draw gaussian distribution or kernel. and all all graph

→ mean = each point } for kernel (cumbe curv Ex gaussian

σ = bandwidth

hyper para

from sklearn.neighbors import KernelDensity

model = KernelDensity(bandwidth=3, kernel='gaussian')

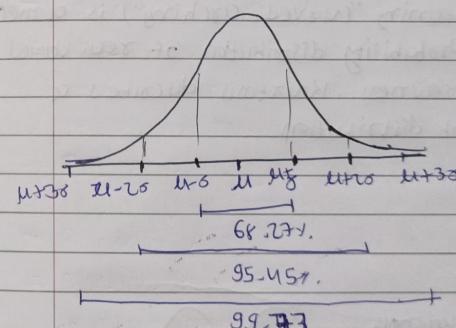
model.fit(x_train)

exponential, linear
and more

sns.Kdeplot (sample, bw_adjust=5)

OR
using
chart

★ Empirical Rule



momentum

- 1st → mean
- 2nd → std
- 3rd → skewness
- 4th → kurtosis

* Non-gaussian distribution

continuous non-gaussian distribution

- uniform distri
- Log-normal distri
- Burrto distri
- T-distri
- chi-square distri

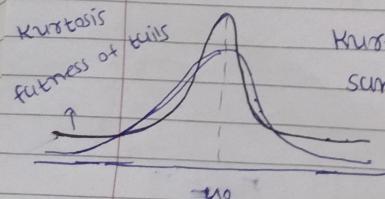
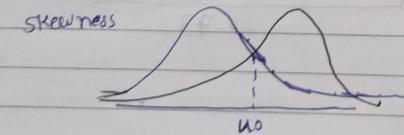
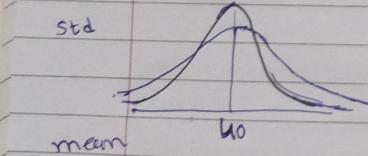
→ non-gaussian distri $\xrightarrow{\text{transformation}}$ gaussian distri

* Kurtosis

→ Kurtosis is 4th statistical moment. In probability theory and statistics, kurtosis (meaning "curved, arching") is a measure of the "tailedness" of probability distribution of real valued random variable. Like skewness, kurtosis describes particularly aspects of prob. distribution.

2007 → 100 matches

2008 → 100 matches



Kurtosis:

$$\frac{\text{sum of } \frac{(x_i - \bar{x})^4}{\sigma^4}}{n(n-1)(n-2)(n-3)} - \frac{3(n-1)^2}{(n-2)(n-3)}$$

PAGE NO.:

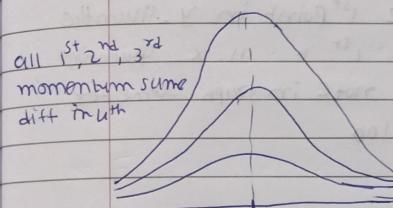
PAGE NO.:

→ Practical use-case of kurtosis

→ In finance, kurtosis risk refers to associated with possibility of extreme outcomes or "fat tails" in the distribution of returns of particular asset or portfolio.

→ If distribution has high kurtosis means there is higher likelihood of extreme event occurring either positive or negative, compared to normal distribution.

⇒ Excess kurtosis & Types



⇒ Excess kurtosis is measurement of how much more peaked or flatter distribution is compared to normal distribution, which is considered to have kurtosis of 0. It is calculated by subtracting 3 from sample kurtosis coefficient.

Types: Leptokurtic
Mesokurtic
Platykurtic

Leptokurtic: ($\text{sample kurtosis} - 3 > 0$)

A distribution with positive excess kurtosis called leptokurtic. "Lepto" means "slender" in term of shape. It has a fatter tails. This indicates that there are more extreme values.

Platykurtic: ($\text{sample kurtosis} - 3 < 0$)

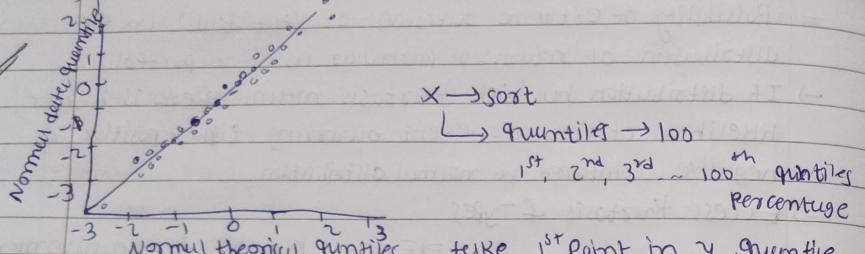
Negative excess kurtosis and fewer outliers

Mesokurtic: ($= 0$)

Normal distri (regardless of any μ, σ)

★ Quantile-Quantile Plot (Q-Q Plot)

useful to determine whether set of data follows normal dist.



$x \rightarrow$ sort

\hookrightarrow quantiles $\rightarrow 100$

1st, 2nd, 3rd ... 100th quantiles
Percentage

take 1st point in y-quantile
& 1st " in x-quantile
and note in graph sum for
all 100