

ROC and AUC

(used for binary classification)

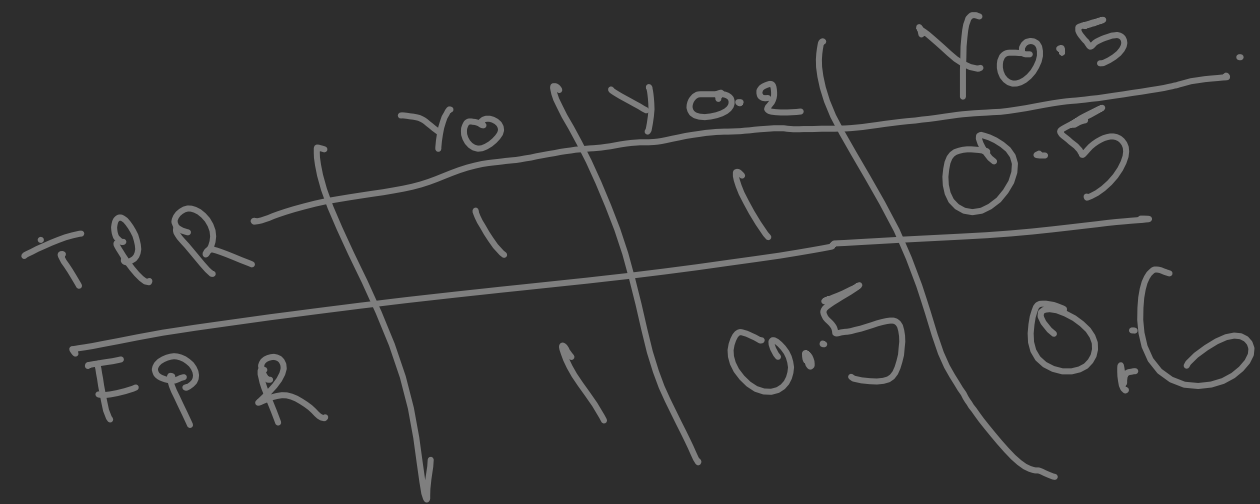
$$\text{True P. rate} = \text{TP} / \text{TP} + \text{FN}$$

used for threshold selection for

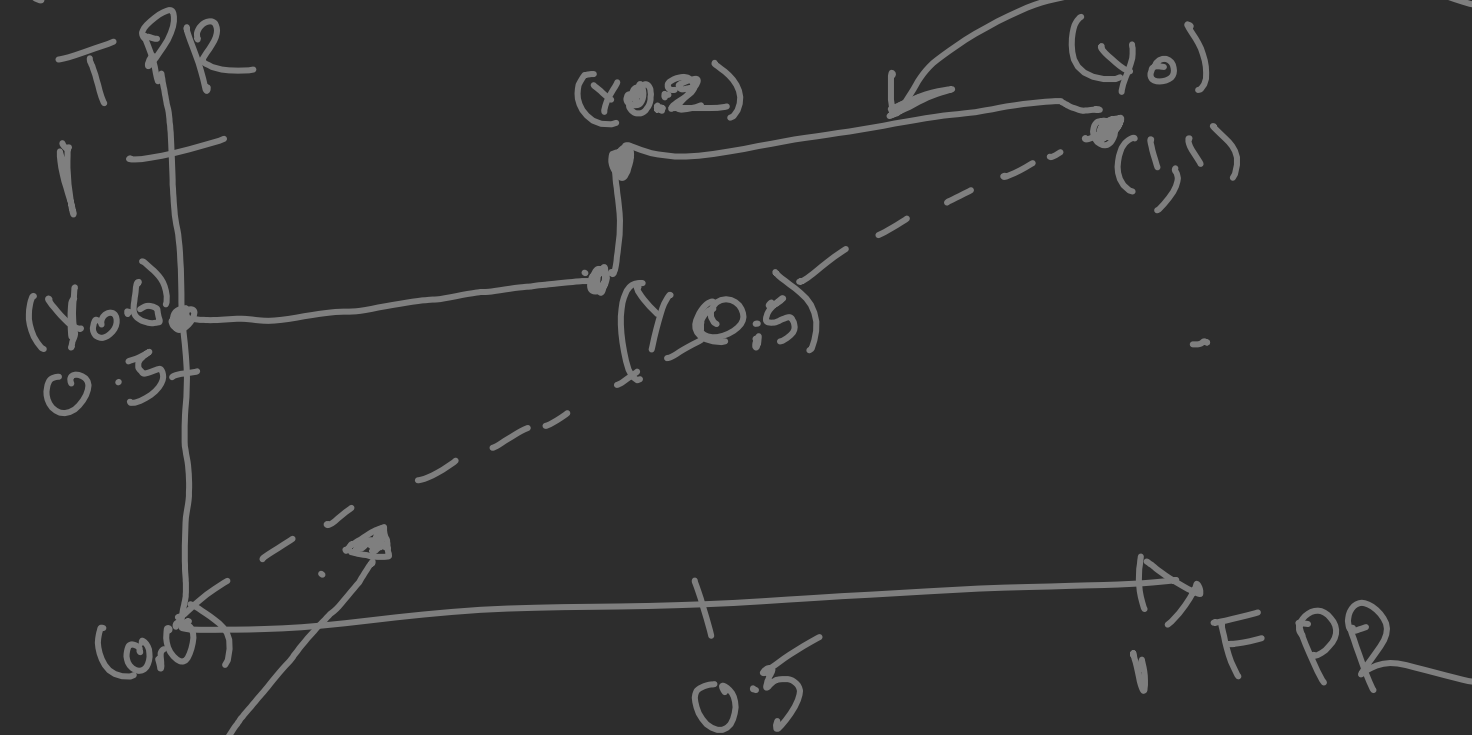
$$\text{False P. rate} = \text{FP} / \text{FP} + \text{TN}$$

also like logistic Regression

Y . \hat{Y} $\hat{Y}(0)$ $\hat{Y}(0.2)$ soon
1 0.8 1 0
0 0.9 1 0



good ROC is which do not come below dashed line



This curve is ROC

and Area under it is AUC

for	TPR only	0.6
"	FPR only	0.5
"	Mixed	0.2

Ensemble Technique

many model are used

Hyperparameter is Num. of tree / models

Bagging - it is the act of observing different models to get a unbiased result (Random forest)

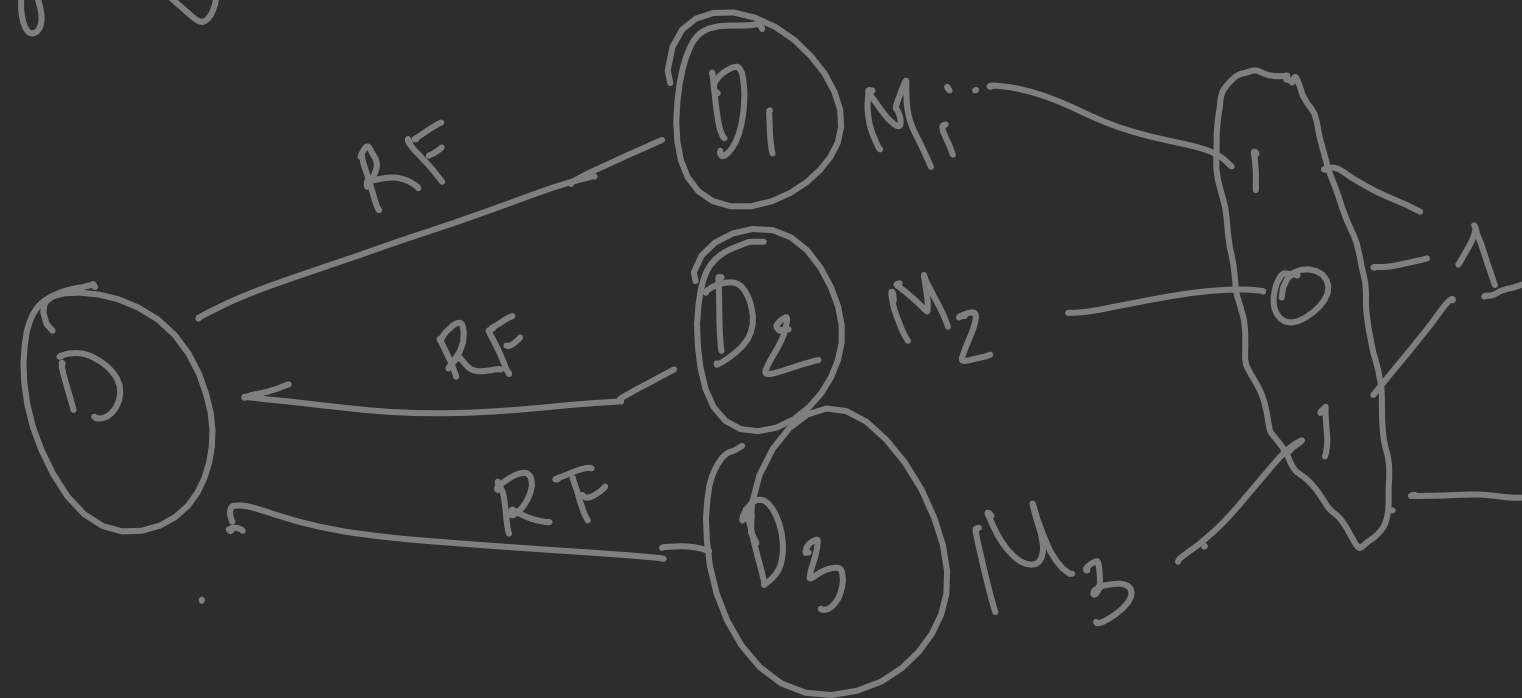
Boosting - Aggregation \Rightarrow it is

High variance
convert to

low variance

Data ($D_1 \neq D_2 \neq D_3$)
Model ($M_1 \neq M_2 \neq M_3$)

Row sampling = RF
Feature " \rightarrow



Note for regression
we just take mean of
these But in
Classification take
mode

Under sampling must be used when have large amount of data
example (nearmiss function in Imbalance Package)

Over sampling - efficient and No data loss
(SMOTETomek in I " ")

SMOTE - it is a oversampling method (minority over sample tech)
Tomek - it is a undersampling "

Curse of Dimensionality - more features / large like 100 is bad
Note every PC is right angle to each previous for ML Algo Aka High Variance

1) PCA - Principal component analysis is unsupervised ML Algo
step 1 → standard Normalize data then apply PCA

IMP ↓

data → covariance → eigen value/vector → sort wrt eigen values in Desc
take top n eigen value and vector that resemble most IMP features
n = no. " " you want

Randomize Search \rightarrow it is like grid search but take permutation
Note \rightarrow used for hyperparameter
Optimization / setting Randomly for best parameter to be found
for to used in any Algorithm for good
results.

Boosting

Base learner are created in which is sequentially increasing Accuracy
(aka models) as compare to predecessor

i) XG Boost - Binary Classifier used in Backend when using multiclass-
Classifier.

extrem gradient boost was created to deal with large data
 $\text{Residual} = \text{original} - \text{pred}$

Pruning must be of leafs not root even they are eligible for it or not

λ (lambda) - regularization parameter

Ada Boost \rightarrow the mislable data of 1st model is passed to that of second and that of second passed to third

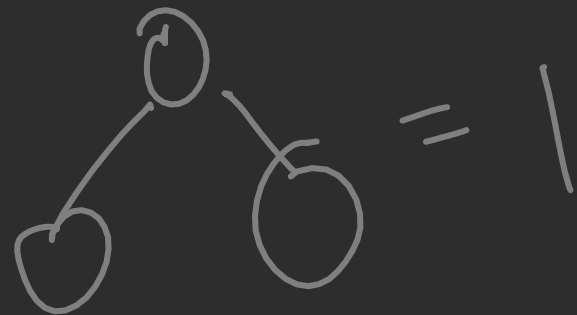
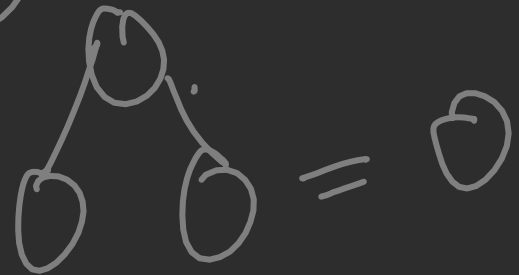
weight are change for mislabeled data by previous

$$\text{new weight} = \text{old weight} \times e^{-\text{Performance (similarity score)}}$$

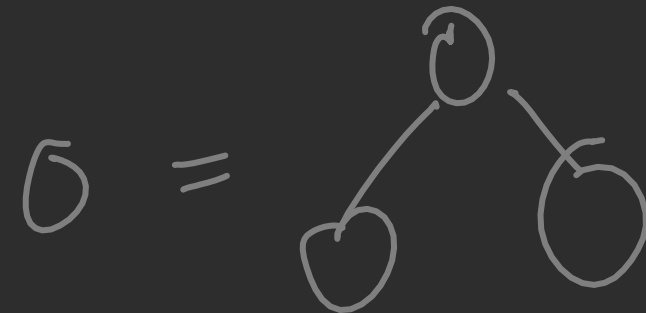
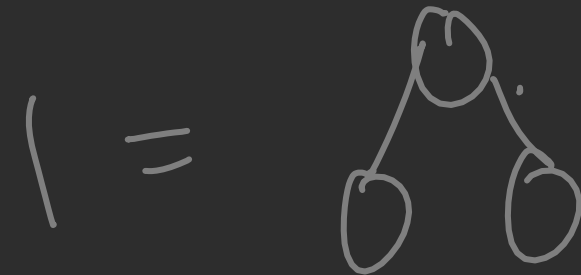
$$G_{avg} = \frac{1}{2} \log \left(\frac{1 - \text{Total error}}{\text{Total error}} \right) \leftarrow \text{row wise error}$$

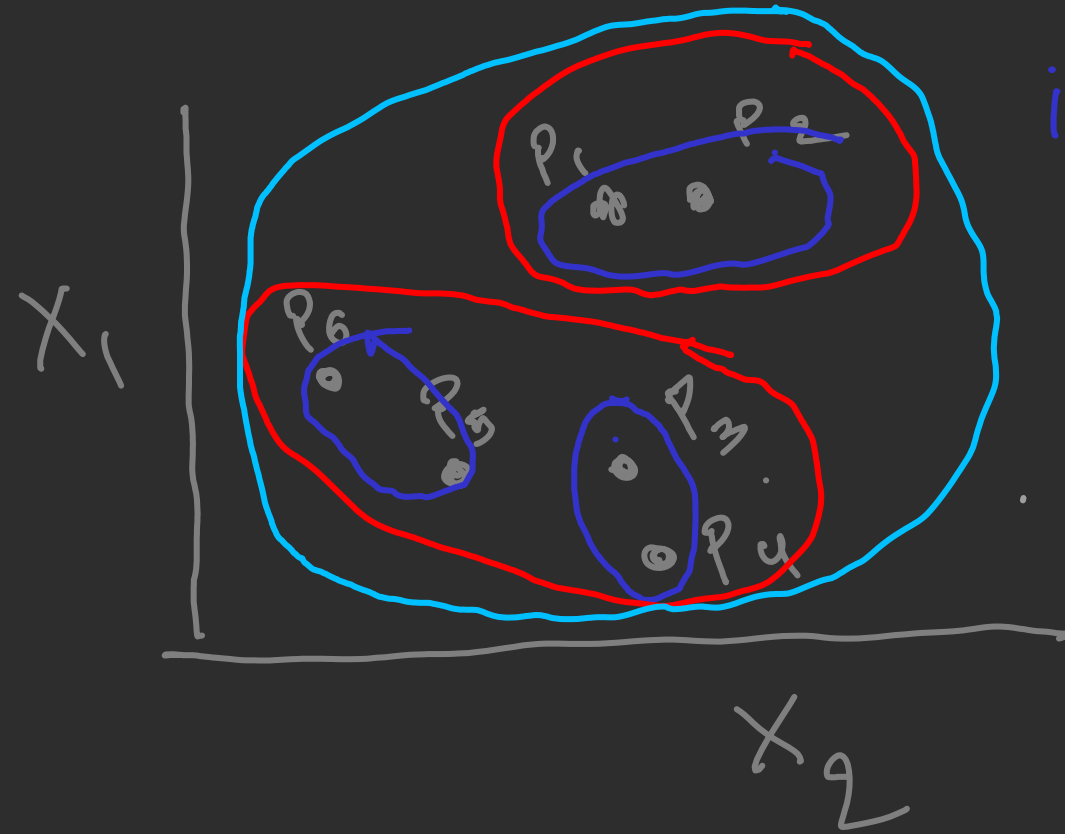
yes

result = yes

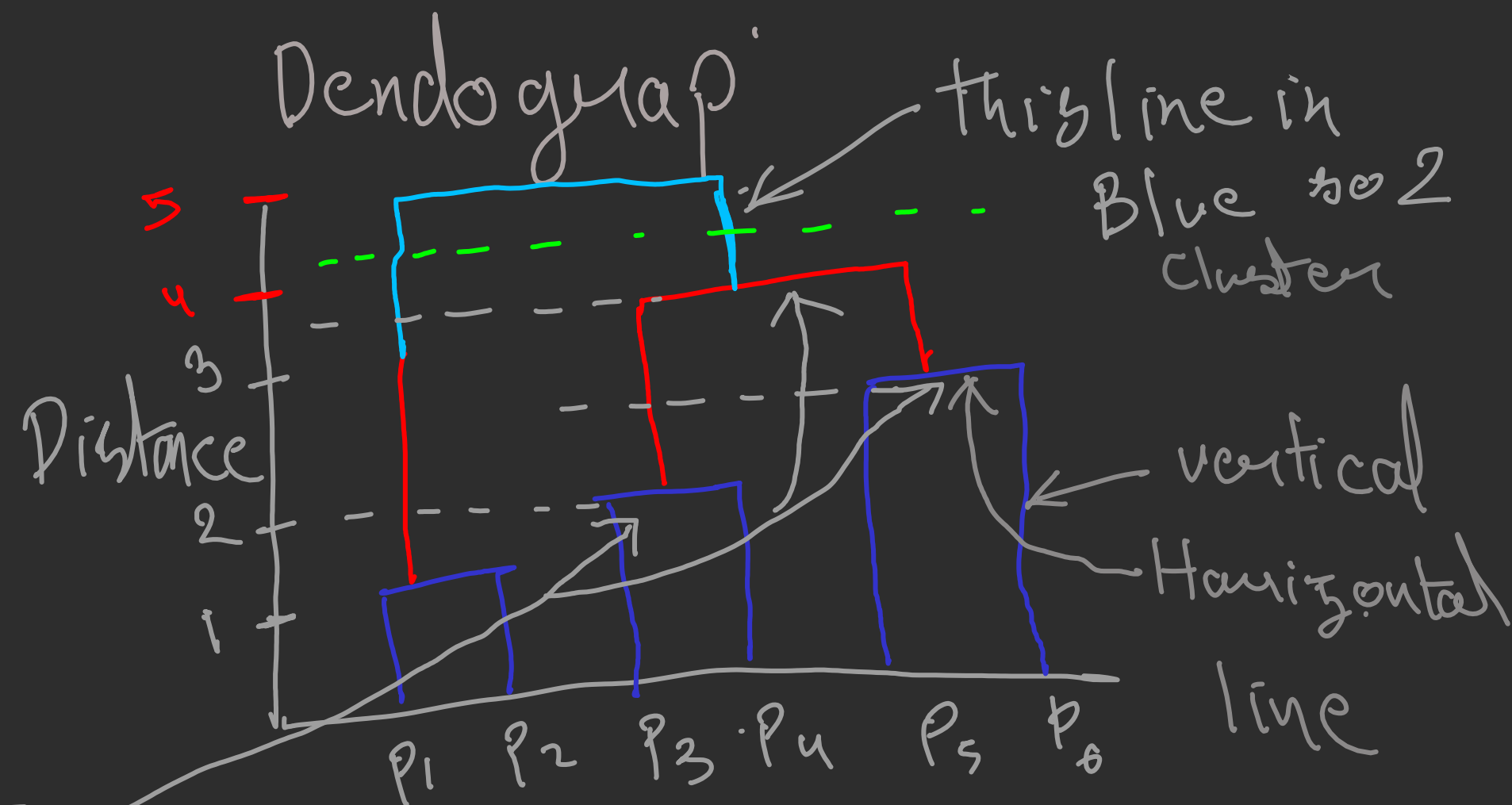


No





iteration 1 - ●
 11 2 = ●
 iteration 3 = ●
 as $6/2 = 3$



This can pass

find number of cluster?

1) longest vertical line with no horizontal line passing through it extends.

DBScan - Density Based Spatial Clustering of applications with Noise

Epsilon \rightarrow some radius around each point of data

Min Points \rightarrow number of mini point that must lie in the boundary of a point with Epsilon(ϵ)

If a point with $0 < \epsilon$ and Minpoint is a core point
if " " $0 < \epsilon$ only one core point

If a point with $0 < \epsilon$ and no point Noise point

DB work good with noise and density based cluster
and collect core point and Border point creat a cluster

* Silhouette Validation technique -

a) find the intracuster distance $\Rightarrow a$

b) find the intercluster " $\Rightarrow b$
 $\hookrightarrow \min(b)$

if $a < b$
good

else Not

Curse of Dimensionality - more features / large like 100 is bad
Note every PC is right angle to each previous for ML Algo Aka High Variance

1) PCA - principal component analysis is unsupervised ML Algo
step 1 \rightarrow standard Normalize data then apply PCA

data \rightarrow covariance \rightarrow eigen value/vector \rightarrow sort wrt eigen values in Desc
take top n eigen value and vector that resemble most IMP features
 $n = \text{no. " " you want}$

Cross Validation (CV)

Problem random state of train test split can change the accuracy of model

Solution

- 1) leave one out CV - only one point for test in a situation then move test point to next (Bad)
low bias but high variance

2) K-fold CV. - K is the number of experiment perform

let $K = 5$

let $\frac{1000 \text{ data}}{5} = 200$ frame/window with 200 stride

Iter 1 = $\frac{\text{Test}}{200} \quad 800 \text{ train}$, Iter 2 = $\frac{\text{Test}}{200} \quad 200 \quad 600 \quad \text{train}$ bo. or

then average mean with K accuracy

bad. for data with low spread of categories or imbalance

no. of A < B

	$K=1$	Train	Test	$K=2$	Train	Test
example class A		50	25		20	10
B		1000	500		1030	515

Stratified CV - consider the random sample is eliminating data imbalance first then its k-fold
like = 60:40, 50:50, 70:30 only

Time Series CV - data must include Time Series means it will use only previous data or time

Baye's theorem - if this happens then what is the probability of happening next if the prior happend
Dependent

$P(A|B) = \frac{P(A \cap B)}{P(B)}$ ← probability of A happening if B happens
 $P(B|A) = \frac{P(A \cap B)}{P(A)}$

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

→ Prior P (pointing to $P(A)$)
 → Marginal P (pointing to $P(B)$)
 → Posterior P (pointing to $P(A|B)$)
 → Likelihood P (pointing to $P(B|A)$)

Naive Bayes's Classification — feature — x_1, x_2

$$P(y|x_1, x_2) = \frac{P(x_1|y) P(x_2|y) P(y)}{P(x_1) P(x_2)}$$

→ This will be Constant (pointing to the denominator $P(x_1) P(x_2)$)

$$P(y|x_1, x_2) \propto P(x_1|y) P(x_2|y) P(y)$$

$$y = \underset{\substack{\text{from} \\ y_1 \rightarrow y_n}}{\text{argmax}} \left(P(x_1|y) P(x_2|y) P(y) \right)$$

example -
Q
Today(Sun, Hot)

$$P(\text{yes} | \text{Sun}, \text{Hot}) = \frac{P(\text{Sun} | \text{yes}) P(\text{Hot} | \text{yes}) P(\text{yes})}{P(\text{today})}$$

$$\begin{array}{cc} \nwarrow & \searrow \\ P(\text{Sun}) & P(\text{Hot}) \end{array}$$

$$\text{for yes } \text{Ans} = 0.031$$

$$P(\text{No} | \text{Sun}, \text{Hot}) = \frac{P(\text{Sun} | \text{No}) P(\text{Hot} | \text{No}) P(\text{No})}{P(\text{today})}$$

$$\begin{array}{cc} \nwarrow & \searrow \\ P(\text{Sun}) & P(\text{Hot}) \end{array}$$

$$\text{for No } \text{Ans} = 0.08571$$

Normalize

$$P(\text{yes}) = \frac{0.031}{0.031 + 0.08571} = 0.27$$

$$P(\text{No}) > P(\text{yes})$$

$$P(\text{No}) = 1 - 0.27 = 0.73$$

classified \rightarrow No

Naive Bias for NLP

Data

Sentiment analysis

This is good
This is not good

$$P(+ve | \text{not good}) = \frac{P(\text{not} | +ve) P(\text{good} | +ve) P(+ve)}{P(\text{sentence})}$$

$$P(+ve | \text{not good}) = 0.032, \quad P(-ve | \text{not good}) = 0.078$$

This	is	good	Not	+ve
1	1	1	0	1
1	1	1	1	0

$$P(+ve) = \frac{0.032}{0.032 + 0.078} = \frac{0.032}{0.110} = 0.29$$

$$P(-ve) = 0.81$$

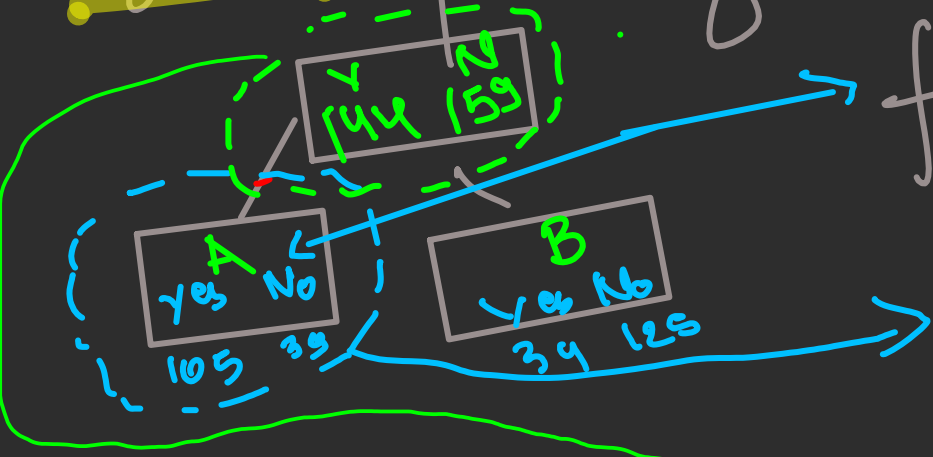
$$\text{as } P(-ve) > P(+ve)$$

The word "Not good" is not good

SVM - support vector Machine

- 1) three Hyperplane (one touch a near +ve, one for -ve, one for buffer)
- 2) marginal distance must be maximum b/w that of +ve and -ve
- 3) support vector are those lie on marginal plane (+ve, -ve)
- 4) SVM kernel used to convert lower dimension data points to higher

Gini impurity - is used to get how impure the split created by a feature is of certain node



Gini =
for leaf

$$1 - (\text{probab of yes})^2 - (\text{probab of No in node})^2$$

(lower Gini is good)

and for root node

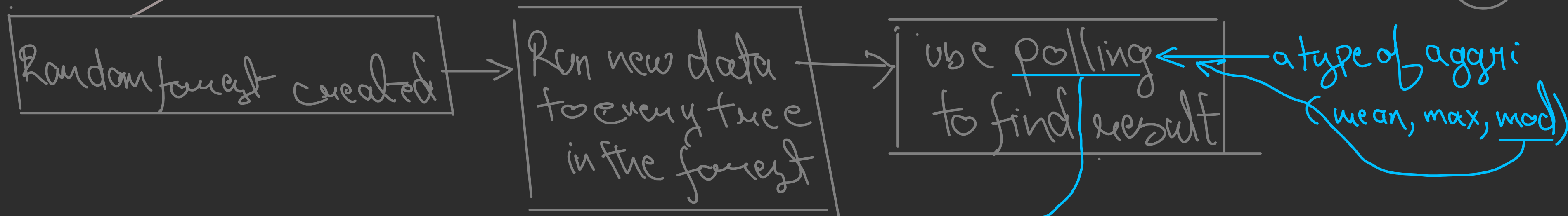
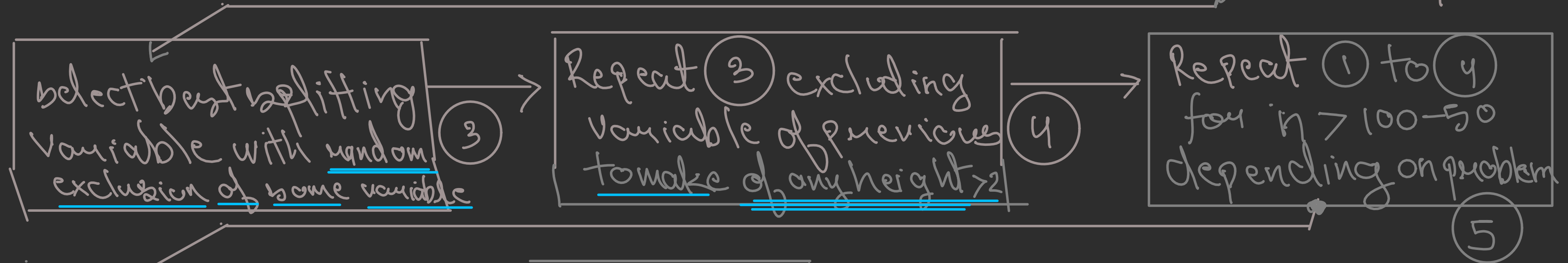
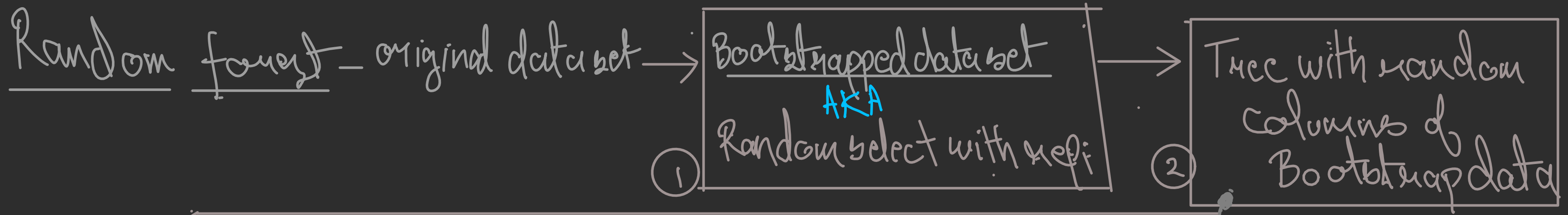
Gini for root = weighted average of Gini of leaf

Alternative of Gini

$$= \left(\frac{144}{144+159} \right) \times \text{Gini of A} + \left(\frac{159}{144+159} \right) \times \text{Gini of B}$$

Information gain not usefull as computation heavy

Entropy - measure the impurity for a feature of split
(lower the Better)

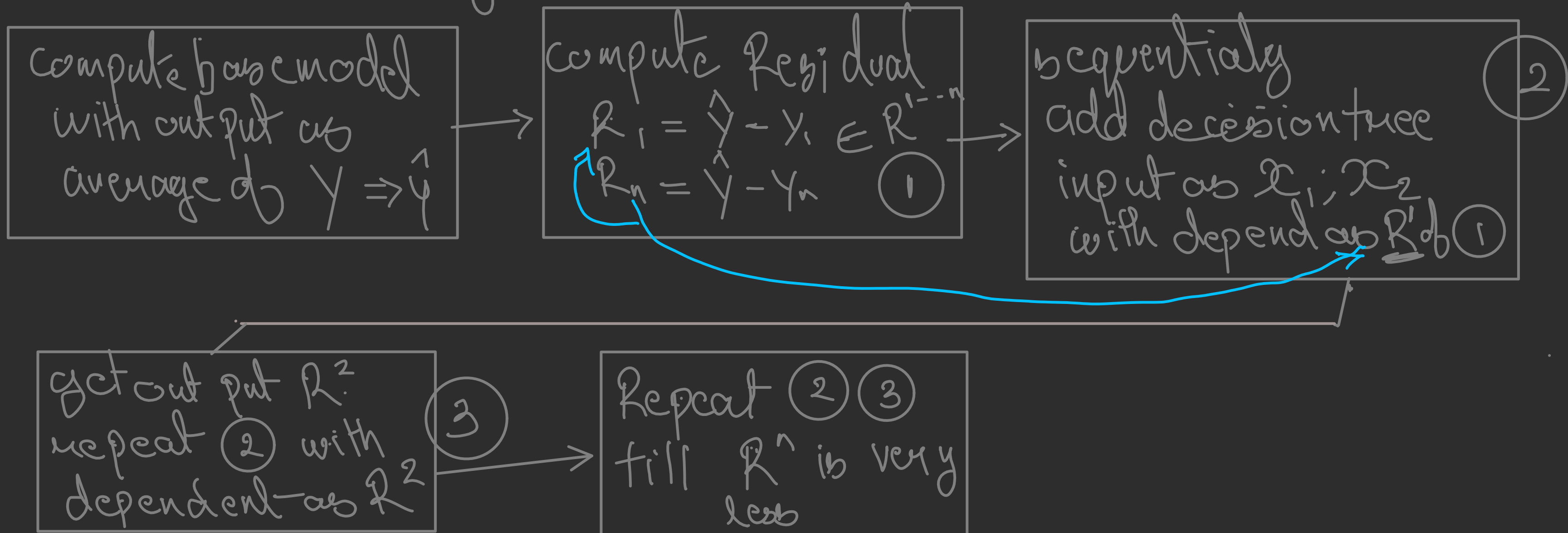


Bagging → process of ① is done then aggregating to make decision

Out of bag data → data from original that is not used in Bootstrap

when this data used to find error it is called out of bag error
this is then used to find accuracy
then tweak Hyperparameter based on forest accuracy

Gradient Boosting — $\{x_1, x_2\}$ indep, $\{y\}$ depend let

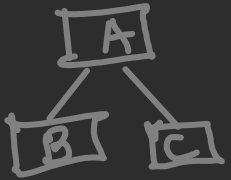


XG Boost

$\{x_1, x_2\}$ indep, $\{y\}$ depend

let \rightarrow Classification

compute base model
with output as \hat{y}
 \hat{y} - total probability = R

Construct tree 
Calculate similarity
 $= \sum R / \sum TP(1-TP)$

for A, B, C repeat ① ②
using other root

①

get output R^2
repeat ② with
dependent as R^2

③

Repeat ② ③
till R^2 is very
low

