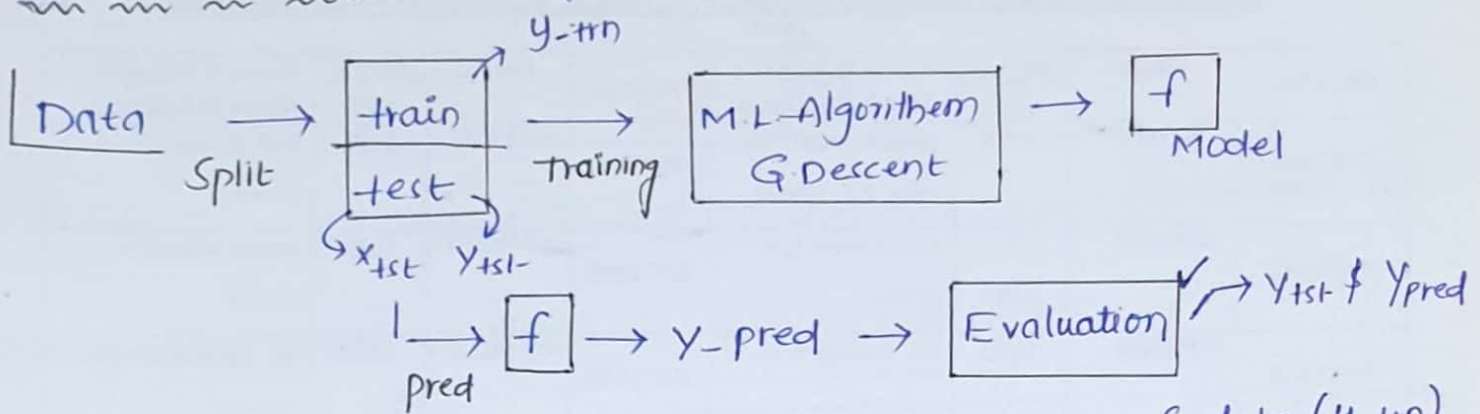


Evaluation Metrics-1 :-

Discrete \rightarrow Classification
(or)
Pipe line of Evaluation :- $\in \mathbb{R} \rightarrow$ Regression



There are two types of Evaluations based on type of data. (y-trn)

if it a 1) classification task,

Perform diff matrices — ① Accuracy

like — ② Confusion Matrices

③ Precision & Recall

④ F1 Score

⑤ ROC - AUC

⑥ Log-loss

2) Regression task we have to perform Evaluation Matrices like —

① MAE — Mean Absolute Error

② MSE — Mean Square "

③ RMSE — Root Mean " "

④ R^2 Score \rightarrow R^2 Score

⑤ MAD \rightarrow Median Absolute deviation

Regression Evaluation Matrix :-

MAE :- is nothing but measure of errors between paired Observations

like (y_{act}, y_{pred}) ie mean difference b/w $y_{act} - y_{pred}$

$$MAE = \frac{\sum_{i=1}^n |y_{act,i} - y_{pred,i}|}{n}$$

In other terms we can say — the Sum of the absolute difference between actual and predicted values.

M.S.E :- it is defined as mean or Avg of the Square of the difference b/w actual and predicted values

$$M.S.E = \frac{1}{n} \sum_{i=1}^n (y_{act,i} - y_{pred,i})^2$$

R.M.S.E :- is defined as root of mean of Square of the diff b/w actual and predicted values (or) Simply root of M.S.E.

$$R.M.S.E = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_{act,i} - y_{pred,i})^2} \text{ or } \sqrt{M.S.E.}$$

M.A.D :- is defined as median of Sum of absolute differences b/w y_{actual} to $y_{predict}$ (or) Median absolute deviation \rightarrow Error

$$M.A.D = \text{Median} | y_{act,i} - y_{pred,i} |$$

R^2 -score :- is a Coefficient of determination it is the total variance explained by the model based on correlation b/w actual & predicted

value

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

$$0 \leq R^2 \leq 1$$

Where - $SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \Rightarrow$ residual Sum of Square
 $SS_{tot} = \sum_{i=1}^n (y_i - \bar{y}_i)^2 \Rightarrow$ Total Sum of Square
 \bar{y}_i mean

Note : In SS_{tot} we get \bar{y}_i (pred) from Simple mean model

In SS_{res} we get \hat{y}_i (pred) from linear regression model

Case - I :- If we make a best-model ie where there is no errors. in that $SS_{res} = 0$

$$R^2 = 1 - \frac{0}{SS_{tot}} = 1$$

Case - II :- if a model (regression) behave like Simple mean model in that $SS_{res} = SS_{total}$

$$R^2 = 1 - \frac{SS_{tot}}{SS_{tot}} = 0$$

Case - III :- if our model contain more errors ie the model behaves like worst than Simple mean model $SS_{res} > SS_{tot}$

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = -ve$$

Case - IV :- if Simple mean model contain more errors as compare residual ie $SS_{tot} > SS_{res}$

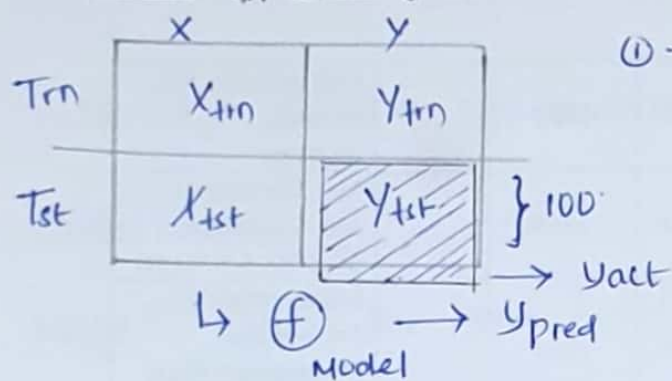
$$R^2 = \frac{SS_{res}}{SS_{tot}} = \text{lies b/w } [0, 1]$$

Classification Evaluation Matrix :-

① Accuracy :- it is the most intuitive performance measure and it is simply a ratio of correctly classified observation to total observations.

if we have high Accuracy then the model is best Model.

Accuracy Calculations :-



① After Split if Suppose test data contain 100 rows in that 60 are +ve classified & remaining 40 -ve classified

② After model making predicted values shows as ~~53~~ In positively classified data. Consist 53(+) Correctly classified and 7(-) are misclassified, Similarly in Negatively classified data. Correctly classified 35(-), miss classified (+ve) 5 Out of 40 shown below.

Actual values

~~$X_{tst} \rightarrow Y_{tst}$~~

Total Data points 100 \rightarrow 60 +ve \xrightarrow{M}
 \rightarrow 40 -ve \xrightarrow{M}

Predicted values

$X_{tst} \rightarrow Y_{tst-pred}$

Correct Miss classifier
 53 +ve, 7 -ve

35 -ve, 5 +ve

Correctly \rightarrow 88 Mis \rightarrow 12

$$\text{Accuracy} = \frac{\text{No of Correctly Classified}}{\text{Total No of data points}} = \frac{88}{100}$$

$$= 88\% \text{ or } 0.88$$

Imbalance data set :-

if we have a discrete type of data set it consist two type classifications like +ve & -ve for Suppose we have 1000 data points on that +ve type of data is 950 and -ve type of data is 50 there is huge mismatch of data. this type of data called as Imbalanced Data.

Case II - for Suppose +ve class consist 500 data points &
-ve class contain 500 data points then we called it as balanced

Data Set

Note: So if we have Imbalanced data set then we have a big problem
to predictions before going to any machine learning model we have to handle
this type of data (ie balance the data set).

Ex:- Generally we got this type of Imbalanced data sets mostly
in Health Care Organizations, BSFI, Banking Domains, Credit Card frauds

if we have a balanced data set then Accuracy will good option, in other
factors like Imbalanced data then Accuracy matrix is not good option
So that's why we go with other matrices like -

Confusion Matrix :- it is a table that is often used to describe
the performance of a model (classification) on set of test data for
which the true values are known.

if we take the same example of
100 data points set there are

correctly classified as - 88

Mis classified as - 12

if we put this data into Confusion matrix terminology.

Actual →

	+ve	-ve
Predict ↑ +ve	TP 53	FP 5
-ve ↓	FN 7	TN 35
	Mis classified	Correctly classified

Step - I - We have to divided table as two classes namely Actual and predicted in that also divide classes as binary-format ie two classes only (like - Yes/No, True/false Male/Female etc.) Confusion Matrix Only Supports to binary format.

Step - II - According to that classes we divided the classifier as shown in figure as (+ve, -ve) labels ie total positive points 60 Correctly classified ie In Actual / prediction give true value in this case 53 (+ve) put in first cell, in below cell we have put misclassified data point ie 7 (-ve) Similarly in actual data Set and predicted data (40) correctly classified (-ve) class put beside cell 35 (-ve) finally above the cell put misclassified 5 (+ve)

Step - III - Now if we observe correctly classified data points and misclassified data points Set in diagonally ie the diagonal frame of data gives us high values (Maximum) always.

Step - IV - Now giving the names to each cell according to we put the names above (TP, FN, TN, FP)

Note :- In case of Imbalanced data set Confusion Matrix are

Good option to calculate Accuracy.

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{Total (TP + FN + TN + FP)}} = \frac{35 + 53}{53 + 7 + 35 + 5} = \frac{88}{95}$$

Terminologies in Confusion Matrix :-

54

Dataset $\rightarrow D_{Test}$ values $\rightarrow 1000$ $\rightarrow y \in \{+ve, -ve\}$

$\rightarrow 900 (+ve) \xrightarrow{(N)} 850 +ve, 50 +ve$
 $\rightarrow 100 (-ve) \xrightarrow{(M)} 94 +ve, 6 -ve$

	A $\rightarrow +ve$	-ve
P $\rightarrow +$	TP 94 \uparrow	FP 50 \downarrow
-	FN 6 \downarrow	TN 850 \uparrow

Total $\rightarrow 100$; $\frac{900}{N}$
 $\rightarrow \frac{100}{P}$; $\frac{900}{N}$

T.P.R = True positive Rate = $\frac{TP}{P} = \frac{94}{100} = 94\%$
 \Rightarrow T.N.R = True Negative Rate = $\frac{TN}{N} = \frac{850}{900} = 94\%$
 F.N.R = False Negative Rate = $\frac{FN}{P} = \frac{6}{100} = 6\%$
 F.P.R = False positive Rate = $\frac{F.P}{N} = \frac{50}{900} = 50\%$

Note :- TPR is also called as Sensitivity &

TNR is called as Specificity

Precision & Recall :- These matrices are derived from Confusion Matrix

(+ve)
 (1) Precision :- one of all the points the model predicts to be positive

$$\text{Precision} = \frac{TP}{TP + FP}$$

(2) Recall :- Also called as Sensitivity (T.P.R)

	Actual \rightarrow	
	+	-
Predicted \uparrow	+	<div>TP</div> <div>FP</div> <div>Precision</div>
	-	<div>FN</div> <div>TN</div> <div>Recall</div>

Out of all the actual positive points, what percentage of them have you predicted to be positive.

$$\text{Recall} = \text{T.P.R} = \frac{TP}{P} \approx \frac{TP}{TP + FN}$$

Both precision & Recall lies in the range of $[0, 1]$

Note :- if any body asking you precision & Recall of your model both of them high and maximum is close to 1,

F1 Score :- F_1 Score also called as F Score it is measure

of model's Accuracy on dataset., F_1 Score is way of combining the Precision & recall of the model and it is defined as harmonic mean of models precision & Recall.

$$F_1 = \frac{2 \times (\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}}$$

ROC - AUC :- Stands for Receiver operating characteristic curve

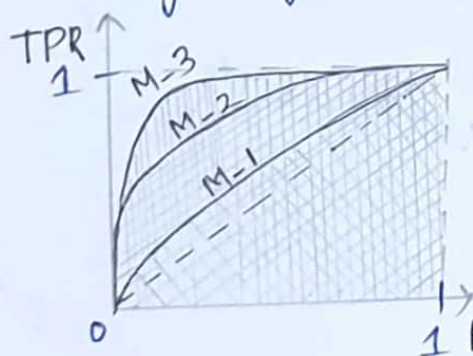
AUC Stands for Area Under the Curve.

ROC - AUC Matrix only works on Binary classification (ie if we have only two classes) if we have more than two classes this matrix is not good.

ROC :- It is an evaluation matrix for binary classification problems. It is the probability curve that plots the TPR against FPR at various threshold values and essentially separates the "Signal from the noise"

AUC :- It is measure of the ability of a classifier to distinguish between classes and is used as a summary of the ROC curve.

Note :- The higher the AUC, the better the performance of the model distinguishing b/w the positive and negative classes.



M-1 = Model 1

from graph we can say

M₃ is better than

as compare with

(M₁ & M₂)

from Confusion Matrix —

$$FPR = 1 - TNR$$

$$FNR = 1 - TPR$$

Log loss :- can also called as logarithm loss, it is a ⁵⁶ most

Important classification metric based on probabilities for any given problem. A lower logloss value means better predictions, log-loss is a slight twist on something called as likelihood function, In fact logloss is $-1 \times$ the ~~logarithm~~ log of the likelihood function.

Note :- Small logloss value gives better performance.

$$\text{logloss} = -\frac{1}{n} \sum_{i=1}^n \left\{ (y_i * \log(p_i)) + ((1-y_i) * \log(1-p_i)) \right\}$$

Short Review of Performance Metrics

Regression

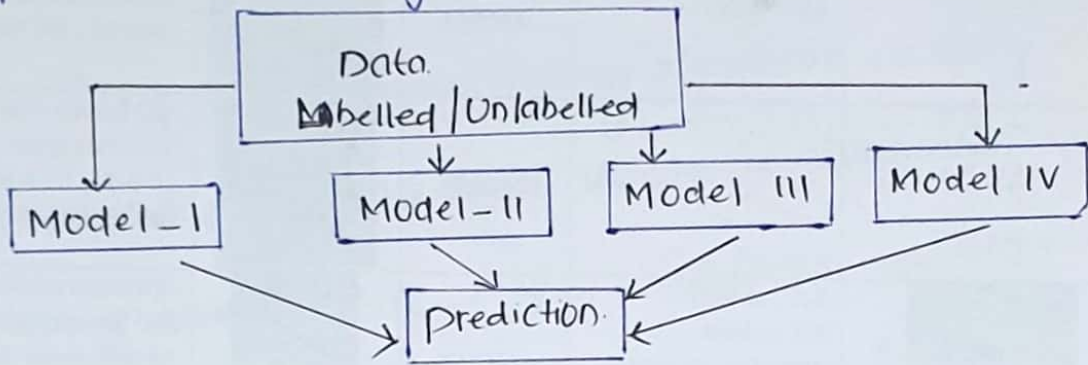
- ① M.A.E.
 - ② M.S.E.
 - ③ R.M.S.E.
 - ④ MAD.
 - ⑤ $R^2 \rightarrow$ As Near to 1
- if '1' then better performance
- Errors made by model so it is as if possible maintain low value

Classification

- ① Accuracy \rightarrow big problem of Imbalance dataset
 - ② Confusion Matrix
 - ③ Precision & Recall
 - ④ F1 score
 - ⑤ ROC-AUC \rightarrow only for binary classification.
 - ⑥ logloss \downarrow its impact with Imbalance \hookrightarrow Smaller the best (Near to '0')
- Best Metrics for Multiple classification.

Note : In classification the performance metrics (C.M, P & R, F1 score) are best metrics because it is well perform even in multiple classification as well handle the Imbalanced datasets.

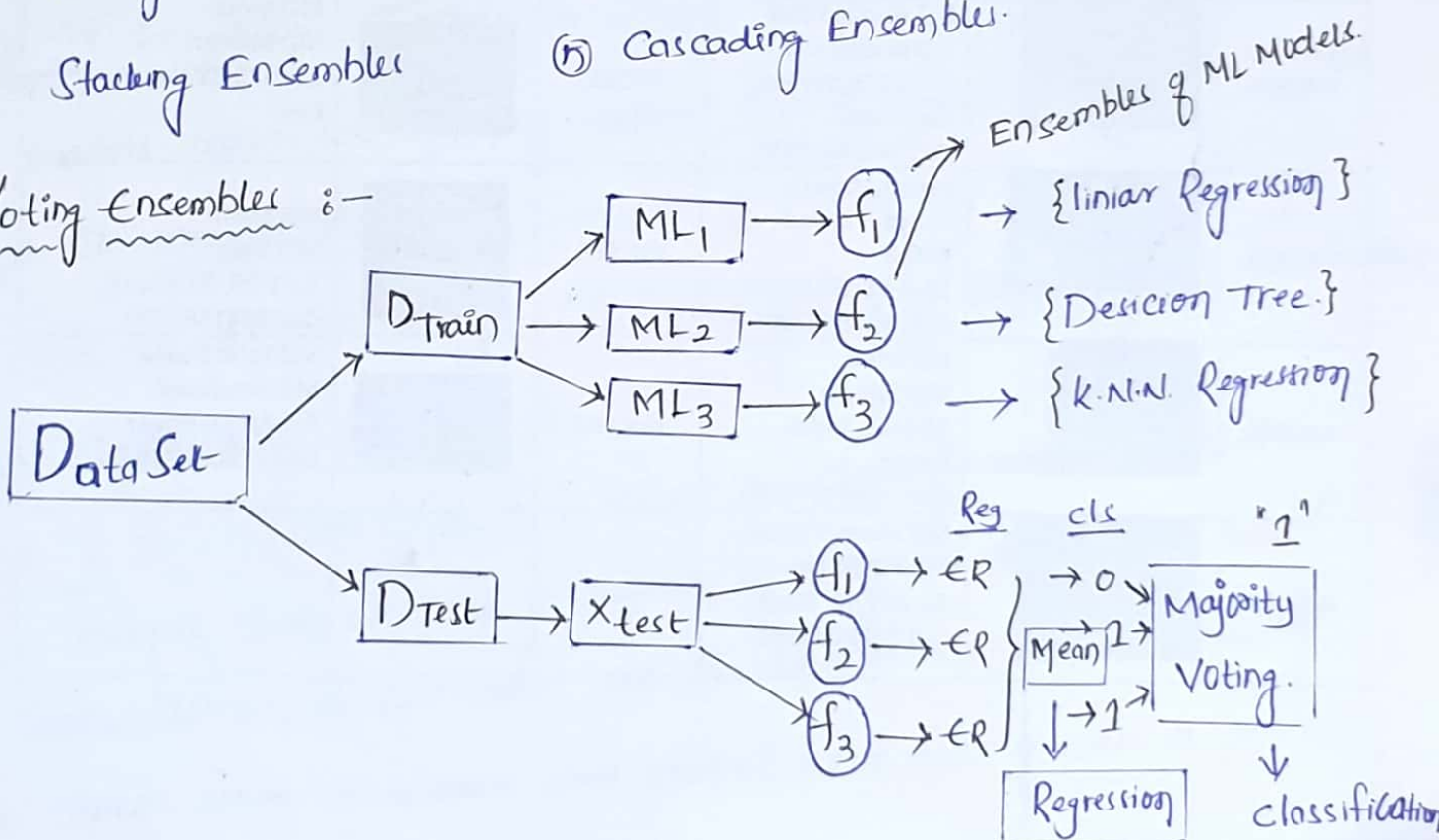
Ensemble Technique :- This technique is used in multiple learning algorithms to obtain better predictive performance than could be obtained from any of the constituent learning algorithms alone.



Classifications :- As name suggested ensemble means groups, previously for certain task completions we used one algorithm only, but here in ensembles we used different type of machine learning algorithms to perform a better result. So there are different techniques available to do this task. They are -

- ① Voting Ensembles
- ② Bagging Ensembles
- ③ Boosting.
- ④ Stacking Ensembles
- ⑤ Cascading Ensembles.

Voting Ensembles :-



Procedure for Voting Ensemble :-

(2) (3)

Step-I :- first of all we have to split a dataset in-to two parts

as $[D_{\text{Train}}, D_{\text{Test}}]$ train and test data sets

Step-II :- Apply Some machine learning algorithms as shown in diagram.

So here we decided whether our target variable is classification type or Regression type of variable. Case-I if it is classification type then do some classification algorithms apply like (Logistic Reg, KNN classification, DT classification etc) to train data set, Case-II if it is regression type of algorithm apply like (Linear Regression, KNN Regression, Decision tree regression etc)

Step-III :- if it is a case-I $\{ML-1, ML-2, ML-3\}$ ensemble of

Machine learning models as (f_1, f_2, f_3) , in case-II also we get it.

Step-IV :- Now Consider a D_{Test} Data applying some X_{test} values

in models which we applying different type of algorithms as shown in graph.

as per it we get (f_1, f_2, f_3) models of $\{ \text{Linear Reg, DT Reg, KNN Reg} \}$

in case-II as well we get (f_1, f_2, f_3) models of (Logistic Reg, DT classification, KNN classification) in case-I.

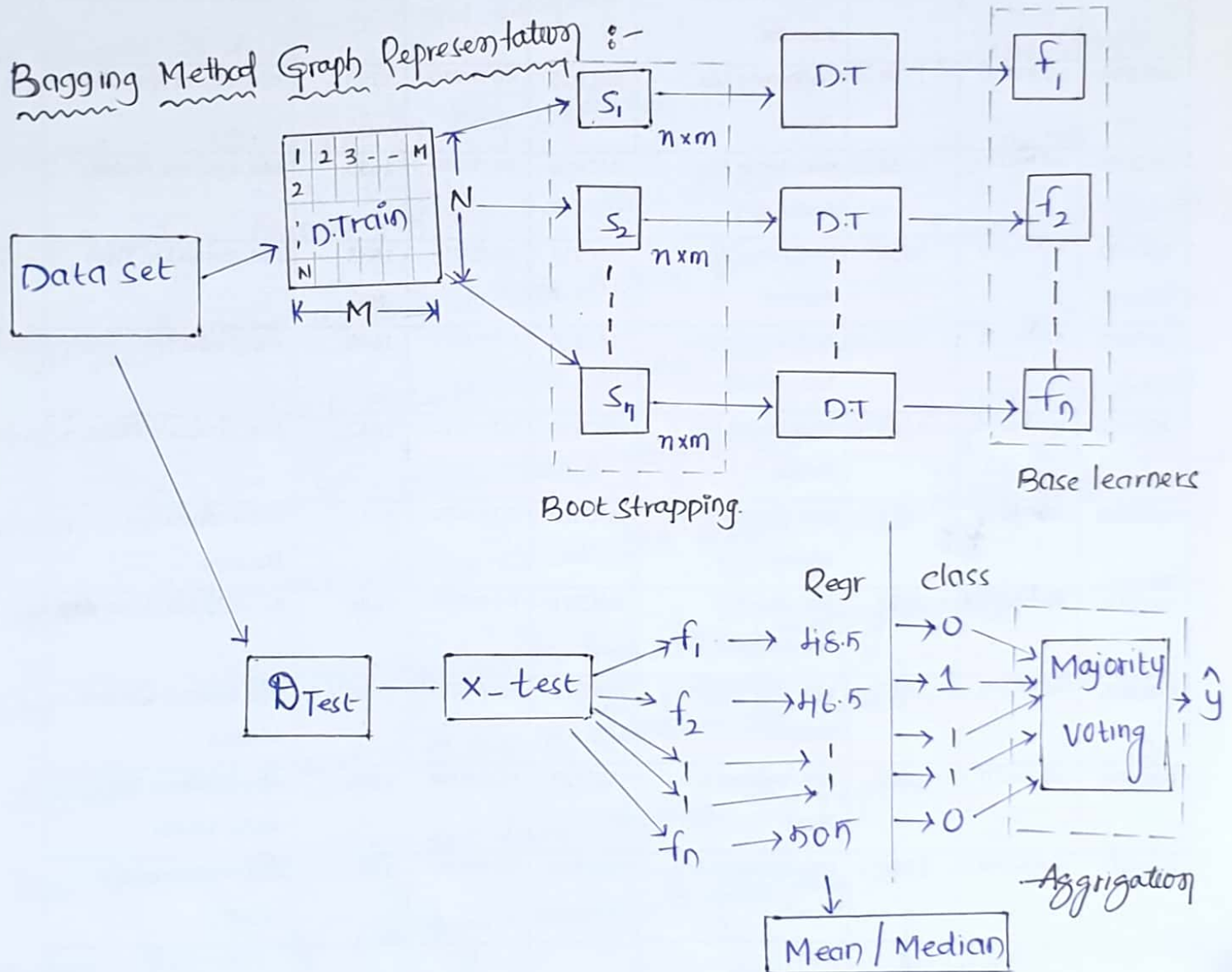
Step-V :- if our (f_1, f_2, f_3) values are classifiers and predicted as

$(0, 1)$ output then we take a majority of voting and decided the result as either (pro), if our values are Real Numbers $\in \mathbb{R}$ then we consider

this three mean (or) median and predict which type of class or prediction.

2) Bagging :- it is also called as Bootstrap Aggrigation, is a ML Ensemble, Meta-Algorithm designed to Improve the Stability and accuracy of machine learning Algorithms used in Statistical Classification & Regression. It also reduces variance and help to avoid Overfitting.

Random forest :- it is a Supervised Machine Learning algorithm, The "forest" it builds is an ensemble of decision trees usually trained with the "bagging" method. the general idea of the bagging method is that a combination of learning models increases the overall model.



Bootstrapping :- it is Nothing but Sampling of row & column.

Procedure for Bagging :-

Step-I :- First of all we have to consider a data set and split that data set into two parts as train & test Data set. Now consider a train data of a size (N, M) represent no of Columns, rows as respect.

Step-II :- Now divided this train data into different 'n' no of samples as shown in figure of a dimensionality of $(n \times m)$ represent no of Columns, rows respectively. After dividing the samples do sampling of row & columns of samples of 'n' numbers this sampling is called as "Bootstrapping".

Step-III :- Now applying only decision Tree algorithm to each sample of 'n' numbers because this algorithm make decision yields only as per this bagging technique we only consider decision-tree algorithm, Similarly in Random forest also we used decision-tree algorithm only because as per name random means select sample randomly from a trained data. and forest means trees.

Note :- In this procedure we discuss both "bagging method as well Random forest".

Step-IV :- After applying Decision-tree Machine learning algorithm we create different kind of models of size 'n' numbers. These models are called as base learners also in random forest algorithm. This algorithm do two types of tasks they are if our target variable is (ER) any real number we do regression tasks, if it is categorical type of data we do classification task. So Random forest algorithm do regression, classification tasks.

Step-V :- Now Consider test data (X_{test}) from X_{test} Data (5)

we applying to our models (f_1, f_2, \dots, f_n) if it is $\in \mathbb{R}$ any real number do regression, if it is categorical do classification task. after this we calculate mean or median to give predicted data (ie prediction), if cat then do majority voting So in this majority voting considered as that type of Data and that give prediction (\hat{y}) this represent category of data.

Note :- After Sampling when ever we apply decision tree algorithm from that step to final prediction (ie outcome) internally doing "Randomforest" process

Terminology description of Random forest :-

The models (f_i) is called as Base learner (Internally we applying "DT" Algorithm) what ever models we used here (f_1, f_2, \dots, f_n) their height (or) depths should be high as possible (ie the base learners should be over-fitted models).

Over-fitting :- is bias terminology Overfitting is nothing but high Variance and low bias, (ie the train data should be well fitted in the boundary)

Bootstrapping - (Row sampling, column sampling)

Aggrigation - classification - Majority voting; Regression - Mean (or) median.

