

Classification → labelling unlabelled data points.

KNN - K-nearest neighbour.

- ① Supervised
- ② non-parametric algorithm - doesn't take any assumption on underlying data.

③ Algorithm

- ① Select  $K$
- ② calculate euclidean distance
- ③ Take  $K$  nearest neighbour.
- ④ Count no. of data points in each category
- ⑤ Assign new data to maximum count.

- |  |  |
|--|--|
| <p><u>Pros</u></p> <ol style="list-style-type: none"> <li>① Simple</li> <li>② effective</li> <li>③ robust to the noisy data</li> </ol> | <p><u>Cons</u></p> <ol style="list-style-type: none"> <li>① High memory storage</li> <li>② need to calculate <math>K</math></li> </ol> |
|--|--|

Application

- ① Banking System
- ② Calculating credit rating

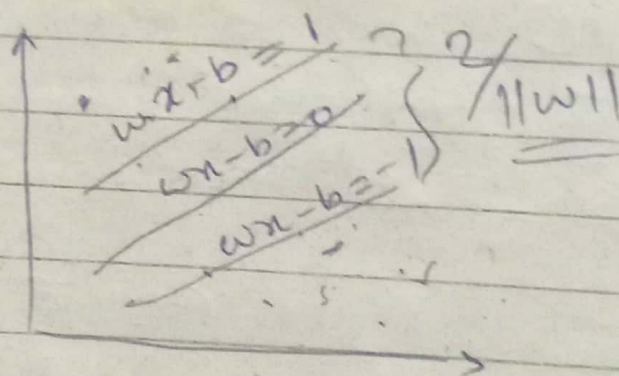
SVM - Support Vector Machine.

- It uses critical boundary datapoints to create a hyperplane that differentiates the data points.

Hyperplane - Decision Boundary.

Space  $\left\{ \begin{array}{l} 2D \rightarrow 1D, \text{ hyperplane} \\ 3D \rightarrow 2D \end{array} \right\}$

Hyperplane -  $w \cdot x - b = 0$



Kernel function

non linear data  $\rightarrow$  linear data  
low dimension  $\rightarrow$  high dimension

Linear	
Polynomial	
Sigmoid	
Radial Basis function	



## multiclass classification

{ more than 2 classes

### techniques

(1) one vs one (OVO) - heuristic method.

multi  $\rightarrow$  binary

Ex. red, blue, green & yellow.

(1) red vs blue

red vs green

red vs yellow

blue vs green

blue vs yellow

green vs yellow

$$\frac{n(n-1)}{2}$$

$$\frac{4(4-1)}{2} = 6$$

(2) one vs rest (OVR) (One vs All)

multi  $\rightarrow$  Binary

red vs [blue, green, yellow]

blue vs [red, green, yellow]

green vs [red, blue, yellow]

yellow vs [red, blue, green]



## Imbalanced Multi-class problem

→ when data points of one class is ~~to~~ much higher than other classes. Such as 100:1

Causes	Technique to handle imbalance
<ul style="list-style-type: none"> <li>↳ Sampling error</li> <li>↳ problem domain</li> </ul>	<p>(1) <u>Random Resampling</u></p> <ul style="list-style-type: none"> <li>↳ Removing sample from majority (under-sampling)</li> <li>↳ Adding sample to minority (over-sampling)</li> </ul> <p>Good choice (no info. loss)</p>

### (2) Tomek Links

↳ pair of very close instances but opposite classes

→ removing it increases the space b/w 2 classes.

→ (under sampling)

### (3) SMOTE (Synthetic Minority oversampling technique)

→ It generate artificial data for minority class



→ The Synthetic points are added b/w the chosen points and its neighbour.

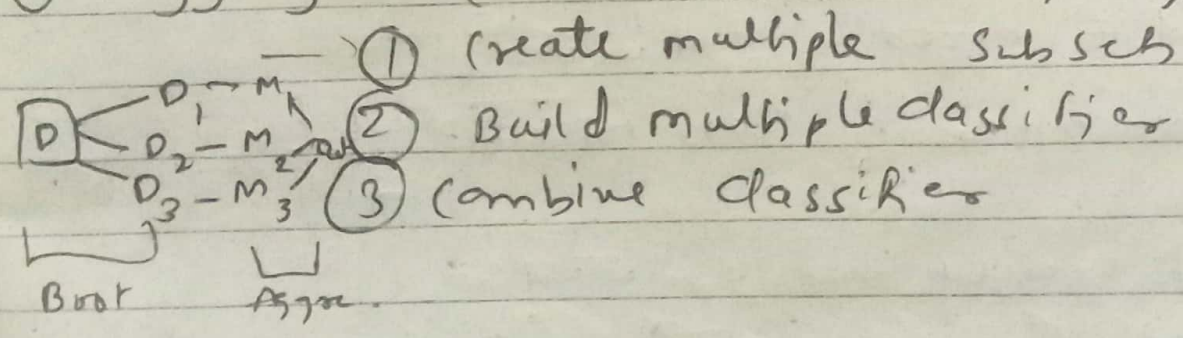
### (4) Class weights

- we can add class weight to add more weightage to the minority class and less weightage to the majority class.

## Ensemble Learning

↳ Combine prediction from two or more model.

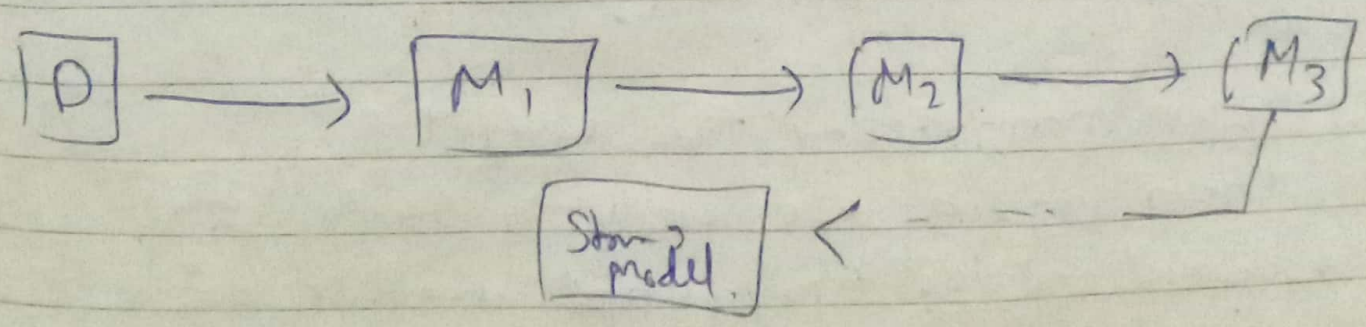
① Bagging → (Bootstrapping + Aggregation)



② Subcogging (Subsample + bagging)

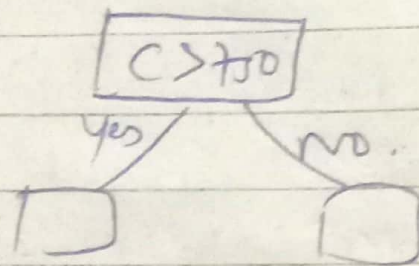
Boosting → Sequential.

→ combine weak classifier into a strong



### ③ Stumping

→ It is one level decision tree



### ④ AdaBoost (Adaptive Boosting)

Steps

- (1) equal weight apply to all data,  $(w = \frac{1}{n})$
- (2) A weak classifier is first trained.  
and  $\epsilon$  error is calculated.
- (3) performance of classifier  $\alpha = \frac{1}{2} \log_e \left( \frac{1-\epsilon}{\epsilon} \right)$

and now weight are updated

$$w_{\text{new}} = w_{\text{old}} e^{-\alpha} \quad (\text{for correct data})$$

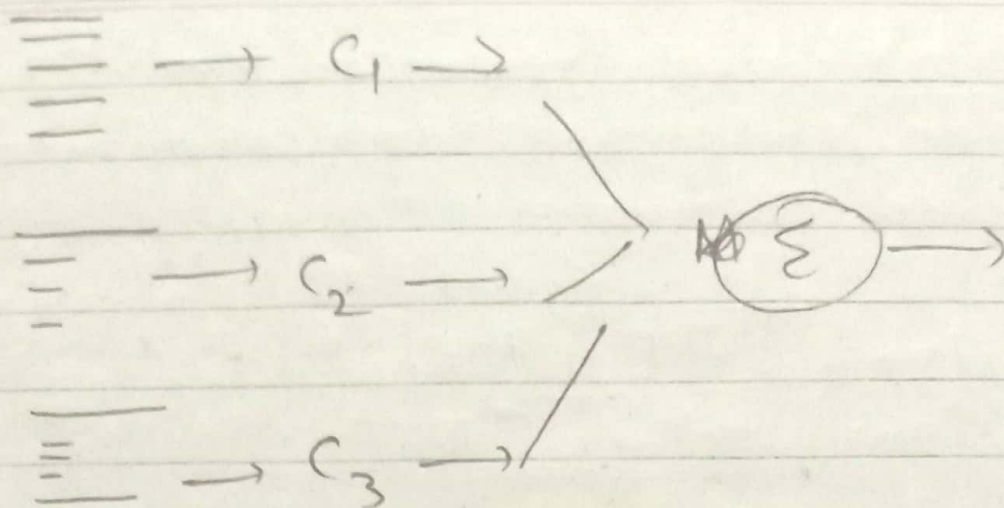
$$w_{\text{new}} = w_{\text{old}} e^{\alpha} \quad (\text{for incorrect data})$$

now, we increase the weight for incorrect data.

now, we normalized the weight:

- (4) then again we apply the same thing  
a weak classifier until we get error 0, or  
some user defined value.





## ⑤ Random Forest

→ combine several decision tree to reduce error and to build a more accurate prediction model.

→ Bagging

Steps ①  $n$  no. of random records are taken from  $k$  no. of data

② Individual decision tree is constructed for each sample

③ Each DT will generate output

④ And at last voting or averaging.

Confusion Matrix → Comparison b/w the predicted classification and actual classification various Performance Measure.

		Actual	
Predicted	Class 1	TP	FP
	Class 2	FN	TN

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN} \times 100$$

TPR, Recall, Sensitivity  $\quad TPR = \frac{TP}{TP + FN}$

FPR  $\quad FPR = \frac{FP}{FP + TN}$

TNR, Specificity   
 selecting  $\quad TNR = \frac{TN}{TN + FP}$

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

F<sub>1</sub> Score  $\quad = \frac{2TP}{2TP + FP + FN}$



⑥ micro Averaging Precision =  $\frac{TP_{total}}{TP_{total} + FP_{total}} = \frac{TP_1 + TP_2 + TP_3}{TP_1 + TP_2 + TP_3 + FP_1 + FP_2 + FP_3}$

⑦ MacroAveraging Precision =  $\frac{1}{3} (Precision_1 + Precision_2 + Precision_3)$

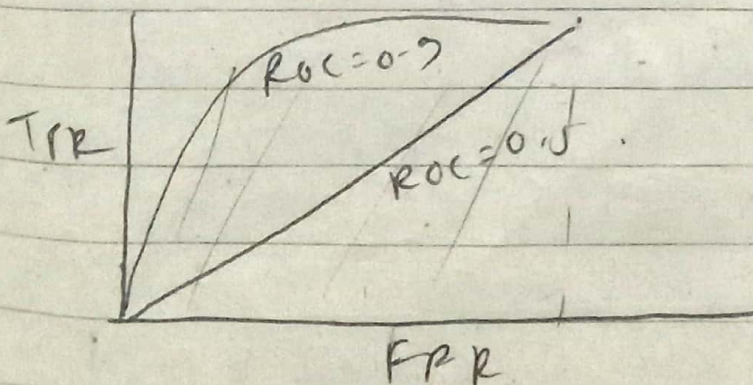
⑧ micro Averaging Recall =  $\frac{TP_{total}}{TP_{total} + FN_{total}}$

⑨ MacroAveraging Recall =  $\frac{1}{3} (Recall_1 + Recall_2 + Recall_3)$

④ Micro average f1 Score =  $\frac{TP_1 + TP_2 + TP_3}{(TP_1 + TP_2 + TP_3) + \frac{1}{2} (FP_1 + FP_2 + FP_3 + FN_1 + FN_2 + FN_3)}$

⑤ Macro Average f1 =  $\frac{f_1 + f_2 + f_3}{3}$

ROC → TPR vs FPR | AUC - Area under curve



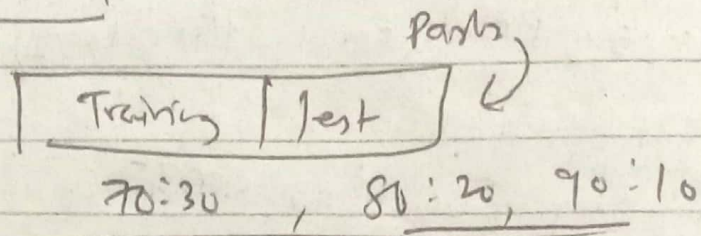
→ higher area  
better classifier

→ low area  
Poor classification



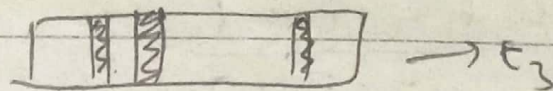
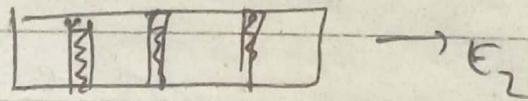
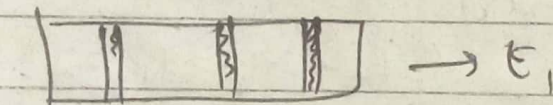
## Cross validation

① Holdout method → dividing data into 2 parts



② Repeat holdout method (iteration)

① Random Sampling → (divide testing & training set)



$$E = \frac{1}{K} \sum_{i=1}^K E_i$$

③ K fold cross validation

① divide Set into ~~2~~  $K$  Part with equal size

3 fold

iteration ①	1	2	3	4	5	6	7	8	9
②	1	2	3	4	5	6	7	8	9
③	1	2	3	4	5	6	7	8	9

training                      test



leave P out (cross validation)

P no of point out from n. total.

we train  $(n-P)$  data. and test P. points

You repeat this process for all possible combinations.

---