

	S	M	T	W	T	F	S	S
2024	1	2	3	4	5	6	7	
2025	8	9	10	11	12	13	14	
2026	15	16	17	18	19	20	21	
2027	22	23	24	25	26	27	28	
2028	29	30	-	-	-	-	-	

8 Friday  
AUGUST

2024-2028 Week 32

## KNN

- KNN — K Nearest Neighbor.
- As the name suggests it considers K Nearest Neighbors (Data points) to predict the class or continuous value for the new Datapoint.
- In this, we don't learn the weights from training data to predict output but use entire training instances to predict output.
- The distance metric and k value are two important considerations while using the KNN algorithm.
- Euclidean distance is most popular. You can also use Hamming distance, Manhattan distance, Minkowski distance as per your need.
- For using KNN in classification:  
We have to check the k-nearest neighbours and predict value for the test case.
- For using KNN in regression:  
We have to find the mean/median of the k-nearest neighbours and predict the value for test case.

→ For choosing  $k$ , we could draw error curves and pick the best one.

→ For classification<sup>(binary)</sup>,  $k$  has to be odd

→ Data pre-processing is required.

→ Code for this is as follows:

1) Find optimal  $k$  from error curves.

2) Calculate distance b/w  $(x, x_i)$  from  
 $i = 1, 2, \dots, n$

3) Sort them in ascending order

4) Take the first  $k$  values.

5) Find out which class has most frequency  
That is predicted class.

## K-Means

- K-means clustering is one of the simplest unsupervised ML algorithm.
- In this, we define a number  $k$ , which is the number of centroids i.e. the center of the cluster.
- The 'means' in k-means refers to averaging of the data; i.e. finding the centroid.
- The prediction is made based on these clusters, i.e. point is allocated to cluster through reducing the in-cluster sum of squares.
- First a group of random points are taken as centroids from each cluster, then iterative calculations are done to optimize the position of centroids.

### ALGORITHM:

Randomly initialize  $K$  centroids  $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

for  $i = 1$  to  $m$

$c^{(i)} := \text{index (from 1 to } K) \text{ of cluster centroid close to } x^{(i)}$

Cluster assign-  
ment

Moving  
centroid

for  $k=1$  to  $K$

$\mu_k :=$  average of points assigned to cluster  $k$

}

-  $\mu_{c(i)}$  = cluster centroid of cluster to which example  $x^{(i)}$  has been assigned

OPTIMIZATION:

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

$$\min_{\substack{c^{(1)}, \dots, c^{(m)} \\ \mu_1, \dots, \mu_K}} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K).$$



## DECISION

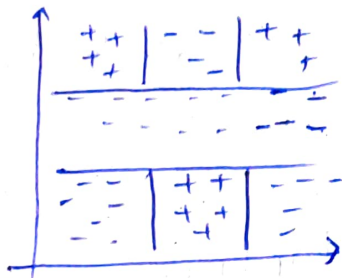
### TREES

→ Decision trees is gonna follow Greedy, Top-Down Recursive partitioning.

→ While partitioning a region  $R_p$ , we are looking for a split  $S_p$ ,

as

$$S_p(j, t) = \left( \begin{array}{l} \{x \mid x_j < t, x \in R_p\}, \quad \xrightarrow{R_1} \\ \{x \mid x_j \geq t, x \in R_p\} \end{array} \right) \quad \xrightarrow{R_2}$$



→ How to choose splits?

Define  $L(R)$  : Loss on  $R$ .

For Given  $C$  classes, define

$\hat{P}_c$  to be the proportion of examples in  $R$  that are of class  $C$ .

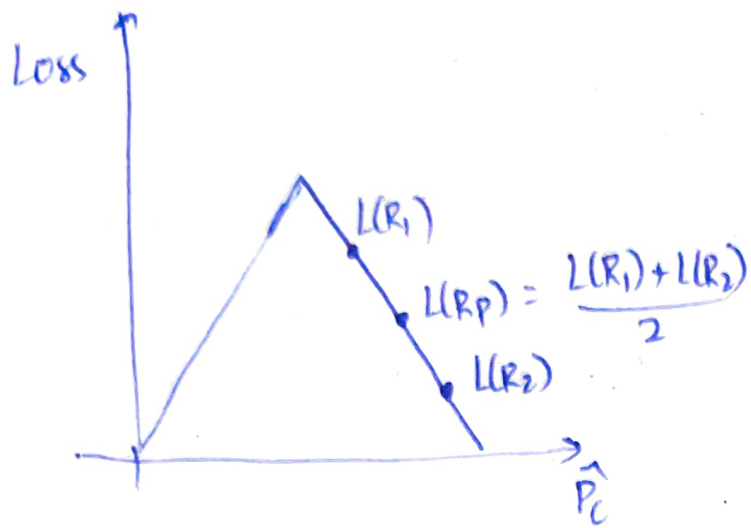
$$L_{\text{misclass}} = 1 - \max_c \hat{P}_c$$

We have to

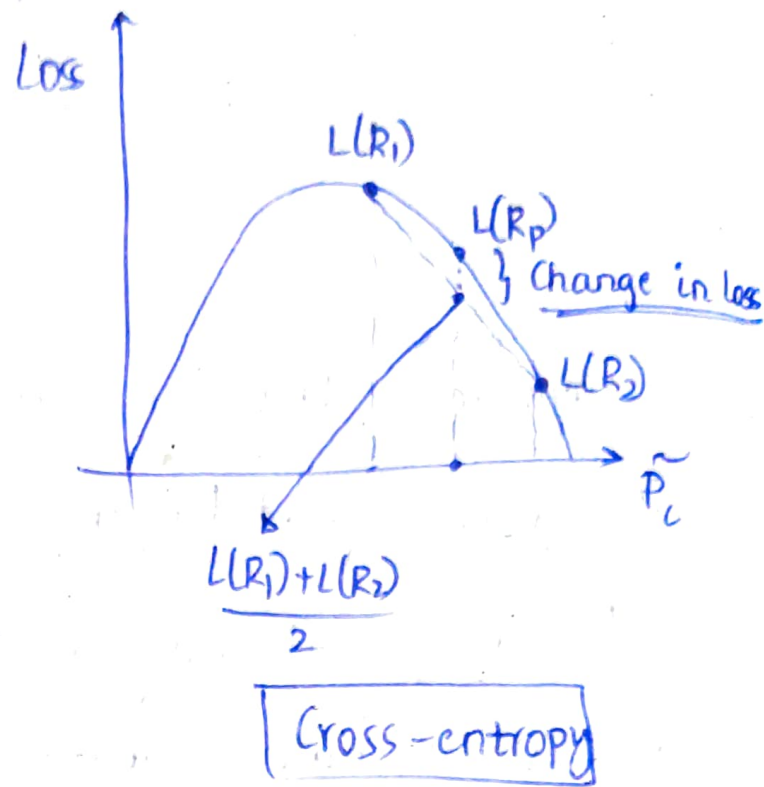
$$\max_{j, t} \underbrace{L(R_p)}_{\text{Parent loss}} - \underbrace{(L(R_1) + L(R_2))}_{\text{Children loss}}$$

\* Misclassification loss has issues, so we define cross entropy loss.

$$L_{\text{cross}} = - \sum_c \hat{P}_c \log_2 \hat{P}_c$$



✓ Misclassification



Run Time:

For  $n$  examples,  $f$  features and depth  $d$  of DT.

Test time:

$O(d)$

Train Time

Each point is part of  $O(d)$  nodes,  
Cost of point at each node is  $O(f)$

Total cost =  $O(nfd)$ .

# RANDOM

## FOREST

→ Random forest consists of a large number of individual decision trees that operate as an ensemble

### Ensembling:

Take  $X_i$ 's which are Random variables that are independent identically distributed (iid)

$$\text{Var}(X_i) = \sigma^2 \quad \text{then} \quad \text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_i X_i\right) = \frac{\sigma^2}{n}$$

If they are not independent and have a correlation constant  $p$

$$\text{then} \quad \text{Var}(\bar{X}) = p\sigma^2 + \frac{(1-p)\sigma^2}{n}$$

### Ways to Ensemble

- Different algorithms
- Different training sets
- Bagging
- Boosting

(Bagging

Ideal for bagging  
as have high var + low bias

+

DT

⇒ Random Forest

### Bagging - Bootstrap Aggregation

True population  $P$ .

and Training set  $S \sim P$

Assume  $P = S$  Bootstrap samples  $Z \sim S$ .

So, for

Bootstrap Samples  $Z_1, \dots, Z_m$

Train Model  $G_m$  on  $Z_m$

$$h(m) = \frac{\sum_{i=1}^M G_i(x)}{M}$$

→ For Random Forests;

At each split, consider only a fraction of your total features.