

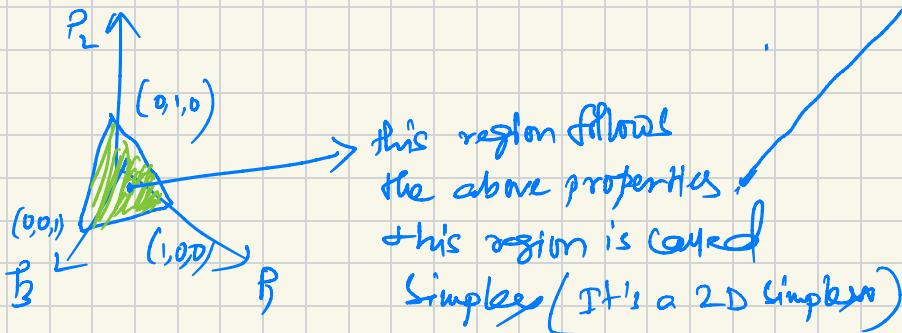
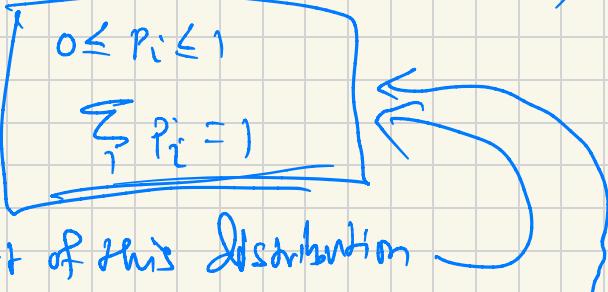
MLE can badly overfit in high dimension.

## Dirichlet distribution

$$P(\boldsymbol{\alpha} | \boldsymbol{\alpha}) = \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \alpha_1^{\alpha_1-1} \cdots \alpha_K^{\alpha_K-1}$$

Sample from probability distribution.

$$p_1, p_2, \dots, p_K \sim \text{Dir}(\alpha_1, \alpha_2, \dots, \alpha_K)$$



$$P_2 = \frac{N_c}{N}$$

Dim  $\Rightarrow C$

$$\mu_i = \frac{1}{N_c} \sum_{y_i \leq c} x_i$$

Dim  $\Rightarrow C \times D$

$$\Sigma_c = \sum_{i: y_i = c} (m_i - \mu_c) (m_i - \mu_c)^T$$

Dim =  $C \times D \times D$   
 $O(CD^2)$

so we need to estimate these many parameters.  
& its very high & if dim is large then  
no of parameters to estimate is large.

so high chance is it will overfit if data is  
less.

Can we reduce the number of parameters?  
YES, using naive Bayes rule.

It assumes that  $x_i$  &  $x_j$  are conditionally independent  
given  $y$ , for all  $i \neq j$ .

$$P(x_1, x_2, \dots, x_n | y) = \prod_i P(x_i | y)$$

Reducing the number of parameters to estimate

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

$$\begin{aligned} P(x_1, x_2, \dots, x_n | y) &= P(x_1 | x_2, x_3, \dots, x_n, y) * P(x_2 | x_3, \dots, x_n, y) * P(x_3 | x_4, \dots, x_n, y) \\ &\quad \vdots \\ &= P(x_1 | x_2, \dots, x_n, y) * P(x_2 | x_3, \dots, x_n, y) * \dots * P(x_n | y) \\ &\quad * P(y) \end{aligned}$$

Now the "naive" conditional independence assumption comes into play: assume that each feature  $x_i$  is conditionally independent of every other feature  $x_j$  for  $i \neq j$ , given  $y$ .

$$\begin{aligned} \text{so, } P(y|x_1, x_2, \dots, x_n) &\propto P(y) * P(x_1|y) * P(x_2|y) * \dots * P(x_n|y) \\ &\propto P(y) * \prod_{i=1}^n P(x_i|y) \end{aligned}$$

To estimate the parameter it requires 2D parameters.

$$\hat{y} = \operatorname{argmax}_{y_k} P(y=y_k) * \prod_i P(x_i|y=y_k)$$

$$\hat{y} = \operatorname{argmax}_{y_k} \sigma_k \prod_i \delta_{ijk}$$

MLB

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(D|\theta)$$

$$\hat{\pi}_k = \hat{P}(y=y_k) = \frac{\#D \{ y=y_k \}}{|D|}$$

$$\hat{\theta}_{ijk} = \hat{P}(x_i=x_j | y=y_k) = \frac{\#D \{ x_i=x_j \wedge y=y_k \}}{\#D \{ y=y_k \}}$$

MAP (Beta, Dirichlet prior)  $\rightarrow$  for binary classification  
 $\rightarrow$  when you have  $k$ -classes

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(\theta|D) = \underset{\theta}{\operatorname{argmax}} \frac{P(D|\theta) P(\theta)}{P(D)}$$

$$\hat{\pi}_k = \hat{P}(y=y_k) = \frac{\#D \{ y=y_k \} + (\beta_k - 1)}{|D| + \sum_m (\beta_m - 1)}$$

$$\hat{\theta}_{ijk} = \hat{P}(x_i=x_j | y=y_k) = \frac{\#D \{ x_i=x_j \wedge y=y_k \} + (\beta_{jk} - 1)}{\#D \{ y=y_k \} + \sum_m (\beta_m - 1)}$$

## NBC Algorithm

$$N_c = 0, N_{jc} = 0$$

for  $i = 1:N$  do

$c = \gamma_i$  // class i

$$N_c := N_c + 1$$

  for  $j = 1:D$  do

    if  $x_{ij} = 1$  then

$$N_{jc} := N_{jc} + 1$$

$$\hat{\pi}_c = \frac{N_c}{N}$$

$$\hat{\theta}_{jc} = \frac{N_{jc}}{N}$$

Complexity:

$O(ND)$  time.

## Naïve Bayes Classifier

- Features are conditionally independent given the class label.

$$p(\mathbf{x}|y=c, \boldsymbol{\theta}) = \prod_{j=1}^D p(x_j|y=c, \theta_{jc})$$

V.W.F.  
R

### • class-conditional density $p(\mathbf{x}|y)$

- In the case of real-valued features, we can use the Gaussian distribution:  $p(\mathbf{x}|y=c, \boldsymbol{\theta}) = \prod_{j=1}^D \mathcal{N}(x_j|\mu_{jc}, \sigma_{jc}^2)$ , where  $\mu_{jc}$  is the mean of feature  $j$  in objects of class  $c$ , and  $\sigma_{jc}^2$  is its variance.
- In the case of binary features,  $x_j \in \{0, 1\}$ , we can use the Bernoulli distribution:  $p(\mathbf{x}|y=c, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Ber}(x_j|\mu_{jc})$ , where  $\mu_{jc}$  is the probability that feature  $j$  occurs in class  $c$ .
- In the case of categorical features,  $x_j \in \{1, \dots, K\}$ , we can model use the multinomial distribution:  $p(\mathbf{x}|y=c, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Cat}(x_j|\boldsymbol{\mu}_{jc})$ , where  $\boldsymbol{\mu}_{jc}$  is a histogram over the  $K$  possible values for  $x_j$  in class  $c$ .

We know  $\hat{p} \equiv \frac{N_c}{N} \sim \text{Beta}(\alpha, \beta)$

$$P(\hat{p}|y) = \text{Beta}(N_c + \alpha, N - N_c + \beta)$$

mode of Beta distribution.

$$\text{mode} = \frac{N_c + \alpha - 1}{N + \alpha + \beta - 2}$$

when  $\alpha = 1$   
 $\beta = 1$

It's same as  
MLE which  
is uniform  
distribution

$$\text{mode} = \frac{N_c}{N} = \text{MLE}$$

when  $\alpha = 2$   $\beta = 2$

$$\text{mode} = \frac{N_c + 1}{N + 2}$$

Same as MAP.  
So, it's like Laplace  
smoothing.

Input:

Training dataset T,

$F = (f_1, f_2, f_3, \dots, f_n)$  // value of the predictor variable in testing dataset.

Output:

A class of testing dataset.

Step:

1. Read the training dataset T;
2. Calculate the mean and standard deviation of the predictor variables in each class;
3. Repeat

Calculate the probability of  $f_i$  using the gauss density equation in each class;

Until the probability of all predictor variables ( $f_1, f_2, f_3, \dots, f_n$ ) has been calculated.

4. Calculate the likelihood for each class;
5. Get the greatest likelihood;

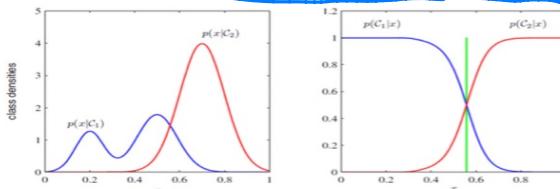
# ~~How~~ Generative model is useful for detecting outliers

## Discriminative vs Generative

- **Generative Models**: model the joint distribution  $p(x, C_k)$  directly and then normalize to obtain the posterior probabilities.

$$p(C_k|x) = \frac{p(x|C_k)p(C_k)}{p(x)}$$

- Most demanding because it involves finding the joint distribution over both  $x$  and  $C_k$ . For many applications,  $x$  will have high dimensionality, and consequently we may need a large training set. If we only wish to make classification decisions, then it can be wasteful of computational resources but can be useful for detecting outliers or novel classes.



Let  $x$  be your data and  $y$  be some sort of label or decision variable that this data can be mapped to. Discriminative models learn (some approximation of)  $p(y|x)$ .

2 Generative models, depending on who you ask, learn  $p(x)$ ,  $p(x|y)$  or  $p(y|x)$ . These three are of course related. We can get the joint from the conditional as  $p(x, y) = p(x|y)p(y)$ , and we can get the marginal from the joint as  $p(x) = \int p(x, y)dy$ . In any case, the point is that they learn something about the probability of the data.

Now let's think about what an outlier is. An outlier is an observation in your data that is unlikely. That is,  $p(x)$  is small (or  $p(x|y)$  is small, if you are able to condition on an observed value of  $y$ ). So you need some way of calculating  $p(x)$  or  $p(x|y)$ . A generative model allows you to do that, through the math I described above. In words, a generative model was trained to accurately gauge the probability of your training data (assigning higher probabilities to more common values), and so given a new piece of data a generative model can tell you how probable that value is (under the model).

For a discriminative model, this is much harder, as we have just learned  $p(y|x)$ . How are you going to get from that to  $p(x)$  or  $p(x|y)$ ? To get the latter, you need:

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

(This is basically Bayes' rule in reverse.) Note that this depends on  $p(x)$ , which you don't know. Rearranging the same formula, we find:

$$p(x) = \frac{p(y|x)}{p(x|y)p(y)} = \frac{p(y|x)}{p(x, y)}$$

Which again isn't helpful as it requires us to fill in terms on the r.h.s. that we don't know because they're not part of the model. In other words, because you never trained your discriminative model to tell you the probabilities of your data (only the probabilities of the labels *given* the data), you can't get it to do so at test time either.

For a generative model, you plug in either  $x$  or the combination of  $x$  and  $y$ , and the resulting probability, if it's small, will tell you that that particular observation is an outlier.

## Parametric Supervised Learning

A parametric algorithm has a fixed number of parameters.

A parametric algorithm is computationally faster, but makes stronger assumptions about the data. The algo may work well if the assumptions turn out to be correct, but it may perform badly if the assumptions are wrong.

A learning model that summarises data with a set of parameters of fixed size (predefined mapped function) (independent of the numbers of training examples).

No matter how much data you throw at a parametric model, it won't change its mind about how many parameters it needs.

Ex of parametric algorithm is Linear Regression, Linear SVM, Perceptron, Logistic Regression.

Pros - need few data points to learn parameters.

Cons - strong distributional assumptions, not satisfied in practice.

## ◻ parametric models :

It assumes some finite set of parameters  $w$ . Given the parameters, future predictions,  $x$ , are independent of the observed data,  $D$ .

$$P(x|w, D) = P(x|w)$$

So,  $w$  capture everything there is to know about the data.

So, the complexity of the model is bounded even if the amount of data is unbounded. This makes them not very flexible.

## ◻ Non parametric models:

It assumes that the data distribution can't be defined in terms of such a finite set of parameters. But they can often be defined by assuming an infinite dimensional  $w$ . Usually we think of  $w$  as a function.

The amount of information that  $w$  can capture about the data  $D$  can grow as the amount of data grows. This makes them more flexible.

Ex: If non-parametric algorithm is KNN, DT, ANN, SVM with Gaussian kernels.

A model is called parameteric if it has finite and fixed number of parameters. Otherwise, it is called as non-parameteric.

For example, let us say that  $x_1, x_2, \dots, x_n$  are real valued samples from some hidden distribution. Now, I can model densities in two different ways:

- Parameteric Model:** I can choose a parameteric distribution to model the density, for example, Gaussian. I can compute mean as  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$  and variance as  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$ . Thus, using these parameters we can model the density as a Gaussian distribution i.e.  $f(x) \propto \exp\left(-\frac{(x-\hat{\mu})^2}{2\hat{\sigma}^2}\right)$ . Note that, no matter how many data points are given, we just have to estimate two parameters mean and variance for this model. Thus, the number of parameters for this model are finite and fixed and independent of the number of samples available to us.

- Non Parameteric Model:** We can also choose a non-parameteric model like Kernel Density Estimator.

Kernel Density Estimator defines density as:  $f(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$  where  $h$  is the bandwidth parameter. Note that, the samples are directly part of the model which means they are the parameters of the model (different samples will give rise to different models). The number of parameters of this model grow linearly with the number of samples and theoretically, they could be infinite in number.

# Non Parametric Density Estimation Methods in Machine Learning

Read Discuss

Non-parametric methods: Similar inputs have similar outputs. These are also called instance-based or memory-based learning algorithms. There are 4 Non – parametric density estimation methods:

1. Histogram Estimator
2. Naive Estimator
3. Kernel Density Estimator (KDE)
4. KNN estimator (K – Nearest Neighbor Estimator)

## Histogram Estimator

It is the oldest and the most popular method used to estimate the density, where the input space is divided into equal-sized intervals called **bins**. Given the training set  $X = \{x^t\}_{t=1}^N$  an origin  $x_0$  and the bin width  $h$ , the histogram density estimator function is:

$$\hat{p}(x) = \frac{\#\{x^t \text{ in the same bin as } x\}}{Nh}$$

Histogram estimator

The density of a sample is dependent on the number of training samples present in that bin. In constructing the histogram of densities we choose the origin and the bin width, the position of origin affects the estimation near the boundaries.

if  $\Delta$  too small density model is too spiky

if  $\Delta$  too large density model is too smooth.

\* model complexity is decided by its bin size.

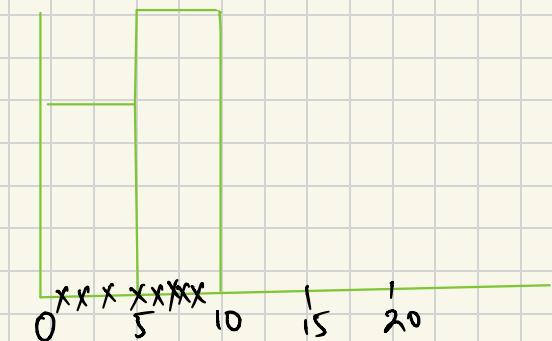
## Ex of Histogram Estimator

Let's say values are

[2, 3, 5, 7, 9, 4, 6, 8]

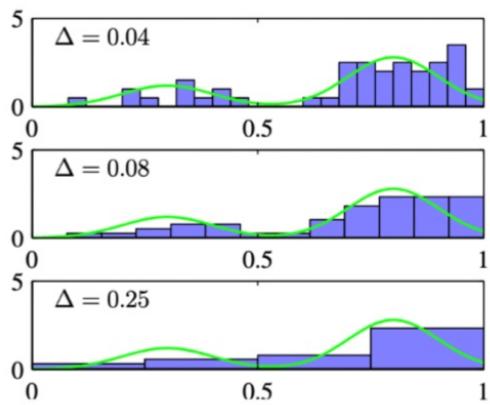
origin = 0

bin-width = 5



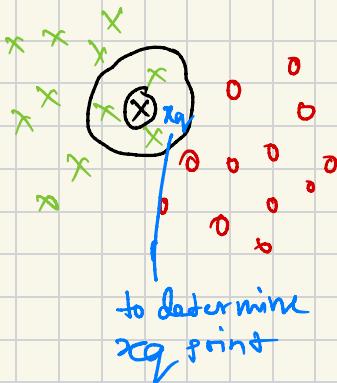
$$p(x_i) = \frac{\text{# of observations belong to same bin } i}{N \Delta i}$$

$\Delta i$  = bin-width



# KNN [Classification and Regression]

$$D =$$



$$D = \{(x_i, y_i) \mid x_i \in \mathbb{R}^n, y_i \in \{0, 1\}\}$$

steps:

1. find K-nearest points to  $x_q$  in  $D$ . let  $K=3$  here.

$$\{x_1, x_2, x_3\} \text{ 3NN to } x_q$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$y_1, y_2, y_3$$

2.  $\{y_1, y_2, y_3\} \rightarrow$  majority

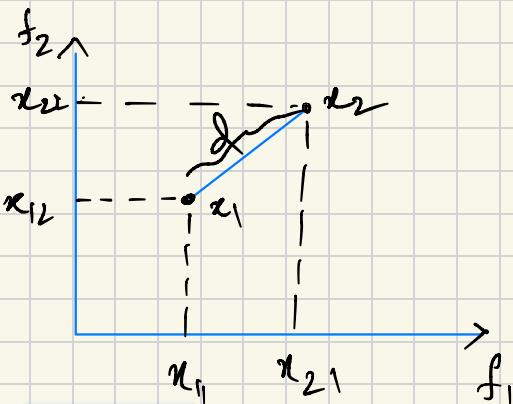
$$\{+, +, +\} \quad y_q = '+'$$

$$\{+, +, -\} \quad y_q = '+'$$

$K$  should be odd no so that no draw happens.

Distance measures:

## ① Euclidean Distance



$$d = \sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2}$$

$$= \|x_1 - x_2\| \quad \begin{cases} \text{length of} \\ x_1, x_2 \end{cases}$$

$$x_1 \in \mathbb{R}^d$$

$$x_2 \in \mathbb{R}^d$$

$$\|x_1 - x_2\| = \left( \sum_{i=1}^d (x_{1i} - x_{2i})^2 \right)^{1/2}$$

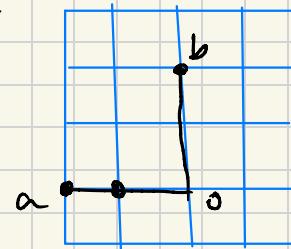
↓  
L<sub>2</sub> norm of a vector.

$\|x_1\|$  = Euclidian distance of  $x_1$  from origin  $\left( \sum_{i=1}^d x_{ii}^2 \right)^{1/2}$

## ② Manhattan Distance :

$$\sum_{i=1}^d |x_{1i} - x_{2i}|$$

↑      absolute value



$\|x_1 - x_2\|_1$ , L<sub>1</sub> norm of a vector.

$$\begin{aligned} \text{dist}(a, b) &= \text{dist}(a, 0) + \\ &\quad \text{dist}(0, b) \end{aligned}$$

Generalization to L<sub>1</sub> & L<sub>2</sub> norm —

## ③ L<sub>p</sub> norm / Minkowski dist :

$$\|x_1 - x_2\|_p = \left( \sum_{i=1}^d |x_{1i} - x_{2i}|^p \right)^{1/p}$$

if  $p=2$  then Minkowski dist = Euclidian dist

if  $p=1$  then Minkowski dist = Manhattan dist.

4

## Hamming Distance:

XOR between two Boolean vectors.

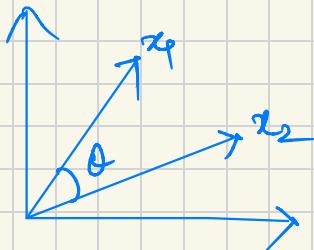
Applicable to strings also,

$$x_1 = \text{a b } \boxed{\text{c}} \text{ a d e f } \quad \text{h i k}$$
$$x_2 = \text{a c } \boxed{\text{b}} \text{ a d e g } \quad \text{h i k}$$

$$\text{hamming dist}(x_1, x_2) = 4.$$

5

## Cosine-Similarity:-



$$\cos(\theta) = \frac{x_1 \cdot x_2}{\|x_1\|_2 \|x_2\|_2}$$

if  $x_1, x_2$  is unit vector  
then  $\cos(\theta) = x_1 \cdot x_2$

6

## Edit Distance:-

- The minimum edit distance between two strings.
- 3 operations allowed (Insert, Delete, substitute)
- Needed to transform one into the other.

Rx:

INTE\*NTION  
|| || || | | |  
\*EXECUTION

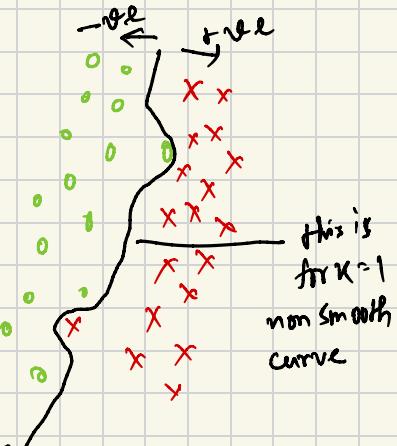
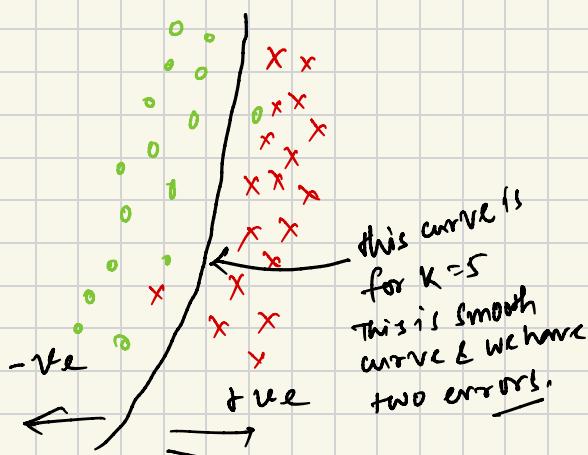
d s s i s

So, distance between these two is 5.

How to decide K?

[K] in KNN is hyperparameter.

geometric intuition if  $K=1$ ,  $K=5$



Such curves are called decision surfaces.

If  $K$  is too large, let  $K=n$

then if  $n_1 > n_2$  [ $n_1 + n_2 = n$ ]

In this case, any query point  $y_m$  take, it will be +ve. because you have to consider all  $n$  points & among  $n_1 < n_2$ ,  $n_1 > n_2$ . so new print will be assigned to  $n_1$  side.

### Overfit

- It tries to cover noisy points also
- non-smooth
- no mistakes

$$K=1$$

### Underfit

- Imperfect
- It leaves about majority prints. every points belongs to majority class.
- Imperfect

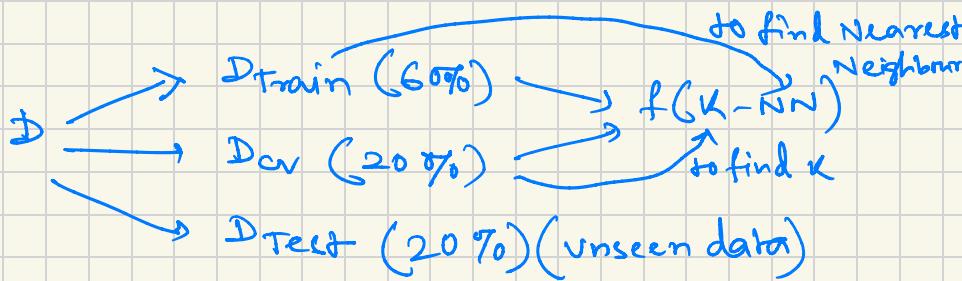
$$K=n$$

### Well fit

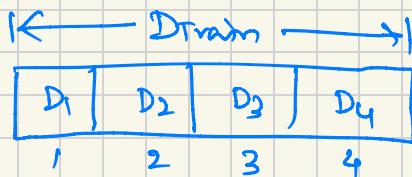
- less prone to noise
- robust
- $K=5$
- Smooth.

## How to determine K?

① we will devide the whole dataset into 3 parts.



②

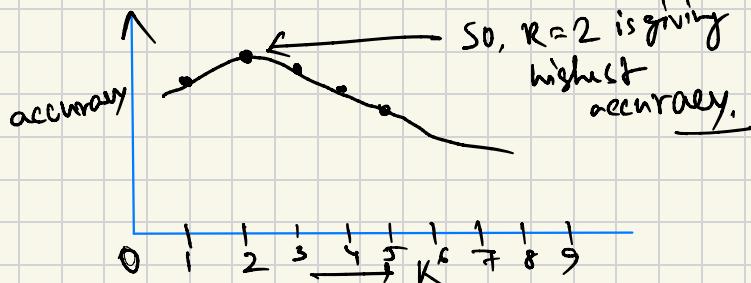


→ randomly divide into 4 parts with equal size. (4-Fold)

③

	Train <sub>m</sub>	D <sub>cv</sub>	Accuracy
K=1	D <sub>1</sub> D <sub>2</sub> D <sub>3</sub>	D <sub>4</sub>	a <sub>1</sub>
K=1	D <sub>1</sub> D <sub>2</sub> D <sub>4</sub>	D <sub>3</sub>	a <sub>2</sub>
K=1	D <sub>1</sub> D <sub>3</sub> D <sub>4</sub>	D <sub>2</sub>	a <sub>3</sub>
K=1	D <sub>2</sub> D <sub>3</sub> D <sub>4</sub>	D <sub>1</sub>	a <sub>4</sub>
K=2			a <sub>1</sub>
K=2			a <sub>2</sub>
K=2			a <sub>3</sub>
K=2			a <sub>4</sub>

Considering 4-fold cross-validation we need to plot average accuracy.



## Test/Evaluation Time & Space complexity:

Input:  $D_{\text{train}}, k, x_q \in \mathbb{R}^d$

Output:  $y_q$   $\downarrow$   $k$  is generally small, its 5 or 10.

$KNNpts = []$

for each  $x_i$  in  $D_{\text{Train}}$ :

Compute  $d(x_i, x_q) \rightarrow d_i$

Keep smallest  $k$  distance  $(x_i, y_i, d_i)$

Count-prs = 0 ; Count-neg = 0

$\rightarrow n$  points,  $d$ -dim each  
so, time complexity  
 $\geq O(nd)$

$\rightarrow$  considering  
 $k$  is very small  
it will be  $O(1)$

for each  $x_i$  in  $KNNpts$ :

if  $y_i$  is +ve :

Count-prs + 1

else Count-neg + 1

if Count-prs > Count-neg:

return  $y_q = 1 (+ve)$

else:  
return  $y_q = 0 (-ve)$

Time Complexity  
 $= O(nd)$

Space Complexity

$\geq O(nd)$  to keep  
training dataset

## \*\* Limitations of KNN:

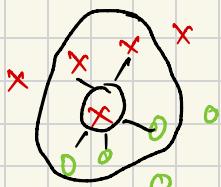
It takes huge time & space. To reduce it we have some techniques

→ KD-tree

→ LSH (Locality sensitive Hashing)

## Weighted KNN:

Instead of majority votes, sometimes we can go with shortest distance.



Let  $K=5$   
for given  $x_q$ ,

The way to assign  
weights is

$$w_i = \frac{1}{d_i}$$

$d_i \uparrow \rightarrow w_i \downarrow$   
 $d_i \downarrow \rightarrow w_i \uparrow$

$$\begin{aligned} \text{Sum} &= 1.5 \\ (1+0.5) & \end{aligned}$$

$$\begin{aligned} \text{Sum} &= 1.75 \\ (1+\frac{1}{2}+\frac{1}{4}) & \end{aligned}$$

Here,  $1.5 \gg 1.75$  so, it will be assigned  
to -ve class.

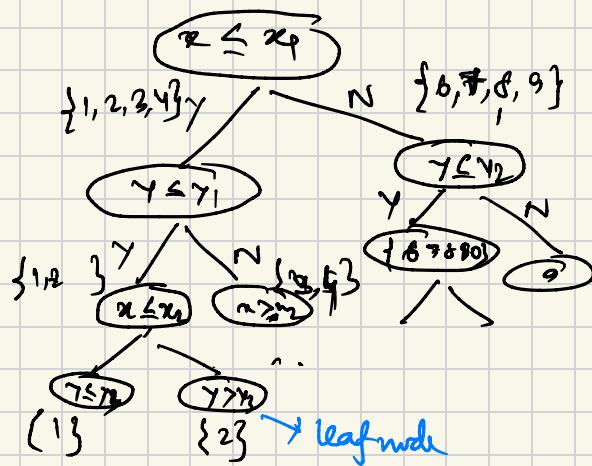
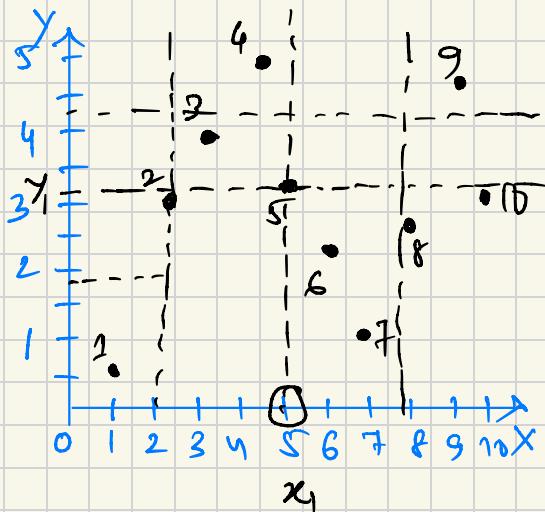
These points are  
closest to  $x_q$ ,  
but if we consider  
majority votes,  
it will be assigned  
to +ve. Now  
need to consider  
distance also?  
Why only majority  
votes?  
To eliminate this  
we can go with  
Weighted KNN.

## How to implement kd-tree using BST?

BST is  $\rightarrow$  1D [Scalar]

Extending BST to K-dimensional tree  
is called Kd tree.

Step 1: Pick X-axis, project points into X-axis, Compute median, split data into median.



You alternate each axis's one after another till leaf nodes.

so, RD-tree boxes the space using axis parallel lines/ planes into rectangle(2D)/ cuboids(3D)/ hypercubrid( $n$ D).

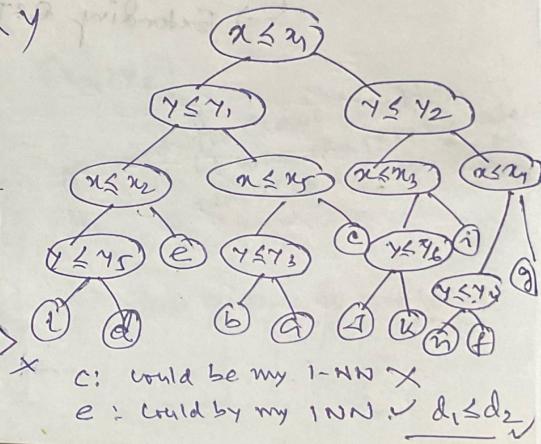
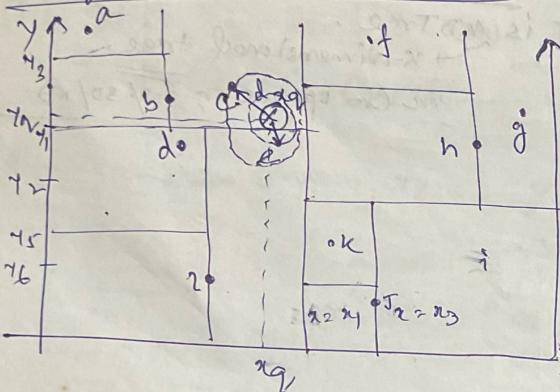
In 2D  $\rightarrow$  axis parallel lines  $\rightarrow$  rectangles

In 3D  $\rightarrow$  axis parallel planes  $\rightarrow$  cuboids.

In  $n$ D  $\rightarrow$  axis parallel hyperplane  $\rightarrow$  hypercubrid.

using KD-tree to find INN:

→ using kd tree to find INN:



# of comparison to find 1-NN

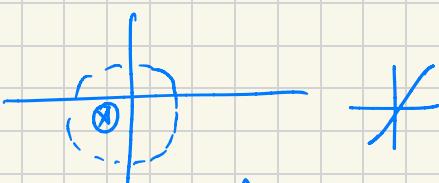
best case:  $O(\log n)$

worst case:  $O(n)$

If it is K-NN, the best case  $O(K \log n)$

worst case  $O(K * n)$ .

Limitations of KD-tree:



for 2D it need 4 regions to check

When d is not small the time complexity  $O(2^d \log n)$