

# Lecture 14

Thursday, October 19, 2023 10:02 AM

Bayes classifier

$$P(c_i | \vec{x}_j) = \frac{\overbrace{P(\vec{x}_j | c_i)}^{\text{likelihood}} \cdot \underbrace{P(c_i)}_{\text{Prior}}}{P(\vec{x}_j)}$$

Full Bayes

→ Parametric: Assume  $P(\vec{x}_j | c_i)$  is normally distributed

$$P(\vec{x}_j | c_i) \sim N(\vec{x}_j | \hat{\mu}_i, \hat{\Sigma}_i)$$

$$N \propto e^{-\frac{(\vec{x}_j - \hat{\mu}_i)^T \hat{\Sigma}_i^{-1} (\vec{x}_j - \hat{\mu}_i)}{2}}$$

$\theta = \{ \hat{\mu}_i, \hat{\Sigma}_i \}_{i=1}^k$

Naive Bayes:

all attributes are independent

$$\Sigma_i = \begin{bmatrix} \sigma_{i1}^2 & 0 & 0 \\ 0 & \sigma_{i2}^2 & 0 \\ 0 & 0 & \sigma_{id}^2 \end{bmatrix}$$

$d = \# \text{ of attributes}$   
 $i = \text{class}$

Non-parametric

$$P(\vec{x}_j | c_i)$$

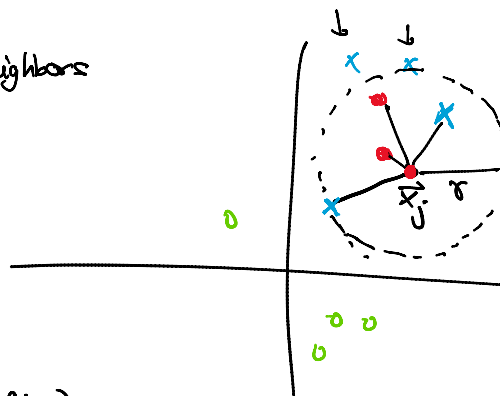
$k = \# \text{ of classes}$

$K = \# \text{ of nearest neighbors}$

K-NN approach

K - nearest neighbors

K=5



# of pts from each	
$K_1 = 2$	$n_1 = 3$
$K_2 = 2$	$n_2 = 4$
$K_3 = 1$	$n_3 = 5$
$K = 5$	

$0, \dots, \rightarrow \boxed{K=5} \rightarrow \dots$

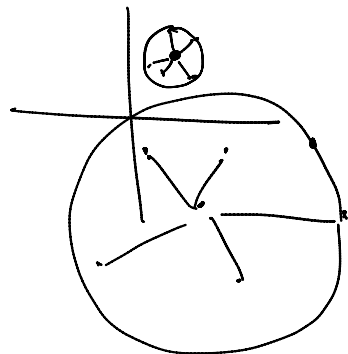
$$P(c_i | \vec{x}_j) \sim \frac{f(\vec{x}_j | c_i) P(c_i)}{P(\vec{x}_j)}$$



$$\frac{1}{K=5}$$

$$B_d(r) = \{x_i \mid \|\vec{x}_i - \vec{x}_j\| \leq r\}$$

$r$ : distance from  $\vec{x}_j$  to its  $K$ -th NN



$$f(\vec{x}_j | c_i) = \frac{\text{fraction of points in } c_i}{\text{vol}(B_d(r))}$$

# of points in  $B_d(r) \equiv K$

$$f(x_j | c_i) = \frac{k_i / n_i}{V}$$

$$f(x_j | c_1) = \frac{2/3}{V}$$

$$f(x_j | c_3) = \frac{1/5}{V}$$

$$f(x_j | c_2) = \frac{2/4}{V}$$

$$P(c_i | \vec{x}_j) = \frac{f(x_j | c_i) P(c_i)}{\sum_{a=1}^K f(x_j | c_a) \cdot P(c_a)}$$

$$= \frac{\frac{k_i / n_i}{V}}{\frac{n_i}{n}} = \frac{\frac{k_i}{n_i}}{\frac{n_i}{n}}$$

$$\frac{k_i}{n_i} \cdot \frac{n}{n_i} = \frac{k_i}{V n}$$

$$P(c_i | x_j) = \frac{k_i}{\sum k_a} = \frac{k_i}{K}$$

Posterior probability = fraction of points in class  $c_i$  within the ball  $B_d(r)$

$$P(c_i | \vec{x}_j) = \frac{2}{5}$$

$K$ -NNI classifier

$$P(C_1 | \vec{x}) = 2/5$$

$$P(C_2 | \vec{x}) = 2/5$$

$$P(C_3 | \vec{x}) = 1/5$$

K-NN classifier

K-NN classifier  $(D, \vec{z}, K)$   
 $\uparrow$  data  $\uparrow$  any point

1) find the K nearest neighbors of  $\vec{z}$

( compute distance  $\|\vec{z} - \vec{x}_i\| \forall \vec{x}_i \in D$   
 sort, pick the smallest K )

2) Count  $K_i \forall i = 1 \dots K$

$\Downarrow$   
 # of points in class  $C_i$

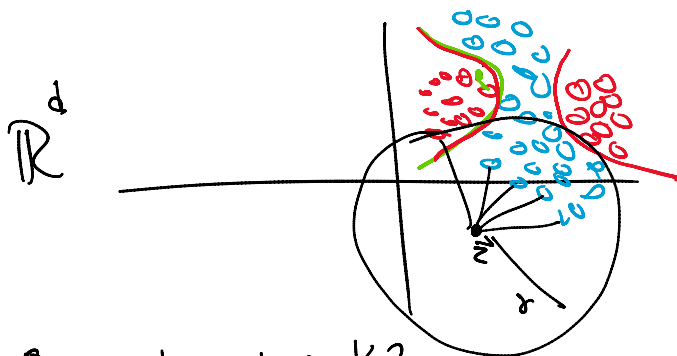
worst case  $\rightarrow O(n \cdot d)$

$O(\log n \cdot d)$

Only true for small d

$$\hat{y} = \arg \max_{i=1}^K \{ K_i \}$$

K-NN is actually very powerful ! Non-linear



K=1

Small K  $\leftarrow$  rough boundary "noisy"

large K  $\leftarrow$  very smooth

Q1: what value of K?

Q2: what norm?

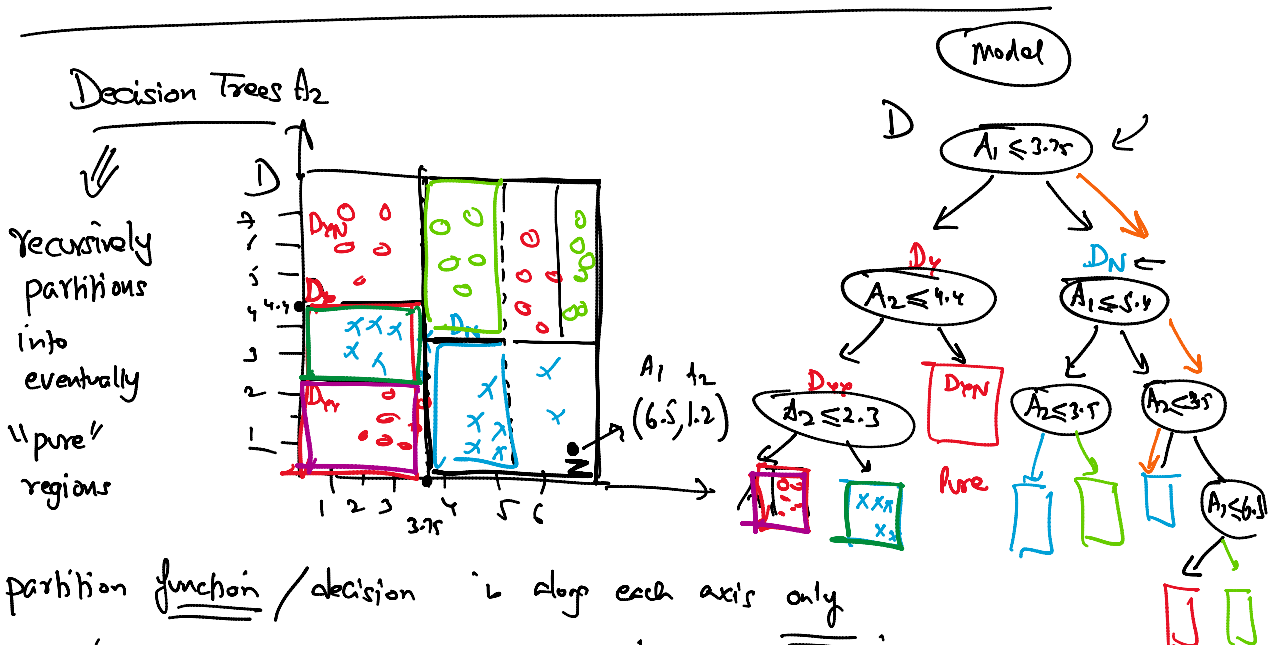
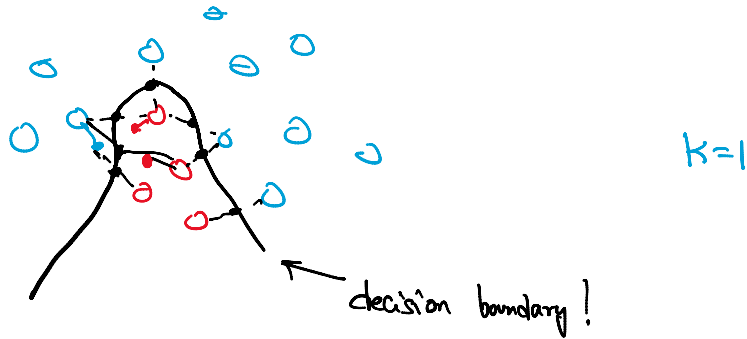
$$L_2 = \|\vec{z} - \vec{x}_i\|^2 \forall \vec{x}_i \in D$$

which distance metric should we use

$L_p$   
6.5.0

Q3: time!

each point requires  $O(n \cdot d)$  time to classify!



partition function / decision is along each axis only.

axis-aligned hyperplane

choose the best hyperplane / decision at each step

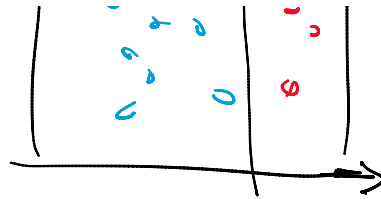
greedy

each axis acts like the 'normal' vector to the hyperplane



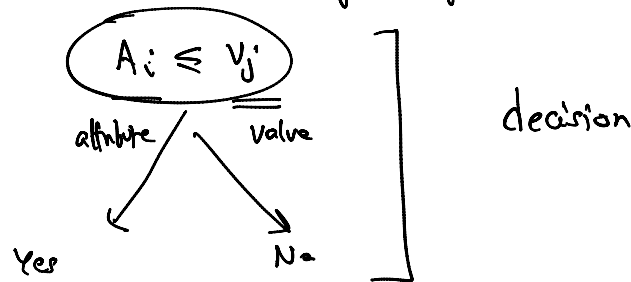


'normal' vector to the hyperplane



Axis aligned hyperplane

↳ all decisions are of the form



at each step: we have a 'dataset' of interest

Decision Tree (D, stopping criteria  $\equiv$  purity or size  $\theta$  or  $m$ )

base { 1) If  $|D| \leq m$  stop (size)  
2) If  $\text{purity}(D) \geq \theta$  stop (majority class)

3) evaluate all splits or decisions of the kind

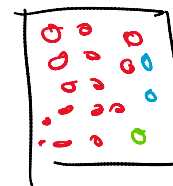
$$A_i \leq v_j$$

$\forall i = 1 \dots d$  (all attributes)

for all mid-points  $v_j$

select the best (greedy)

(try all (attribute, value) combinations)



4) use  $A_i^* \leq v_j^*$  best split to partition D

$$D_r = \{ \vec{x}_a \mid x_{ai} \leq v_j^* \}$$

$$D_N = \{ \vec{x}_a \mid x_{ai} > v_j^* \}$$

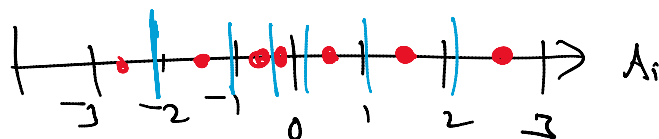
5) Decision Tree (DT)

Decision Tree (DN)

how to choose the best  $(A_i, v_j)$  combination

$A_i \leq v_j$  decision.

D



1) find the distinct values along  $A_i$

-2.7, -1.5, -0.8, -0.2, 0.5, 1.5, 2.6

$v_j \{ -2.1, -1.15, -0.5, 0.5, 1, 2.15 \}$

try all of them

$$\begin{aligned} A_i &\leq -2.1 \\ A_i &\leq -1.15 \\ &\vdots \\ A_i &\leq 2.15 \end{aligned}$$

which of these is the best over all  $A_i \leq v_j$  D

Criteria: Information Gain  $\swarrow \searrow$   
 $D_N \quad D_r$

Criteria: Information Gain  $\begin{matrix} \swarrow & \searrow \\ D_Y & D_N \end{matrix}$

$$IG(D, D_Y, D_N) = \underbrace{H(D)}_{\text{original entropy}} - \underbrace{H(D_Y, D_N)}_{\text{joint entropy after the split}}$$

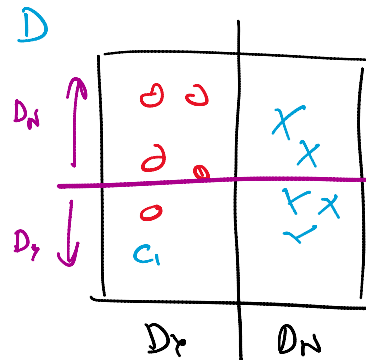
we want this value to be smaller

$H(D) \equiv$  entropy

$$= - \sum_{i=1}^k p_i \log p_i$$

$$p_i = \text{prob of class } c_i \text{ in } D \\ = \frac{n_i}{n} = \frac{\text{\# of points in } c_i}{\text{\# of points in } D}$$

$$H(D_Y, D_N) = \frac{n_Y}{n} H(D_Y) + \frac{n_N}{n} H(D_N)$$



$c_1 = \text{red}$   
 $c_2 = \text{blue}$

$$p_1 = 5/10 = 0.5$$

$$p_2 = 5/10 = 0.5$$

$$- \left( \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right)$$

$$- \log(0.5)$$

worst possible entropy

$H(D)$

$$c_1 = 5/5 = 1 \\ c_2 = 0/5 = 0$$

$$c_1 = 0/5 = 0 \\ c_2 = 5/5 = 1$$

$$H(D_N) = 0$$

$$H(D_Y) = -(1 \cdot \log 1 + 0 \cdot \log 0) \\ = 0$$

$$H(D_Y, D_N) = \frac{5}{10} \cdot 0 + \frac{5}{10} \cdot 0 \\ = 0$$

$$IG = - \log(0.5)$$

higher

$$H(D_Y, D_N) = \frac{5}{10} \cdot 0 + \frac{5}{10} \cdot 0$$

$$= 0$$

higher

$$IG = \underline{H(D)} - \underline{H(D_Y, D_N)}$$

Gini Index:

worst case

$$\underline{G(D)} = 1 - \sum_{i=1}^k p_i^2$$

If all  $p_i = 1/k$

$$1 - \sum_{i=1}^k (1/k^2)$$

$$1 - k/k^2 = \left( \frac{k-1}{k} \right)$$

worst gini index

$$\underline{G(D) - G(D_Y, D_N)}$$

Gini index is 0  
if we have  
"pure" partitions