

## Agenda: ① Naive Bayes Algorithm (Classification Algorithm)

## ② KNN

## 1. Naive Bayes Algorithm

Rolling a dice  $\{1, 2, 3, 4, 5, 6\}$  → Independent events

$P(1) = \frac{1}{6} \quad P(3) = \frac{1}{6} \quad P(5) = \frac{1}{6}$

Picking marbles from box → Dependent events

$$\begin{array}{l} \text{1. } P(B) = \frac{3}{6} = \frac{1}{2} \quad \begin{array}{|c|} \hline \text{Blue} \\ \hline \end{array} \\ \text{2. } P(G/B) = \frac{3}{5} \quad \begin{array}{|c|} \hline \text{Blue} \\ \hline \text{Green} \\ \hline \end{array} \end{array}$$

Conditional Probability: Probability of Green if Blue is already taken once before

$$\begin{aligned} P(\text{B and G}) &= P(\text{B}) \times P(\text{G}/\text{B}) \\ &= \frac{1}{2} \times \frac{3}{5} = \frac{3}{10} = 0.3 \end{aligned}$$

$P(A \text{ and } B) = P(A) \times P(B/A)$

$P(A \text{ and } B) = P(B \text{ and } A)$

$P(A) \times P(B/A) = P(B) \times P(A/B)$

$$P(B/A) = \frac{P(B) \times P(A/B)}{P(A)} \quad \text{Bayes' Theorem}$$

Dataset:  $\underbrace{x_1, x_2, \dots, x_n}_{\text{input features}}$   $\downarrow$   $y$   $\uparrow$   $\text{output}$ 

$$\begin{aligned} P(y/(x_1, x_2, \dots, x_n)) &= \frac{P(y) \times P((x_1, x_2, \dots, x_n)/y)}{P((x_1, x_2, \dots, x_n))} \\ &= \frac{P(y) \cdot P(x_1/y) \cdot P(x_2/y) \cdot \dots \cdot P(x_n/y)}{P(x_1) \cdot P(x_2) \cdot \dots \cdot P(x_n)} \end{aligned}$$

## Binary Dataset

 $x_1 \quad x_2 \quad x_3 \quad \text{o/p}$ 

Yes

No

 $\rightarrow P((x_1, x_2, x_3)/\text{Yes})$ 

$$P(y=\text{Yes}/x_i) = \frac{P(\text{Yes}) \cdot P(x_1/\text{Yes}) \cdot P(x_2/\text{Yes}) \cdot P(x_3/\text{Yes})}{P(x_1) \cdot P(x_2) \cdot P(x_3)}$$

$$P(y=\text{No}/x_i) = \frac{P(\text{No}) \cdot P(x_1/\text{No}) \cdot P(x_2/\text{No}) \cdot P(x_3/\text{No})}{P(x_1) \cdot P(x_2) \cdot P(x_3)}$$

Let's say  $P(\text{Yes}/x_i) = 0.13$

$P(\text{No}/x_i) = 0.05$

$P(\text{Yes}/x_i) = \frac{0.13}{0.13+0.05} = \frac{0.13}{0.18} \approx \% 72 \geq \% 50$

$P(\text{No}/x_i) = \frac{0.05}{0.13+0.05} = \frac{0.05}{0.18} \approx \% 28 + \% 100 <% 50$

## PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Outlook	Temperature		P(Y)	P(N)	P(Y)	P(N)	
	Yes	No					
Sunny	2	3	2/5	3/5	Hot	2	2
Overcast	4	0	4/9	0/9	Mild	4	2
Rain	3	2	3/9	2/9	Cold	3	1
	9	5				9	5
Play		P(Yes)		P(No)			
Yes	9	P(Yes)		P(No)			
No	5	9/14		5/14			

(Sunny, Hot)  $\rightarrow$  O/P Naive Bayes  $\Rightarrow ?$

$$P(\text{Yes} | (\text{sunny}, \text{hot})) = \frac{P(\text{Yes}) \cdot P(\text{sunny} | \text{Yes}) \cdot P(\text{hot} | \text{Yes})}{P(\text{sunny}) \cdot P(\text{hot})} = \frac{9/14}{9/14} \cdot \frac{2/9}{2/9} = \frac{2}{63} = 0.031$$

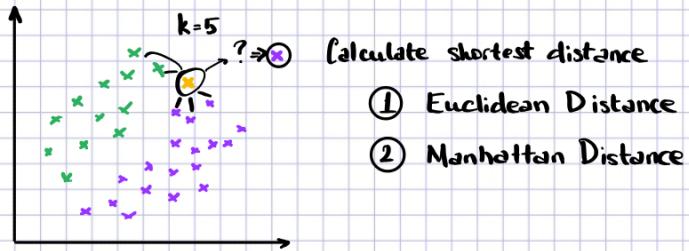
$$P(\text{No} | (\text{sunny}, \text{hot})) = \frac{P(\text{No}) \cdot P(\text{sunny} | \text{No}) \cdot P(\text{hot} | \text{No})}{P(\text{sunny}) \cdot P(\text{hot})} = \frac{5/14}{5/14} \cdot \frac{3/5}{3/5} \cdot \frac{2/5}{2/5} = \frac{3}{35} = 0.0857$$

Real Probability

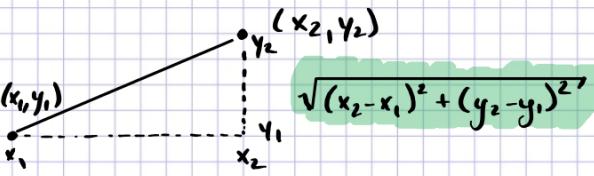
$$P(\text{Yes} | (\text{sunny}, \text{hot})) = 0.031 = \frac{0.031}{0.031 + 0.0857} = \% 27$$

$$P(\text{No} | (\text{sunny}, \text{hot})) = 0.0857 = 1 - \% 27 = \% 73$$

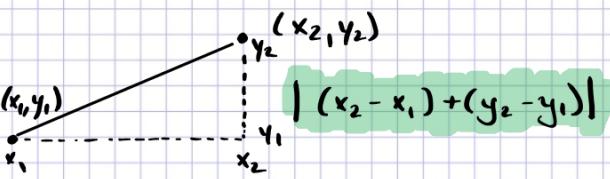
## 2. K-Nearest Neighbors (KNN) [Classification & Regression]



### 1. Euclidean Distance



### 2. Manhattan Distance



The same idea is used in regression problem. Output is calculated as average of k-nearest data points' y-values.