

3rd Day → Machine Learning Algorithms

Agenda

- ① Practicals
- ② Naive Bayes's Intuition
- ③ KNN algorithms

→ Simple Examples

Previous Session

- ① Linear Regression
- ② Ridge & Lasso
- ③ Logistic Regression

⇒ Complex

① Naive Bayes's Intuition { Classification }



{ BAYE'S THEOREM }

Rolling a Dice

{ 1, 2, 3, 4, 5, 6 }

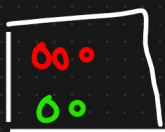
$$P(1) = \frac{1}{6}$$

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

$$P(2) = \frac{1}{6}$$

{ Independent Events }

Dependent Event



First Event ✓
 $P(R) = \frac{3}{5} \rightarrow R$ ✓

Green Marble

$$P(G) = \frac{2}{4} = \frac{1}{2} \rightarrow G$$

Dependents

$$P(G) = \frac{2}{5}$$

$$P(R) = \frac{3}{4}$$

$$P(R \text{ and } G) = P(R) * P(G/R)$$

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

→ Conditional Probability

$$P(A) * P(B/A) = P(B) * P(A/B)$$

$$P(A) * P(B/A) = P(B) * P(A/B)$$

Bayes Theorem

Naive Bayes

T/P

$\boxed{y} \rightarrow \alpha_p$

$$P(y/x_1, x_2, x_3 \dots x_n) = \frac{P(y) * P(x_1, x_2 \dots x_n / y)}{P(x_1, x_2 \dots x_n)}$$

Yes or No

$$= \frac{P(y) * P(x_1/y_1) * P(x_2/y_2) * P(x_3/y_3) \dots P(x_n/y_n)}{P(x_1) * P(x_2) * P(x_3) \dots P(x_n)}$$

DARME

$$\begin{matrix} x_1 & x_2 & x_3 & x_4 & y \\ \parallel & \parallel & \parallel & \parallel & \parallel \end{matrix}$$

→

 Yes ✓

 No ✓

$$P(y=y_0/x_i) = P(y_0) * P(x_1/y_0) * P(x_2/y_0) * P(x_3/y_0) * P(x_4/y_0)$$

Constant Ignore

$$\rightarrow P(x_1) \neq P(x_2) \neq P(x_3) \neq P(x_4) \quad \# \text{ fixed}$$

$$P(y = N_0 | x_i) = \frac{P(N_0) * P(x_1/N_0) * P(x_2/N_0) * P(x_3/N_0) * P(x_4/N_0)}{P(N_0) * P(x_1/N_0) * P(x_2/N_0) * P(x_3/N_0) * P(x_4/N_0) + P(N_1) * P(x_1/N_1) * P(x_2/N_1) * P(x_3/N_1) * P(x_4/N_1)}$$

Constant

→ $P(x_1) \neq P(x_2) \neq P(x_3) \neq P(x_4) \neq \text{fixed}$

x_i $\begin{cases} \rightarrow \text{Yes} \\ \rightarrow \text{No} \end{cases}$

$$P(\text{Yes} | x_i) = \boxed{0.13}$$

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

\downarrow

$$\geq 0.5 \Rightarrow 1$$

$$< 0.5 \Rightarrow 0$$

$$P(\text{Yes} | x_i) = \frac{0.13}{0.13 + 0.05} = 0.72 = 72\%$$

$$P(\text{No} | x_i) = 1 - 0.72 = 0.28 = 28\%$$

DATASET

Binary class

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
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

x_1
Outlook

$P(\text{Sunny} | \text{Yes})$

	Yes	No	$P(Y)$	$P(N)$
Sunny	2	3	2/9	3/5
Overcast	4	0	4/9	0/5
Rain	3	2	3/9	2/5
<u>Total</u>	<u>9</u>	<u>5</u>		

PLAY

Temperature

	Yes	No	$P(Y)$	$P(N)$
Hot	2	2	2/9	2/5
Mild	4	2	4/9	2/5
Cold	3	1	3/9	1/5
<u>Total</u>	<u>9</u>	<u>5</u>		

	Yes	No	$P(Y)$	$P(N)$
Yes	9		$\frac{9}{14}$	$\frac{5}{14}$
No		5		
<u>Total</u>	<u>14</u>			

→ Test (Sunny, Hot) → O/P

$$P(\text{Yes} | \text{Sunny, Hot}) = \frac{P(\text{Yes}) * P(\text{Sunny} | \text{Yes}) * P(\text{Hot} | \text{Yes})}{\dots}$$

$$\begin{aligned}
 & \frac{P(\text{Sunny}) * P(\text{Hot})}{\cancel{1/14} * \cancel{\frac{2}{7}} * \frac{2}{9}} \\
 & = \frac{2}{63} = 0.031
 \end{aligned}$$

$$P(\text{No} \mid \text{Sunny, Hot}) = \frac{P(\text{No}) * P(\text{Sunny} \mid \text{No}) * P(\text{Hot} \mid \text{No})}{\cancel{P(\text{Sunny}) * P(\text{Hot})} \rightarrow \text{Constant}}$$

$$\begin{aligned}
 & = \frac{8}{14} * \frac{3}{5} * \frac{2}{5} \\
 & = \frac{3}{35} = 0.085
 \end{aligned}$$

$$P(\text{Yes} \mid \text{Sunny, Hot}) = 0.031 = 1 - 0.73 = 0.27 = 27\%$$

$$P(\text{No} \mid \text{Sunny, Hot}) = 0.085 = \frac{0.085}{0.031 + 0.085} = 0.73 = 73\%$$

→ (Sunny, hot) → Yes or No

Answer → No ✓

Assignment

(Overcast, Mild) → Naive Bayes?

② KNN Algorithm {K Nearest Neighbour}

Classification

Regression

① Classification



K Nearest Neighbour

$K=5$

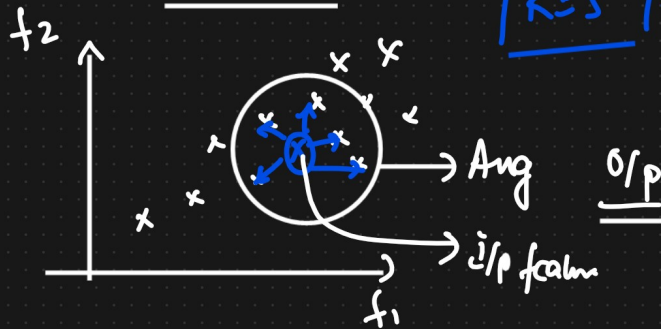
Maximum No

Red = 3

White = 2

$K=1$ to 50

Regression



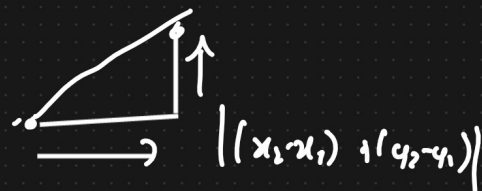
K Nearest

Euclidean Distance

(x_1, y_1) (x_2, y_2)

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Manhattan Distance



$K=5$ \rightarrow Hyperparameters

Error Rate $\uparrow \uparrow$

① Outliers

② Imbalanced

