

01-11-2025

## Agenda:

→ Naive Bayes Algo (2 methods)

→ k-means clustering

## Naive Bayes Algorithm

$$\xrightarrow{\text{Bayes theorem}} \frac{P(y|x)}{P(x)} = \frac{P(x|y) \cdot P(y)}{P(x)}$$

↓  
probability

what if ?

- Data is categorical !
- we want probabilistic reasoning
- we want something that learns from frequency count.
- we want fast & interpretable solution,
- NB → classification → probabilistic model

Dataset: (play, Tennis)

Outlook	Temperature	Play	Prob of Rainy & Hot?
Sunny	Hot	No	$P(Y_0) = \frac{4}{6} \rightarrow 0.66$
Sunny	Mild	No	
Cloudy	Hot	Yes	$P(No) = \frac{2}{6} \rightarrow 0.33$
Rainy	Mild	Yes	
Rainy	Cool	Yes	
Cloudy	Cool	No	

$$P(Y_0) \rightarrow 0.66$$

outlook	O - calc	Temp	T - calc
Cloudy	$2/4 \rightarrow 0.5$	Hot	$1/4 \rightarrow 0.25$
Rainy	$2/4 \rightarrow 0.5$	Mild	$1/4 \rightarrow 0.25$
Sunny	$0/4 \rightarrow 0$	Cool	$2/4 \rightarrow 0.5$
<hr/>			
Cloudy	$0/2 \rightarrow 0$	Hot	$1/2 \rightarrow 0.5$
Rainy	$0/2 \rightarrow 0$	Mild	$1/2 \rightarrow 0.5$
Sunny	$2/2 \rightarrow 1$	Cool	$0/2 \rightarrow 0$

→ prob of Rainy & Hot

$$\begin{aligned}
 P(Y_0 | \text{Rain, Hot}) &\propto P(Y_0) \times P(\text{Rain} | Y_0) \times P(\text{Hot} | Y_0) \\
 &= 0.66 \times 0.5 \times 0.25 \\
 &= 0.0825
 \end{aligned}$$

$$\begin{aligned}
 p(\text{No} \mid \text{Rain, Hot}) &\propto p(\text{No}) \times p(\text{Rain} \mid \text{No}) \times p(\text{Hot} \mid \text{No}) \\
 &= 0.33 \times 0 \times 0.5 \\
 &= 0
 \end{aligned}$$

$$p(\text{Yes} \mid \text{Rain, Hot}) > p(\text{No} \mid \text{Rain, Hot}) \Rightarrow \text{Yes}$$

→ prob of sunny & cool

$$\begin{aligned}
 p(\text{Yes} \mid \text{sunny, cool}) &\propto p(\text{Yes}) \times p(\text{sunny} \mid \text{Yes}) \times p(\text{cool} \mid \text{Yes}) \\
 &= 0.66 \times 0 \times 0.5 \\
 &= 0
 \end{aligned}$$

$$\begin{aligned}
 p(\text{No} \mid \text{sunny, cool}) &\propto p(\text{No}) \times p(\text{sunny} \mid \text{No}) \times p(\text{cool} \mid \text{No}) \\
 &= 0.33 \times 1 \times 0 \\
 &= 0
 \end{aligned}$$

$$\begin{array}{ccc}
 p(\text{Yes} \mid \text{sunny, cool}) & & p(\text{No} \mid \text{sunny, cool}) \\
 \searrow & & \swarrow \\
 \text{Both are } 0
 \end{array}$$

↓  
To solve this we need smoothing

we add 1 to every count & divide by  
# samples + # categories

outlook	$\theta$	-calc	Temp	$T_{\text{calc}}$
cloudy	$2/4 \rightarrow$	$2+1/4+3 \rightarrow 0.429$	Hot	$1/4 \rightarrow 1+1/4+3 \rightarrow 0.286$
Rainy	$2/4 \rightarrow$	$2+1/4+3 \rightarrow 0.429$	Mild	$1/4 \rightarrow 1+1/4+3 \rightarrow 0.286$
Sunny	$0/4 \rightarrow$	$0+1/4+3 \rightarrow 0.143$	Cool	$4/4 \rightarrow 2+1/4+3 \rightarrow 0.429$
<hr/>				
cloudy	$0/2 \rightarrow$	$0+1/2+3 \rightarrow 0.2$	Hot	$1/2 \rightarrow 1+1/2+3 \rightarrow 0.4$
Rainy	$0/2 \rightarrow$	$0+1/2+3 \rightarrow 0.2$	Mild	$1/2 \rightarrow 1+1/2+3 \rightarrow 0.4$
Sunny	$2/2 \rightarrow$	$2+1/2+3 \rightarrow 0.6$	Cool	$0/2 \rightarrow 0+1/2+3 \rightarrow 0.2$

$$p(Y_e) \rightarrow 4/6 \rightarrow 0.66$$

$$p(\text{No}) \rightarrow 2/6 \rightarrow 0.33$$

→ prob of sunny & cool

$$p(Y_e | \text{sunny}, \text{cool}) \propto p(Y_e) \times p(\text{sunny} | Y_e) \times p(\text{cool} | Y_e)$$

$$= 0.66 \times 0.143 \times 0.429 = \underline{\underline{0.040}}$$

$$p(\text{No} | \text{sunny}, \text{cool}) \propto p(\text{No}) \times p(\text{sunny} | \text{No}) \times p(\text{cool} | \text{No})$$

$$= 0.33 \times 0.6 \times 0.2$$

$$= \underline{\underline{0.0396}}$$

Normalise to get the final probabilities:

$$\text{Sum} = 0.04 + 0.0396 = 0.0796$$

$$P(Y_0 | \text{sunny}, \text{no}) = \frac{0.040}{0.0796} \approx 0.5025 \quad \downarrow \\ \text{sum} \approx 1$$

$$P(Y_1 | \text{sunny}, \text{no}) = \frac{0.0396}{0.0796} \approx 0.4971 \rightarrow$$

$$P(Y_0 | \text{sunny}, \text{no}) > P(Y_1 | \text{sunny}, \text{no})$$

$$p(Y_0) = p(\text{no})$$

$\rightarrow$  features  $\rightarrow$  categorical  
 $\rightarrow$  target  $\rightarrow$  categorical  $\rightarrow$  categorical Naive Bayes

Next topic:

Dataset  $\rightarrow$  numeric features  $\rightarrow$  categorical NB?  
 $\rightarrow$  categorical target column

$\rightarrow$  Gaussian NB  
 $\downarrow$   
(normal)

we compute:

$$P(y|x) \propto p(y) \prod N(x | \mu_y, \sigma_y^2) \quad \begin{matrix} \exp(z) \\ 1 \\ e^z \\ \downarrow \\ 2\pi\sigma^2 \end{matrix}$$

$$N(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Data set:

id	$x_1$	$x_2$	$y$
A	1.0	2.0	Ye
B	2.0	1.0	Ye
C	1.5	1.8	Ye
D	3.0	3.0	No
E	3.5	2.8	No
F	2.9	3.2	No

Step 1: Calculate MLE:

Maximum likelihood estimate:  $(\hat{M}_{LE}) \rightarrow$  uses  $n$  instead of  $n-1$  which is used in unbiased sample variance

mean, variance, std

$$\underline{\underline{\mu(x_1 = Ye)}} = \frac{1+2+1.5}{3} \rightarrow \mu_{Ye,1} = 1.5$$

$$\underline{\underline{\mu(y_2 = Ye)}} = \frac{2+1+1.8}{3} \rightarrow \mu_{Ye,2} = 1.6$$

$$\underline{\underline{\mu(x_1 = No)}} = \frac{3+3.5+2.9}{3} \rightarrow \mu_{No,1} = 3.133$$

$$\underline{\underline{\mu(x_2 = No)}} = \frac{3+2.8+3.2}{3} \rightarrow \mu_{No,2} = 3.0$$

Variable  $x_1$ :

$$(\hat{M}_{Le}) : (1-1.5)^2, (2-1.5)^2, (1.5-1.5)^2 \\ : (-0.5)^2, (0.5)^2, 0^2$$

$$\text{sum} : 0.25 + 0.25 + 0 \rightarrow 0.5/3 \rightarrow 0.166$$

$$\sigma^2 = \sum (x - \mu)^2 / N$$

Variance  $\chi_2$ :

$$(\text{Yes}) : (0.4)^2, (-0.6)^2, (0.2)^2 \rightarrow 0.56 \\ : 0.56/3 \rightarrow 0.186$$

Variance  $\chi_1$ : sum  $\rightarrow 0.2066 \rightarrow 0.2066/3 \rightarrow 0.0688$   
(No)

Variance  $\chi_2$ : sum  $\rightarrow 0.08 \rightarrow 0.08/3 \rightarrow 0.0266$   
(No)

$$\underline{x} = \begin{pmatrix} \underline{\chi_1} & \underline{\chi_2} \\ 2.0 & 2.0 \end{pmatrix}$$

Step 2: calculate Gaussian prob per feature:

for class = Yes

$\rightarrow \chi_1$ :

$$\rightarrow \sigma = \sqrt{0.186} \approx 0.4308$$

$$\rightarrow \text{coeff} = 1/\left(\sqrt{2\pi}\sigma^2\right) \approx \frac{0.9772}{\sigma^2}$$

$$\rightarrow \text{exponent exp} : \left( -\frac{(2-1.5)^2}{2 \times 0.186} \right) \approx \exp(-0.75) \approx 0.472$$

$$\rightarrow p(\chi_1 | \text{Yes}) \approx \frac{1}{\sqrt{2\pi}\sigma^2} \cdot \exp\left(-\frac{(\chi_1 - \mu)^2}{2\sigma^2}\right)$$

$$\approx 0.9772 \times 0.472$$

$$\approx 0.4616$$

$$p(x_1 | \text{Yes}) \quad p(x_2 | \text{Yes})$$

$$p(x_1 | \text{No}) \quad p(x_2 | \text{No})$$

$x_2$ :

$$\rightarrow \sigma = \sqrt{0.1866} \approx 0.4320$$

$$\rightarrow \text{coeff} \approx 0.923$$

$$\rightarrow \text{exp}(\text{non exp}) : \exp\left(-\frac{(2-1.6)^2}{2 \times 0.1866}\right) = \exp(-0.428) \approx 0.6510$$

$$p(x_2 | \text{Yes}) \approx 0.4320 \times 0.6510$$

$$\approx 0.6009$$

For (all) = No:

$$p(x_1 | \text{No}) \approx 1.36 \times 10^{-1}$$

$$p(x_2 | \text{No}) \approx 1.68 \times 10^{-8}$$

$$p(x_1 | \text{Yes}) \approx 0.4616$$

$$p(x_2 | \text{Yes}) \approx 0.6009$$

Joint prob:  $p(x | \text{Yes}) = 0.4616 \times 0.6009$   
 $\approx 0.2774$

$$p(x_1 | \text{No}) \approx 1.36 \times 10^{-1} \quad p(x_2 | \text{No}) \approx 1.68 \times 10^{-8}$$

Joint prob:  $p(x | \text{No}) \approx 2.29 \times 10^{-12}$

$$p(x | \gamma_{e1}) \approx 0.2779 \quad p(x | \text{No}) \approx 2.29 \times 10^{-12}$$

↓

Unnormalized prob

$$\begin{aligned} p(\gamma_{e1} | x) &: p(\gamma_{e1}) \cdot p(x | \gamma_{e1}) & p(\text{No} | x) &: p(\text{No}) \cdot p(x | \text{No}) \\ &: 0.5 \times 0.2779 & &: 0.5 \times 2.29 \times 10^{-12} \\ &\approx : 0.1387 & &\approx : 1.145 \times 10^{-12} \end{aligned}$$

Step 3: normalize to get class probabilities.

$$\text{sum } S = 0.1387 + 1.145 \times 10^{-12} \approx 0.1387$$

so normalized!

$$p(\gamma_{e1} | x) \approx \frac{0.1387}{0.1387} \approx 0.9999 \dots . \quad 1$$

$$p(\text{No} | x) \approx \frac{1.145 \times 10^{-12}}{0.1387} \approx 8.26 \times 10^{-12}$$

Decision: predict  $\rightarrow \gamma_{e1}$  (very strong)

↓

$$\log p(y|x) = \log p(y) + \prod \log N(x | \mu, \sigma^2)$$

↑

python implement

when to use Gaussian NB:

- features are continuous
- you want a very fast, interpretable model,
- we are assuming that features are normally distributed.

$$\text{Email} \rightarrow \begin{cases} \text{spam} & p(s) = 0.2 \\ \text{Not spam} & p(n) = 0.8 \end{cases}$$

s : please click on the link below

$$p(\text{click}(s)) = x_1 =$$
$$p(\text{click}(n)) = x_2 =$$

$$n : p(\ ) = p_1 \\ = p_2 \\ p_3$$

→ Tomorrow is holiday  
↓  
empty

$$p(\text{Yes} | x) = p(\text{spam}) \times p(\text{Tom} | s) \times p(\text{click} | s)$$
$$= 0.6$$
$$p(\text{No} | x) = 0.4$$

New Algo:

k-mean clustering !

→ Unsupervised learning → No target column

Data → k-mean → group data based on  
clustering similarity (clustering)



how many  
groups you  
want ?.



K

Initiation:

Stop train

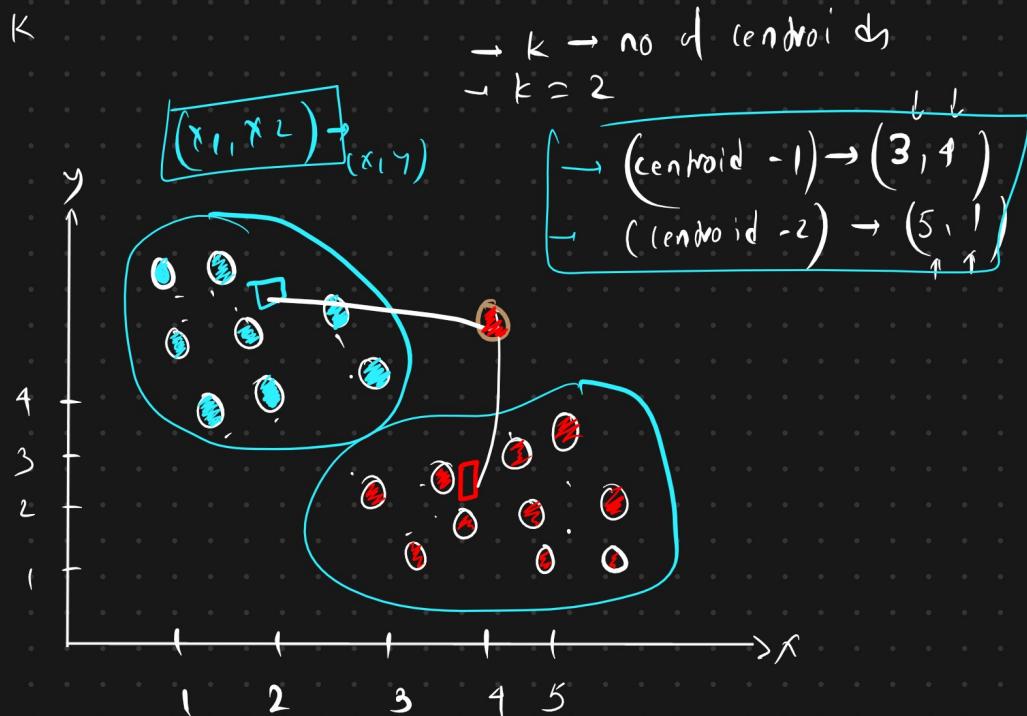
B → 8

R → 10

Iter - 3 :

B → 8

R → 10



Dataset:

$$A = \underline{(1,2)}$$

$$B = (1.5, 1.8)$$

$$C = (5, 8)$$

$$D = (8, 8)$$

$$E = (1, 0.6)$$

$$F = (9, 11)$$

$k=2$

$$c_1 = A = (1, 2)$$

$$c_2 = C = (5, 8)$$

euclidean distance:

$$(x_1, x_2), (A_1, A_2)$$

$$\sqrt{(x_1 - A_1)^2 + (x_2 - A_2)^2}$$

Step 2: update centroids

New  $c_1$ : mean of  $A(1,2)$ ,  $B(1.5, 1.8)$ ,  $E(1, 0.6)$

$$c_{1-x} : (1 + 1.5 + 1) / 3 = 1.1667$$

$$c_{1-y} : (2 + 1.8 + 0.6) / 3 = 1.4667$$

New  $c_2$ :

mean of  $C(5, 8)$ ,  $D(8, 8)$ ,  $F(9, 11)$

$$c_{2-x} : (5 + 8 + 9) / 3 \approx 7.33$$

$$c_{2-y} : (8 + 8 + 11) / 3 \approx 9$$

$$c_1 : (1.1667, 1.4667)$$

$$c_2 : (7.33, 9)$$

Iteration 1 : (Step 1)

1. point  $A(1,2)$

$$\rightarrow c_1 : \sqrt{(1-1)^2 + (2-2)^2} = 0$$

$$\rightarrow c_2 : \sqrt{(1-5)^2 + (2-8)^2} = 7.211$$

→ assign to cluster  $c_1$

2. point  $B(1.5, 1.8)$

$$\rightarrow c_1 : \sqrt{(1.5-1)^2 + (1.8-2)^2} \approx 0.5385$$

$$\rightarrow c_2 : \sqrt{(1.5-5)^2 + (1.8-8)^2} \approx 7.119$$

→ cluster 1

3. point  $C$  → cluster 2

4. point  $D$  → cluster 2

5. point  $E$  → cluster 1

6. point  $F$  → cluster 2

cluster 1 :  $\{A, B, E\}$

cluster 2 :  $\{C, D, F\}$

How to compute loss  $\rightarrow$  WCSS

$\rightarrow$  we calculate WCSS after every iteration of centroid update.

$\rightarrow$  WCSS  $\rightarrow$  within cluster summation square

$\rightarrow$  sum over clusters & point of squared distance to their centroid.

Iteration: 2

$$\rightarrow A(1,2) \approx 0.5587$$

$\rightarrow C_1$

$$\rightarrow B(1.5, 1.8) \approx 0.4714$$

$\rightarrow C_1$

$$\rightarrow C(5,8) \approx 2.539$$

$\rightarrow C_2$

$$\rightarrow D(8,8) \approx 1.2019$$

$\rightarrow C_2$

$$\rightarrow E(9,11) \approx 2.602$$

$\rightarrow C_2$

$$\text{new-}C_1 (1.1667, 1.4667)$$

$$\text{new-}C_2 (7.3333, 9)$$

WCSS :

$C_1 \rightarrow A, B, E$

$$(0.5587)^2 + (0.4714)^2 + (0.1837)^2$$

$$= 1.3163$$

$C_2 \rightarrow C, D, F$

$$(2.539)^2 + (1.2019)^2 + (2.602)^2$$

$$= 14.666$$

$$\text{WCSS} \rightarrow \text{WCSS}(C_1) + \text{WCSS}(C_2)$$

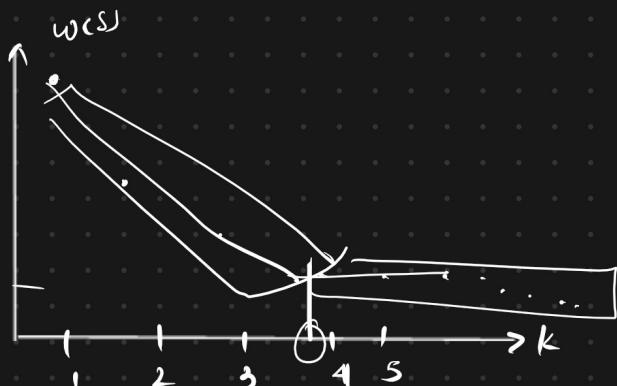
$$\rightarrow \underline{\underline{15.982}}$$

$$\underline{\underline{18.683}}$$

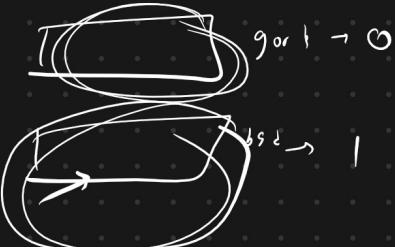
$\rightarrow$  To get optimal value  $k$ , plot elbow graph

x axis  $\rightarrow k$

y-axis  $\rightarrow$  WCSS



increase in  $K$ , decrease  $\omega_{CS}$

dataset  $\rightarrow$  1000 rows  $\rightarrow$  kmean (2)  $\rightarrow$    
 $\downarrow$   
review  
 $L$

  $k=2$

