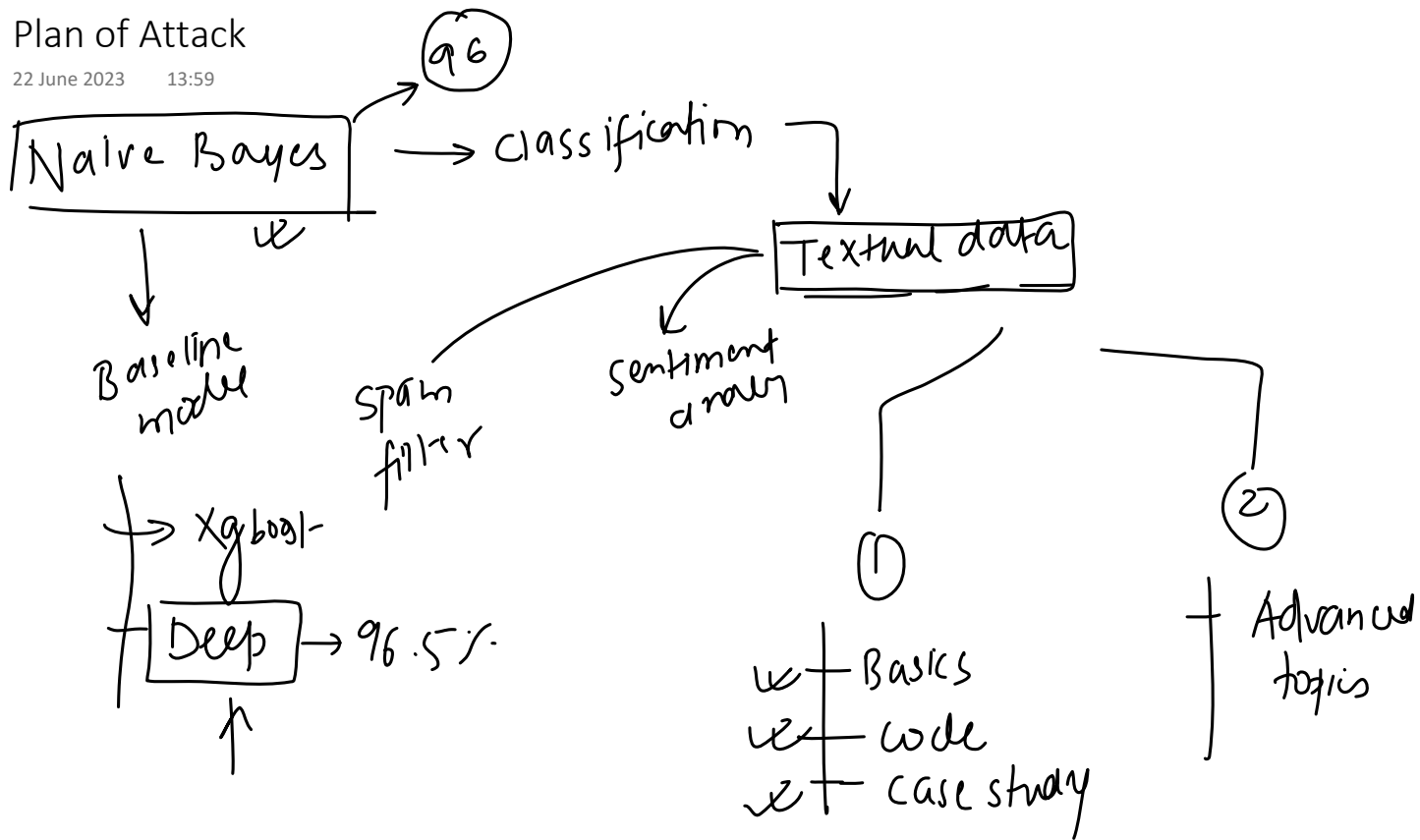


# Plan of Attack

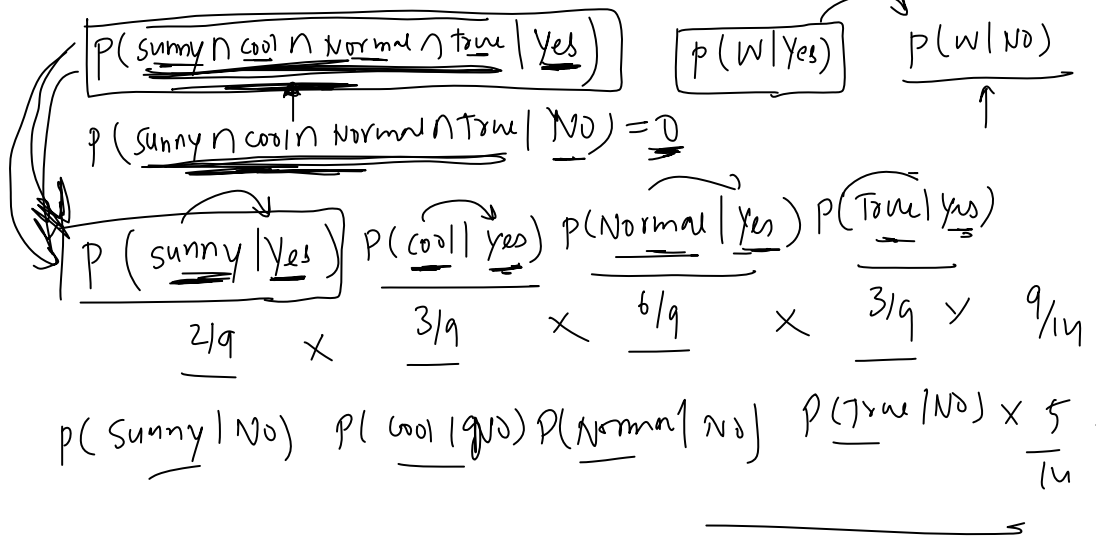
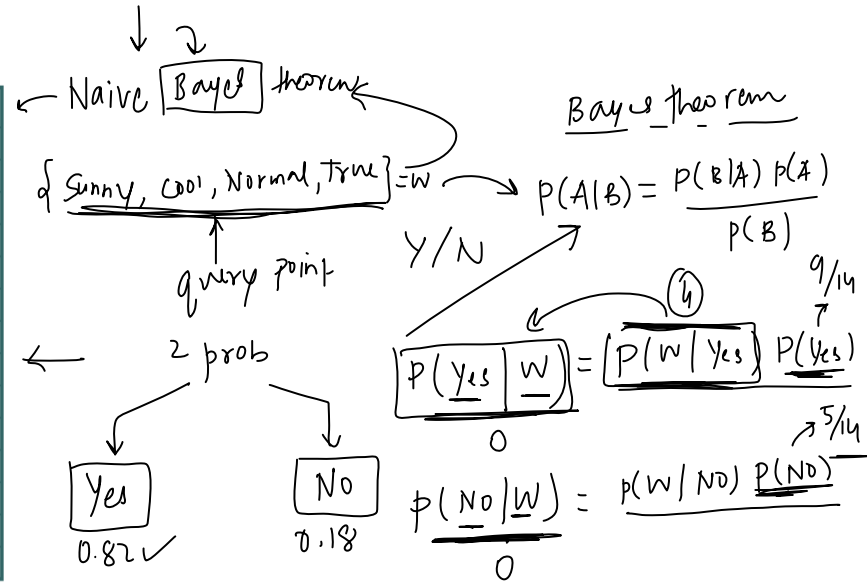
22 June 2023 13:59



# play tennis Toy binary classification

14 days

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No ✗
Sunny	Hot	High	True	No ✗
Overcast	Hot	High	False	Yes ✓
Rainy	Mild	High	False	Yes ✓
Rainy	Cool	Normal	False	Yes ✓
Rainy	Cool	Normal	True	No ✗
Overcast	Cool	Normal	True	Yes ✓
Sunny	Mild	High	False	No ✗
Sunny	Cool	Normal	False	Yes ✓
Rainy	Mild	Normal	False	Yes ✓
Sunny	Mild	Normal	True	Yes ✓
Overcast	Mild	High	True	Yes ✓
Overcast	Hot	Normal	False	Yes ✓
Rainy	Mild	High	True	No ✗



y N

$x_1 \ x_2 \ \dots \ x_n$   $y$  multiclass classification  $\underbrace{1, 2, 3, \dots, k}_{k \text{ classes}}$

$\rightarrow \langle x_1, x_2, x_3, \dots, x_n \rangle \rightarrow$  naive bayes  $k$  prob

$$\begin{cases} p(y_1 | x_T) = p(x_T | y_1) p(y_1) \\ p(y_2 | x_T) \rightarrow p(x_T | y_2) p(y_2) \\ \vdots \\ p(y_k | x_T) \rightarrow p(x_T | y_k) p(y_k) \end{cases}$$

$$(x_T) = \langle x_1, x_2, \dots, x_n \rangle$$

$$P(y_k | x_T) = P(x_T | y_k) P(y_k)$$

$$= P(x_1 \wedge x_2 \wedge x_3 \wedge \dots \wedge x_n | y_k) P(y_k)$$

$$= P(\underbrace{x_1, x_2, x_3, \dots, x_n}_A | \underbrace{y_k}_B) P(y_k)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$= \frac{P(x_1, x_2, x_3, \dots, x_n, y_k)}{P(y_k)}$$

$$P(A \cap B) = \frac{P(A|B) P(B)}{1}$$

$$= P(\underbrace{x_1}_A, \underbrace{x_2, x_3, \dots, x_n, y_k}_B)$$

$$= \underbrace{P(x_1 | x_2, x_3, \dots, x_n, y_k)}_A \underbrace{P(x_2, x_3, \dots, x_n, y_k)}_B$$

$$\downarrow \quad \downarrow$$

$$P(x_1 | x_2, x_3, \dots, x_n, y_k) \quad P(x_2 | x_3, x_4, \dots, x_n, y_k) \quad P(x_3 | x_4, \dots, x_n, y_k) \quad \dots$$

$$= \frac{P(x_1 | x_2, x_3, \dots, x_n, y_k)}{1} \frac{P(x_2 | x_3, x_4, \dots, x_n, y_k)}{1} \dots \frac{P(x_{n-1} | x_n, y_k)}{1} P(x_n | y_k) P(y_k)$$

Naive assumption  $\rightarrow$  features are independent of each other

$x_1$	$x_2$	$x_3$	$\dots$	$x_n$
-------	-------	-------	---------	-------

$$P(A|B) = P(A)$$

$$P(A|B \cap C) = P(A|B, C) = P(A|C)$$

$$\begin{array}{c}
 \begin{array}{|c|c|c|c|c|} \hline x_1 & x_2 & x_3 & \dots & x_n \\ \hline \end{array} \\
 \uparrow \quad \quad \quad \uparrow \\
 y_1 = \underbrace{p(x_1 | y_k)}_{\substack{\uparrow \\ y_k}} \underbrace{p(x_2 | y_k)}_{\uparrow} \underbrace{p(x_3 | y_k)}_{\uparrow} \dots \underbrace{p(x_{n-1} | y_k)}_{\uparrow} \underbrace{p(x_n | y_k)}_{\uparrow} \underbrace{p(y_k)}_{\uparrow}
 \end{array}$$

$y_2$   
 $\vdots$   
 $y_h$   
 $y_k$

$\max \rightarrow$  maximum a posteriori rule  
(MAP)

$p(A|B \cap C) = p(A|B, C) = \frac{p(A|C)}{p(C)}$

Code

22 June 2023 14:00

2

No, Yes

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	<u>Hot</u>	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	<u>Mild</u>	High	False	Yes
Rainy	<u>Cool</u>	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Sunny, Hot, High, False

Training

Probabilities

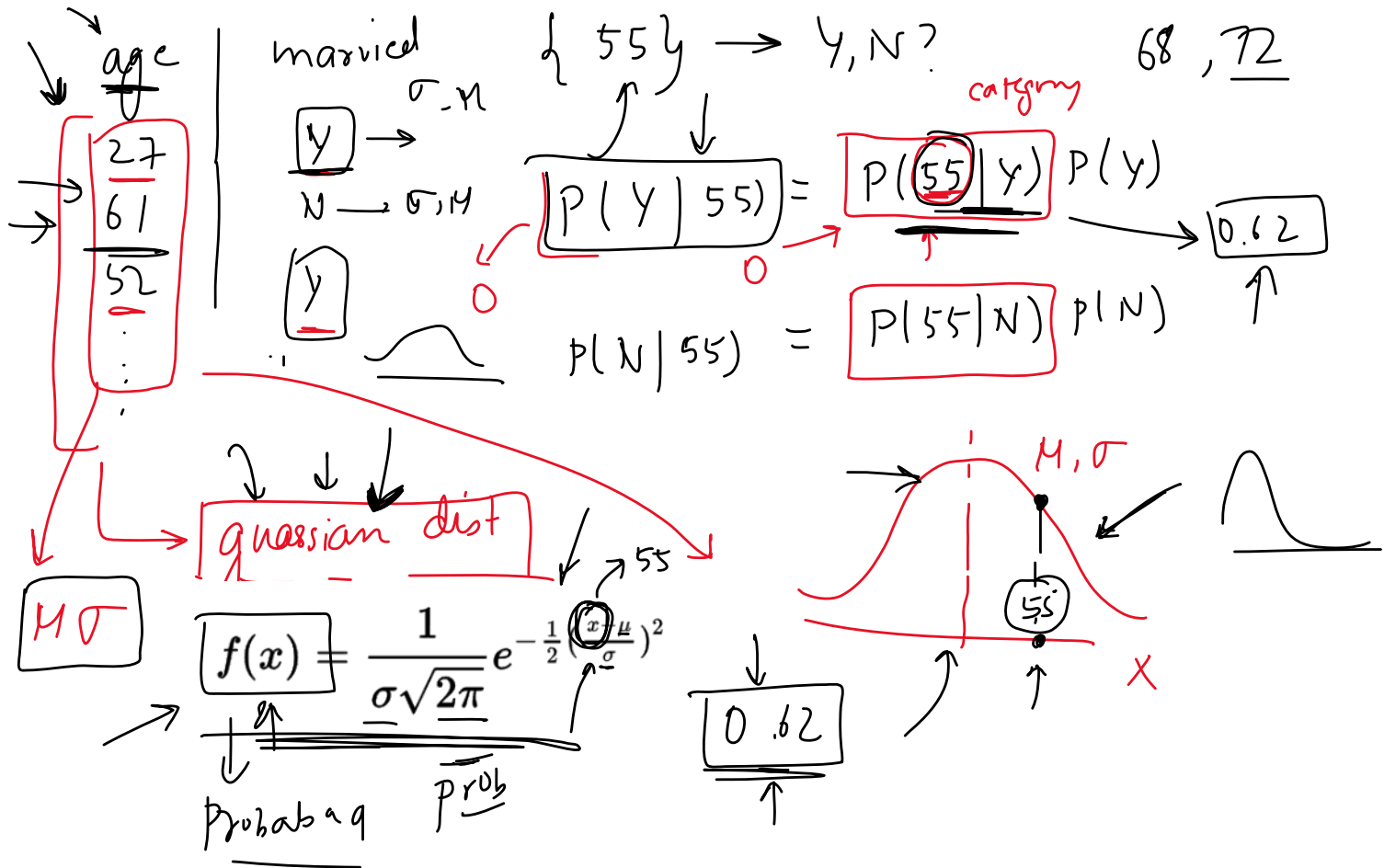
Testing

Dictionary

Outlook  $\rightarrow$  Sunny | Overcast | Rainy $3 \times 2 = 6$  Prob $3 \times 2$  $2 \times 2$  $P(\text{Sunny} | \text{Yes})$  $P(\text{Sunny} | \text{No})$  $P(\text{Overcast} | \text{Yes})$  $P(\text{Overcast} | \text{No})$  $P(\text{Rainy} | \text{Yes})$  $P(\text{Rainy} | \text{No})$  $P(\text{Hot} | \text{Y})$  $P(\text{Hot} | \text{N})$  $P(\text{Mild} | \text{Y})$  $P(\text{Mild} | \text{N})$  $P(\text{Cool} | \text{Y})$  $P(\text{Cool} | \text{N})$  $P(\text{H} | \text{Y})$  $P(\text{H} | \text{N})$  $P(\text{M} | \text{Y})$  $P(\text{M} | \text{N})$  $P(\text{C} | \text{Y})$  $P(\text{C} | \text{N})$

# How Naïve Bayes handles numerical data?

22 June 2023 16:51



# What is data is not Gaussian?

22 June 2023 16:52

1. **Data Transformation:** Depending on the nature of your data, you could apply a transformation to make it more normally distributed. Common transformations include the logarithm, square root, and reciprocal transformations.

y

10-20 | 20-30 | -

2. **Alternative Distributions:** If you know or suspect that your data follow a specific non-normal distribution (e.g., exponential, Poisson, etc.), you can modify the Naïve Bayes algorithm to assume that specific distribution when calculating the likelihoods.

→ Try out

3. **Discretization:** You can turn your continuous data into categorical data by binning the values. There are various ways to decide on the bins, including equal width bins, equal frequency bins, or using a more sophisticated method like k-means clustering. Once your data is binned, you can use the standard Multinomial or Bernoulli Naïve Bayes methods.

binning

4. **Kernel Density Estimation:** A non-parametric way to estimate the probability density function of a random variable. Kernel density estimation can be used when the distribution is unknown.

5. **Use other models:** If none of the above options work well, it may be best to consider a different classification algorithm that doesn't make strong assumptions about the distributions of the features, such as Decision Trees, Random Forests, or Support Vector Machines.

→

# Naïve Bayes on Text Data

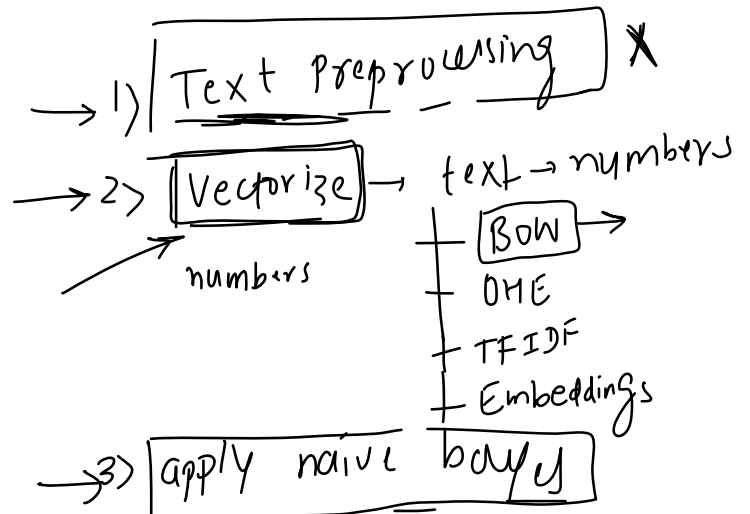
22 June 2023 14:13

## Sentiment Analysis

review | sentiment  
 [ ] | +  
 [ ] | -  
 [ ] | +

5000

transform



BOW

10L words → 35000 uniquely

reviews | sentiment  
 [I liked the movie] | +ve  
 [really hated the movie] | -ve

→ [I liked the movie]  
 [really hated the movie]  
 [ ]  
 [ ]  
 [ ]

(50000, 5000)  
 → shaped

hate	like	movie	actor	...	1
0	1	1	0	...	-
2	0	3	2	...	-
3	2	1	1	...	-
...	...	...	...	...	-

Adj - - great movie

numbers

(0 2 1 ... 5) → (1, 5000)  
 vector

+ve

-ve



$$\begin{array}{l}
 \boxed{p(+ve \mid \text{hate}=0, \text{like}=2, \dots)} \rightarrow 0.37 \\
 p(-ve \mid \text{hate}=0, \text{like}=2, \dots) \rightarrow \boxed{0.23}
 \end{array}
 \rightarrow \text{+ve positive sentiment}$$

$$\boxed{p(\text{hate}=0 \mid +ve) \quad p(\text{like}=2 \mid +ve) \dots}$$

$\hookrightarrow p$