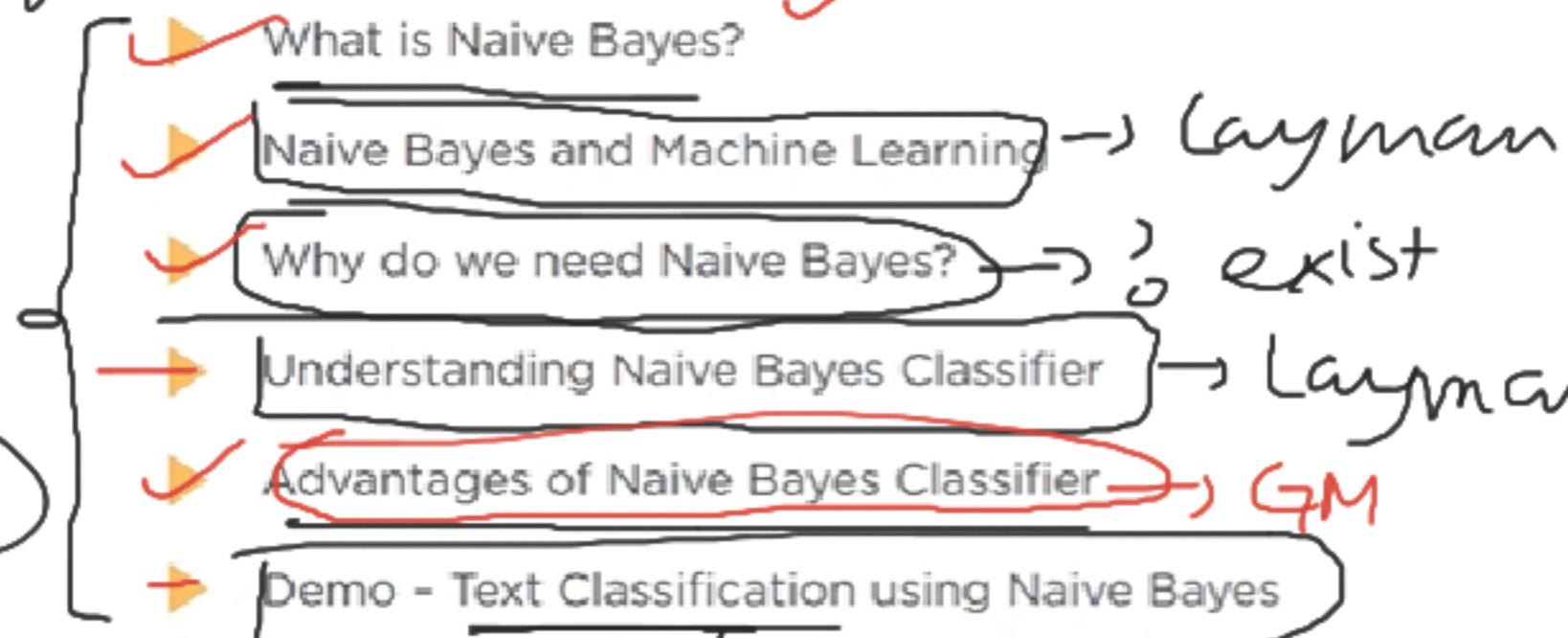


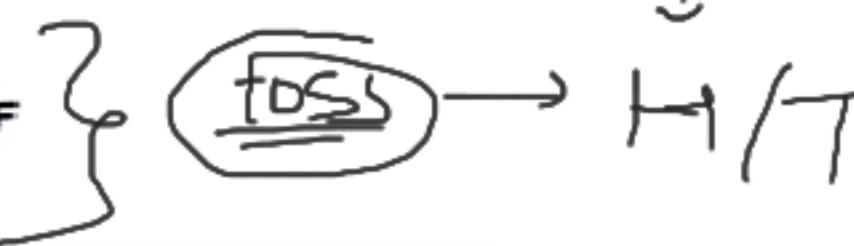
Agenda

brief

Intro NLP



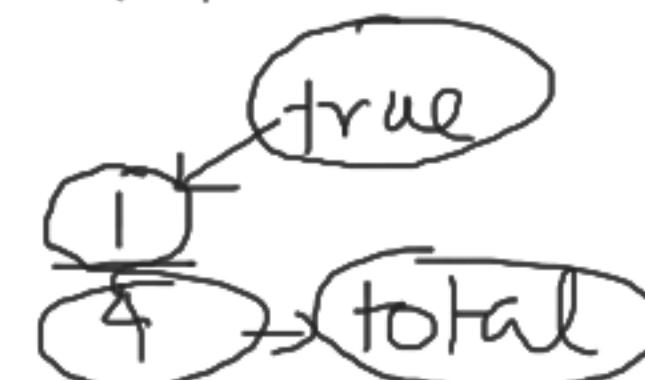
LET US CONSIDER THE FOLLOWING EXAMPLE OF TOSSING TWO COINS



Here, the sample space is:

{HH, HT, TH, TT}

1. $P(\text{Getting two heads}) = 1/4$
2. $P(\text{At least one tail}) = 3/4$
3. $P(\text{Second coin being head given first coin is tail}) = 1/2$
4. $P(\text{Getting two heads given first coin is a head}) = 1/2$



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

Bayes Theorem



→ req

conditional
prob
revision

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

A B

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

where:

$P(A|B)$ = Conditional Probability of A given B

$P(B|A)$ = Conditional Probability of B given A

$P(A)$ = Probability of event A

$P(B)$ = Probability of event A

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

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

|

$$\rightarrow \Rightarrow \boxed{P(A|B) P(B) = P(B|A) P(A)}$$

ACT/ION

P

IN THIS SAMPLE SPACE, LET A BE THE
EVENT THAT SECOND COIN IS HEAD
AND B BE THE EVENT THAT FIRST COIN
IS TAIL.



In the sample space:

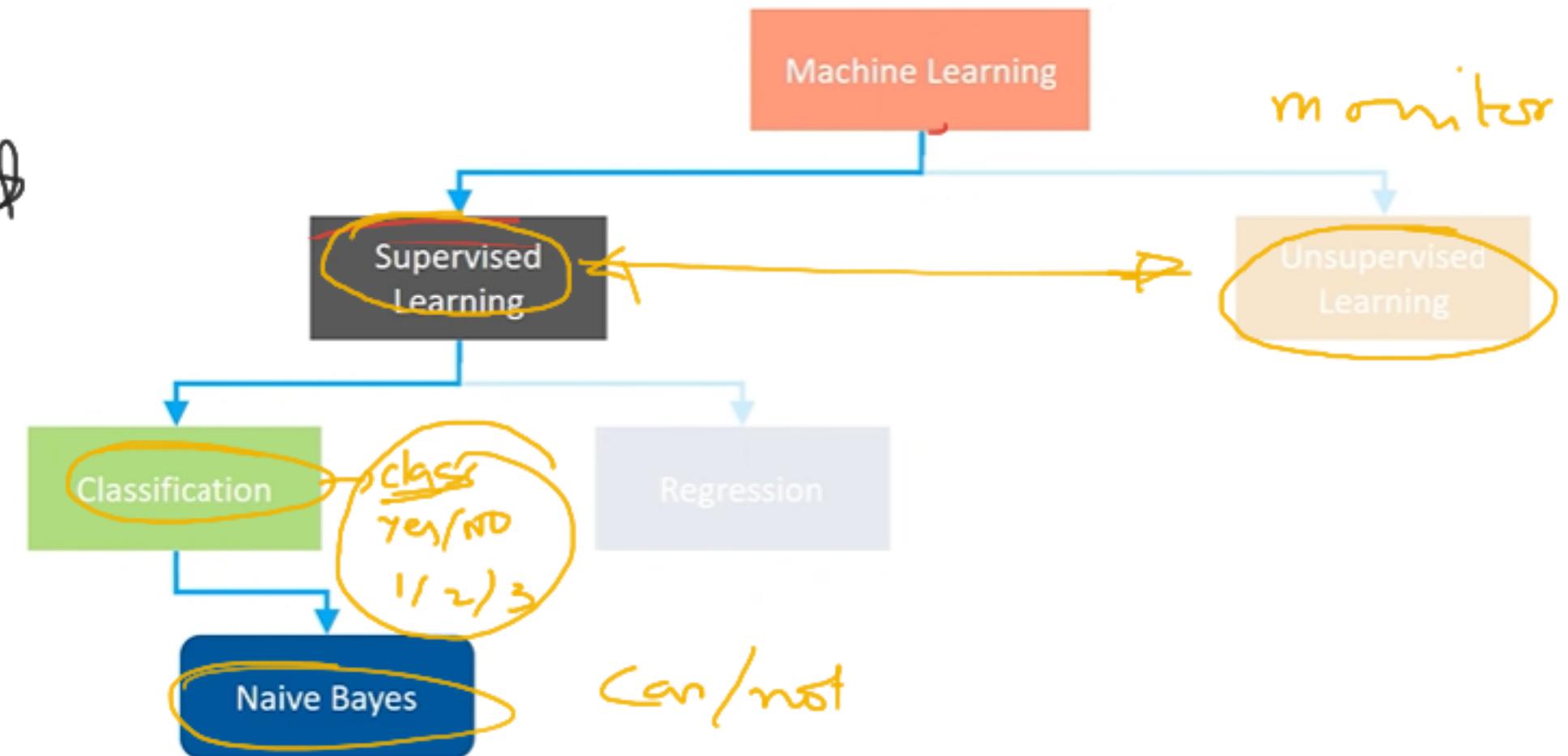
2
4

$$\{\text{HH}, \text{HT}, \text{TH}, \text{TT}\}$$

P(Second coin being head given first coin is tail)

$$\begin{aligned}
 &= P(A|B) \\
 &= [P(B|A) * P(A)] / P(B) \\
 &= [P(\text{First coin being tail} \\
 &\quad \text{being head})] / P(\text{First coin being tail}) \\
 &= [(1/2) * (1/2)] / (1/2) \\
 &= \boxed{1/2} = 0.5
 \end{aligned}$$

$$P(1^{S_1} \text{ tail} / 2^m + 1)$$



$$\frac{1}{2} \times \frac{2}{y}$$

Maine
Buyes

→ Brain

Where is Naive Bayes used?

NLP

P(Ram/Hi T)

Shape

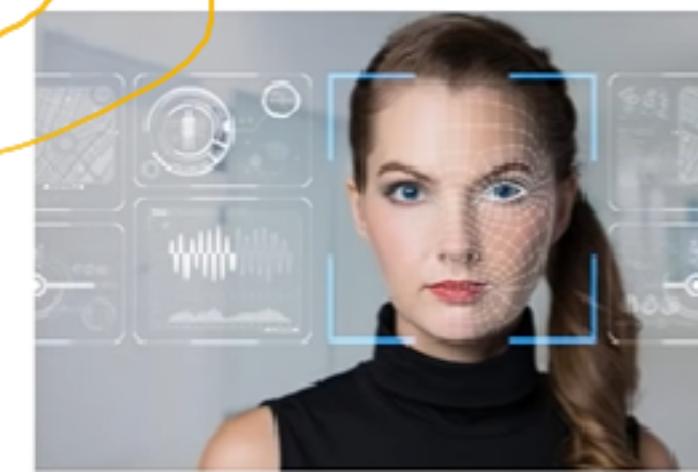
base
yes/no

shape

Video KYC

Neural

Face
Recognition



Weather
Prediction



R
cg

dis
Linear
regular

big
aly
x

pjc

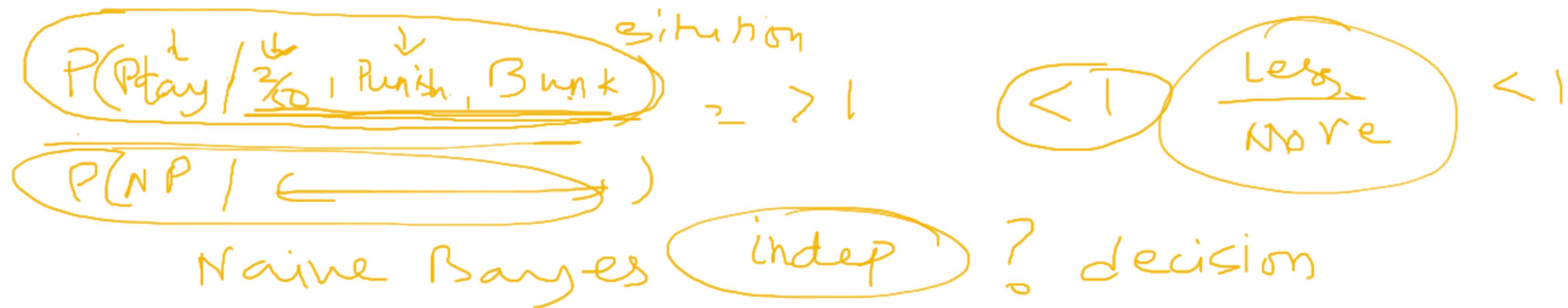
(RGB)

Medical
Diagnosis



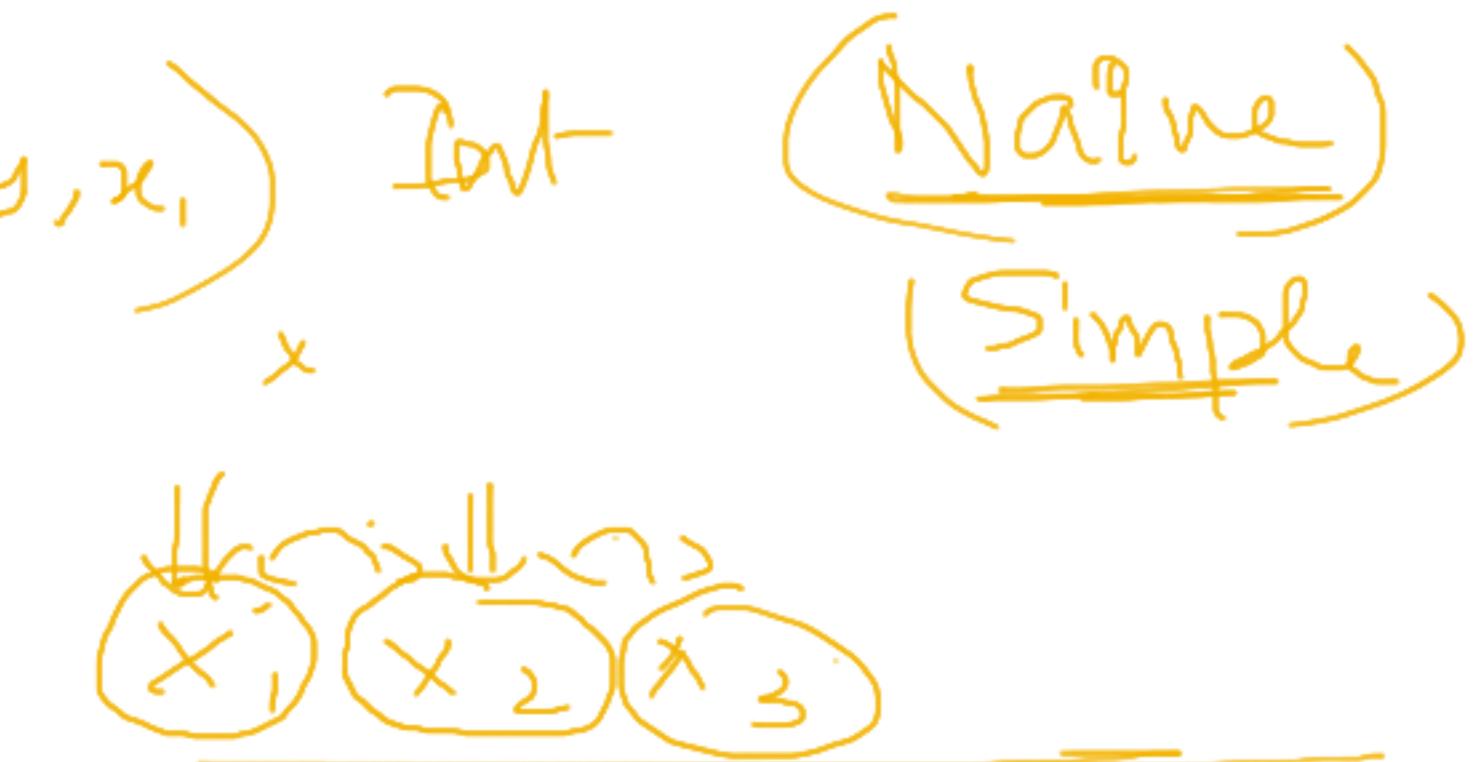
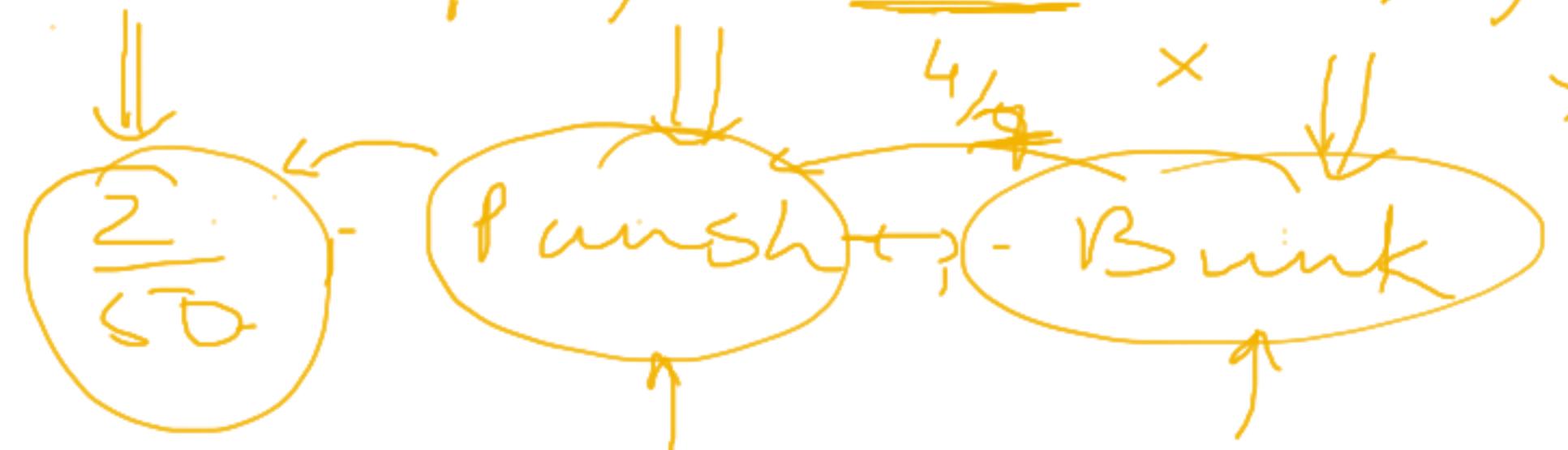
News
Classification



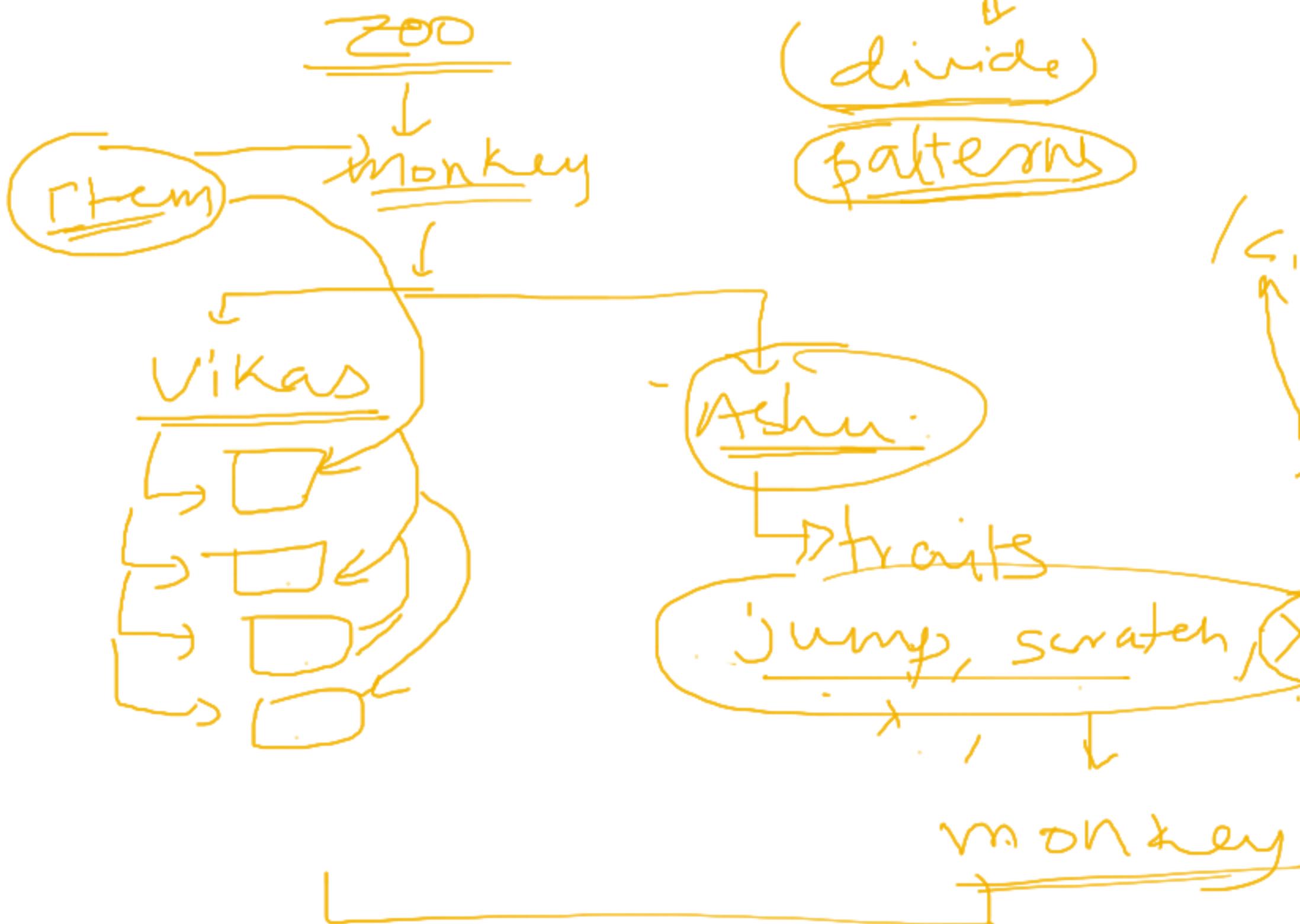


$$P(Y | \underline{X}) = P(y)P(x_1|y)P(x_2|y, x_1) \dots P(x_n|y, x_1, \dots, x_{n-1})$$

$$P(Y/x) = \frac{P(y)}{\underline{x_1 x_2 x_3}} P(x_1, \dots, x_n) \quad \text{Int} \quad \begin{array}{l} \text{(Naive)} \\ \text{(Simple)} \end{array}$$



Discriminative vs Generative



(divide)

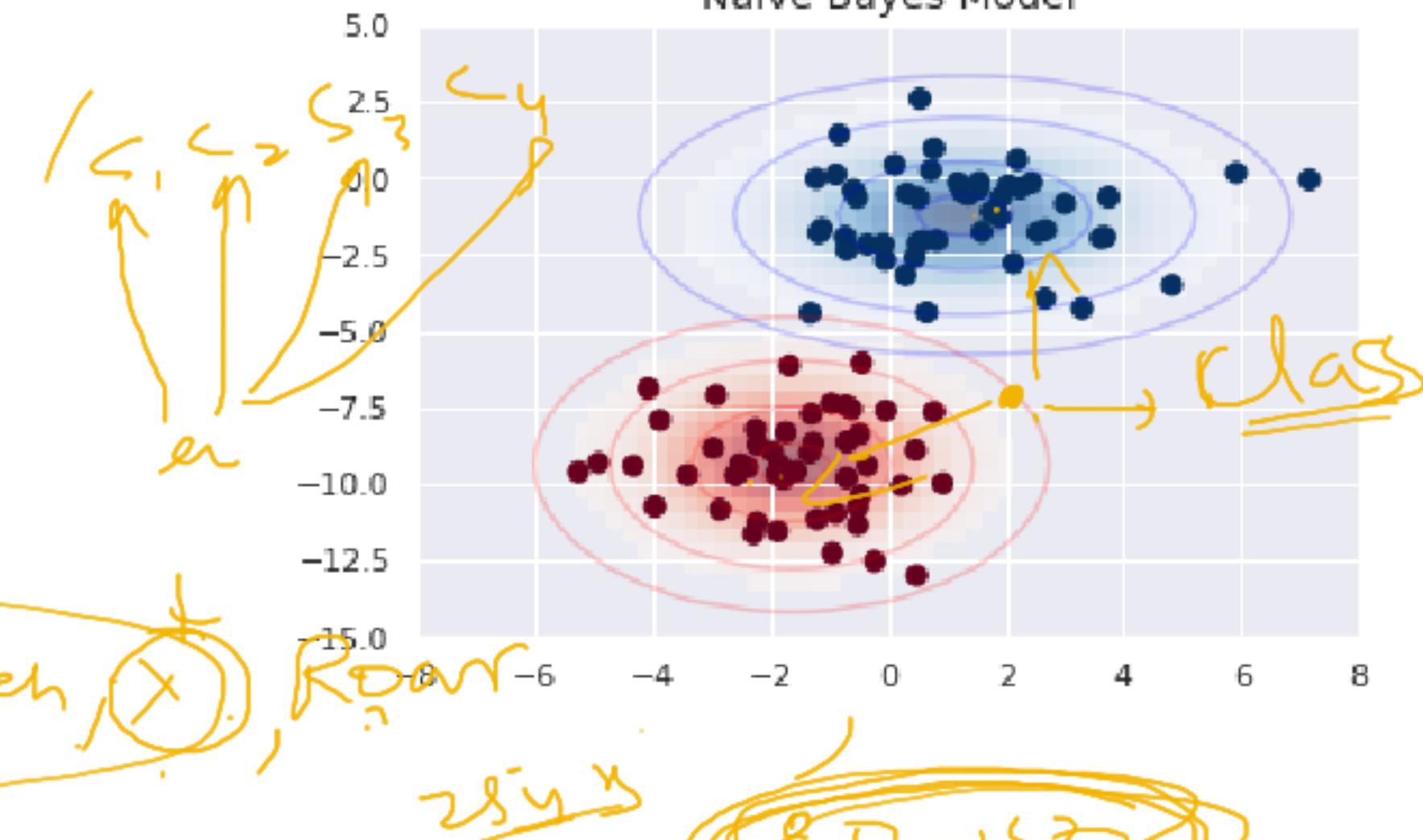
patterns

Generate

scenarios

p rob?

Naive Bayes Model



$$20 \rightarrow 52$$

under



To predict whether a person will purchase a product on a specific combination of Day, Discount and Free Delivery using Naive Bayes Classifier



Generate?

We have a small sample dataset of 30 rows for our demo

Dataset

	A	B	C	D
1	Day	Discount	Free Delivery	Purchase
2	Weekday	Yes	Yes	Yes
3	Weekday	Yes	Yes	Yes
4	Weekday	No	No	No
5	Holiday	Yes	Yes	Yes
6	Weekend	Yes	Yes	Yes
7	Holiday	No	No	No
8	Weekend	Yes	No	Yes
9	Weekday	Yes	Yes	Yes
10	Weekend	Yes	Yes	Yes
11	Holiday	Yes	Yes	Yes
12	Holiday	No	Yes	Yes
13	Holiday	No	No	No
14	Weekend	Yes	Yes	Yes
15	Holiday	Yes	Yes	Yes

Based on this dataset containing three input types of *Day*, *Discount* and *Free Delivery*, we will populate frequency tables for each attribute

Frequency Table		Buy	
		Yes	No
Discount	Yes	19	1
	No	5	5

Frequency Table		Buy	
		Yes	No
Free Delivery	Yes	21	2
	No	3	4

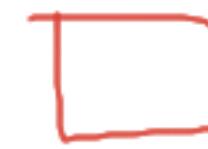
Pick ↓

Frequency Table		Buy	
		Yes	No
Day	Weekday	9	2
	Weekend	7	1
	Holiday	8	3

Pick ↗

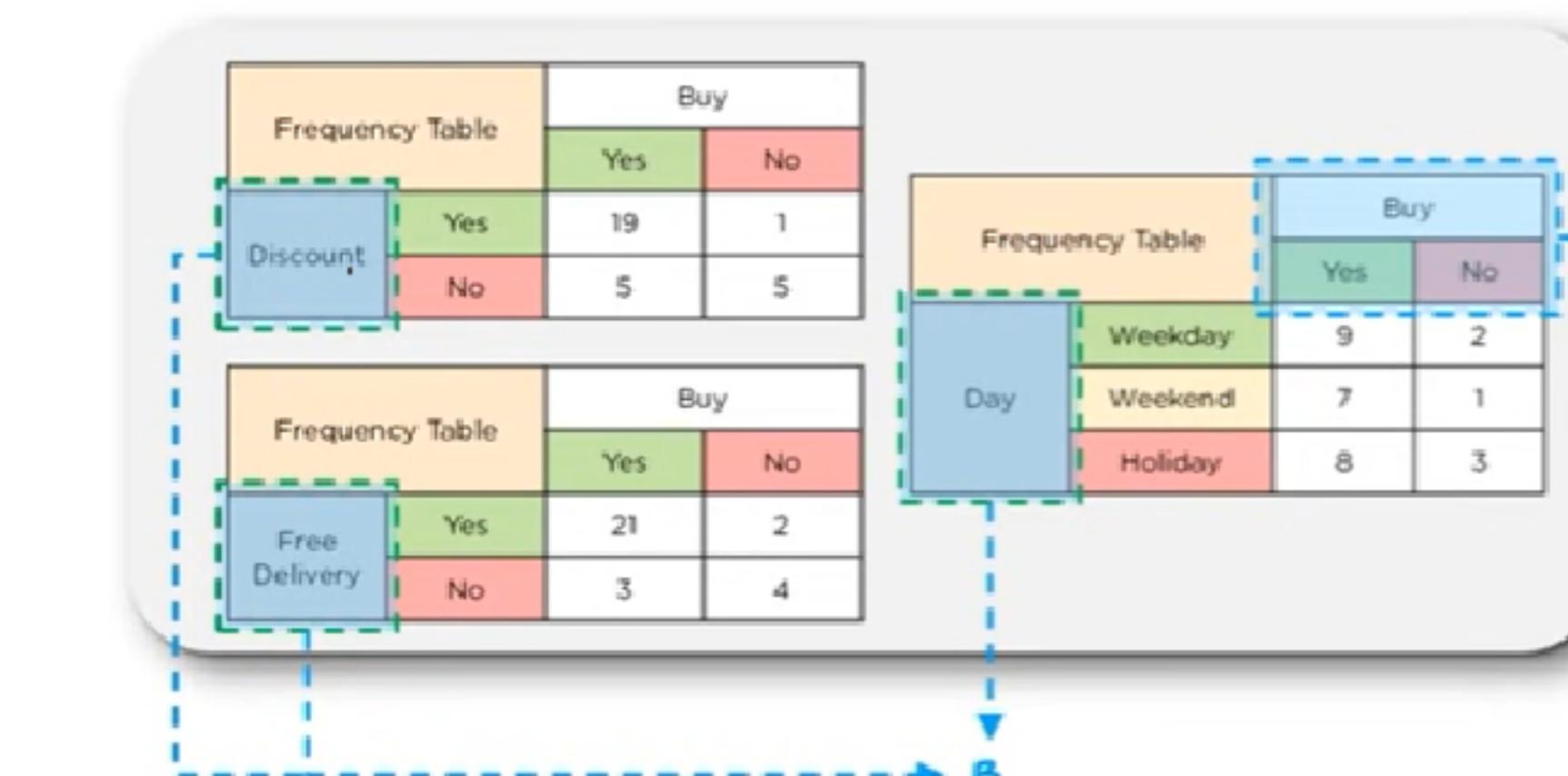
→ step

Pivot tables



Based on this dataset containing three input types of *Day*, *Discount* and *Free Delivery*, we will populate frequency tables for each attribute

FOR OUR BAYES THEOREM, LET THE EVENT **BUY** BE **A** AND THE INDEPENDENT VARIABLES, **DISCOUNT**, **FREE DELIVERY** AND **DAY** BE **B**



Now let us calculate the Likelihood table for one of the variable, Day which includes Weekday, Weekend and Holiday

Step 1

Frequency Table		Buy		
Day	Weekday	9/24	2	11
	Weekend	7/24	1	8
	Holiday	8/24	3	11
		24	6	30

$n(B)$

Likelihood Table		Buy		
Day	Weekday	9/24	2/6	11/30
	Weekend	7/24	1/6	8/30
	Holiday	8/24	3/6	11/30
		24/30	6/30	

$P(B|A)$ \rightarrow target

$$P(B) = P(\text{Weekday}) \\ = 11/30 = 0.37$$

$$P(A) = P(\text{No Buy}) \\ = 6/30 = 0.2$$

$$P(B|A) \\ = P(\text{Weekday} | \text{No Buy}) \\ = 2/6 = 0.33$$

$P(A|B)$

Based on this likelihood table, we will calculate conditional probabilities as below

Frequency Table		Buy		
		Yes	No	
Day	Weekday	9	2	11
	Weekend	7	1	8
	Holiday	8	3	11
		24	6	30

Likelihood Table		Buy		
		Yes	No	
Day	Weekday	9/24	2/6	11/30
	Weekend	7/24	1/6	8/30
	Holiday	8/24	3/6	11/30
		24/30	6/30	

$$P(B) = P(\text{Weekday}) = 11/30 = 0.367$$

$$P(A) = P(\text{No Buy}) = 6/30 = 0.2$$

$$P(B|A) = P(\text{Weekday} \mid \text{No Buy}) = 2/6 = 0.33$$

~~Prob NB~~

$$P(A|B) = P(\text{No Buy} \mid \text{Weekday})$$

$$= P(\text{Weekday} \mid \text{No Buy}) * P(\text{No Buy}) / P(\text{Weekday})$$

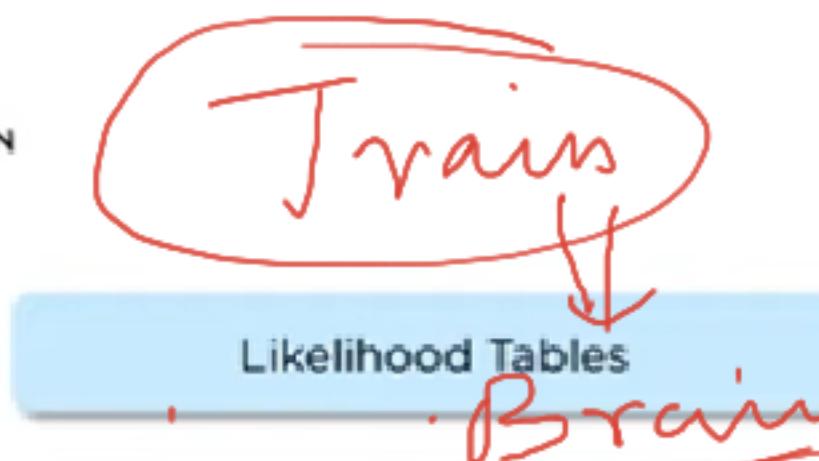
$$= (0.33 * 0.2) / 0.367 = 0.179$$

Z
11

LET US USE THESE 3 LIKELIHOOD TABLES TO CALCULATE WHETHER A CUSTOMER WILL PURCHASE A PRODUCT ON A SPECIFIC COMBINATION OF DAY, DISCOUNT AND FREE DELIVERY OR NOT

HERE, LET US TAKE A COMBINATION OF THESE FACTORS:

- DAY = HOLIDAY
- DISCOUNT = YES
- FREE DELIVERY = YES



P(No Buy)

0.178

P(Buy)

0.822

New

.fit()

~~predict~~

Likelihood \rightarrow w & Target

Calculating Conditional Probability of purchase on the following combination of day, discount and free delivery:

Where B equals:

- Day = Holiday
- Discount = Yes
- Free Delivery = Yes

~~high~~ low

$$P(D \cap F \cap N) \times P(F \cap D \cap N) \times P(H \cap F \cap N)$$

$$\text{Let } A = \text{No Buy}$$

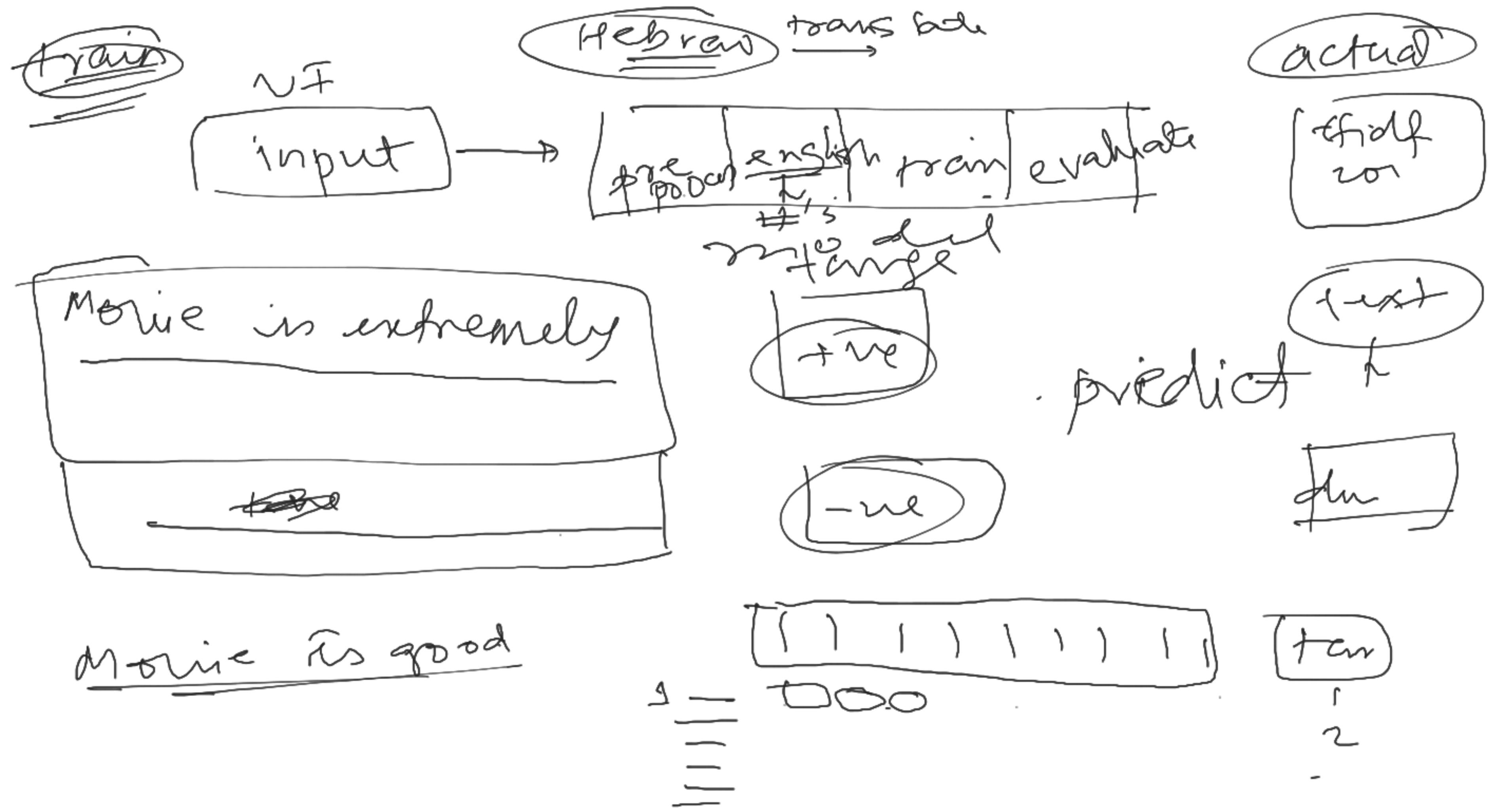
$$P(A|B) = P(\text{No Buy} | \text{Discount} = \text{Yes}, \text{Free Delivery} = \text{Yes}, \text{Day} = \text{Holiday})$$

$$= \frac{P(\text{Discount} = \text{Yes} | \text{No}) * P(\text{Free Delivery} = \text{Yes} | \text{No}) * P(\text{Day} = \text{Holiday} | \text{No}) * P(\text{No Buy})}{P(\text{Discount} = \text{Yes}) * P(\text{Free Delivery} = \text{Yes}) * P(\text{Day} = \text{Holiday})}$$

$$= \frac{(1/6) * (2/6) * (3/6) * (6/30)}{(20/30) * (23/30) * (11/30)}$$

$$= 0.178$$

$$P(Y = 1 | X) = \frac{P(y=1)P(x_1|y=1)\dots P(x_n|y=1)}{P(X)}$$



SUM OF PROBABILITIES

$$= 0.986 + 0.178 = 1.164$$

LIKELIHOOD OF PURCHASE

$$= 0.986 / 1.164 = 84.71 \%$$

LIKELIHOOD OF NO PURCHASE

$$= 0.178 / 1.164 = 15.29 \%$$

PC

PROBABILITY OF PURCHASE = 0.986

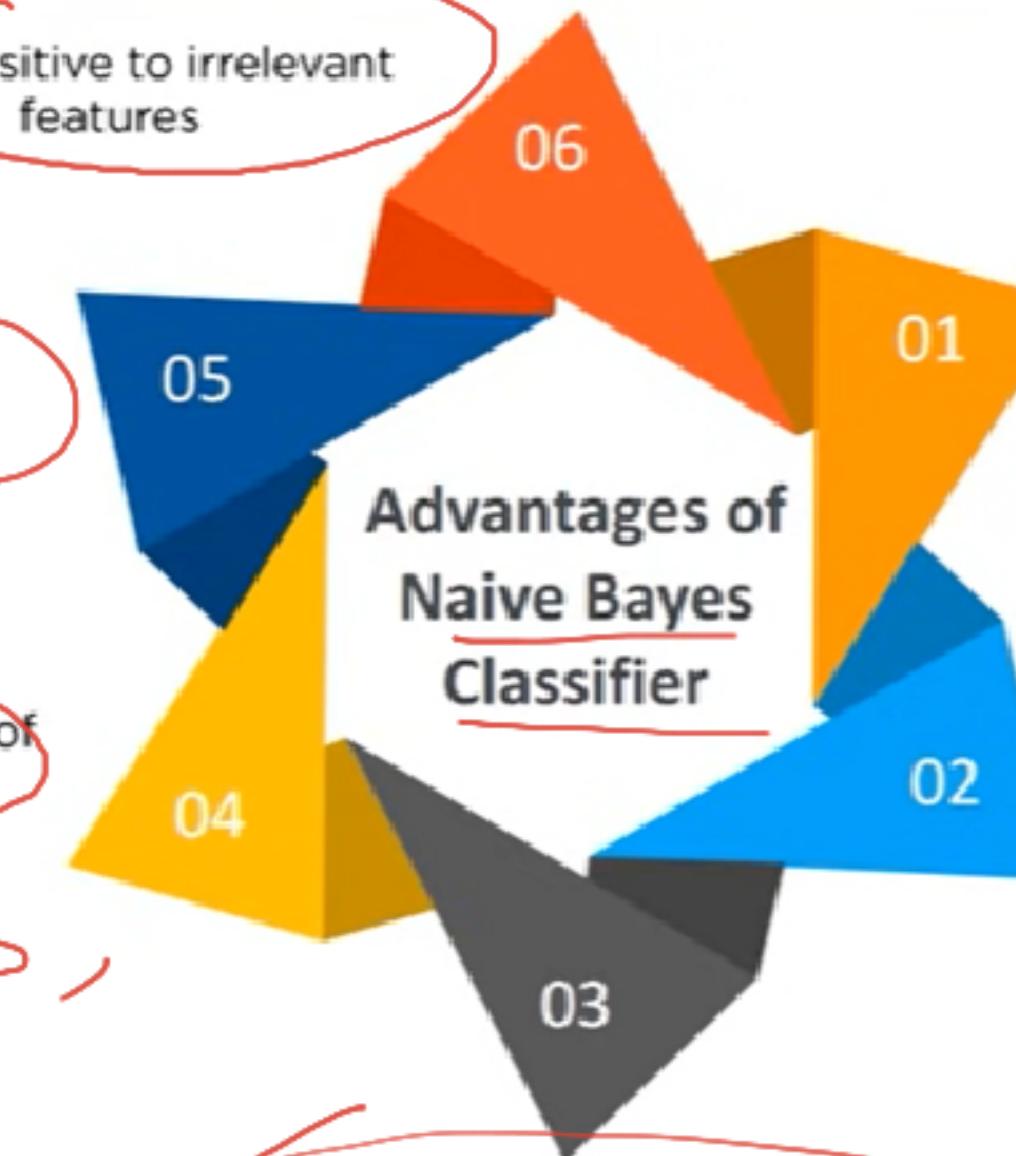
PROBABILITY OF NO PURCHASE = 0.178

Not sensitive to irrelevant features

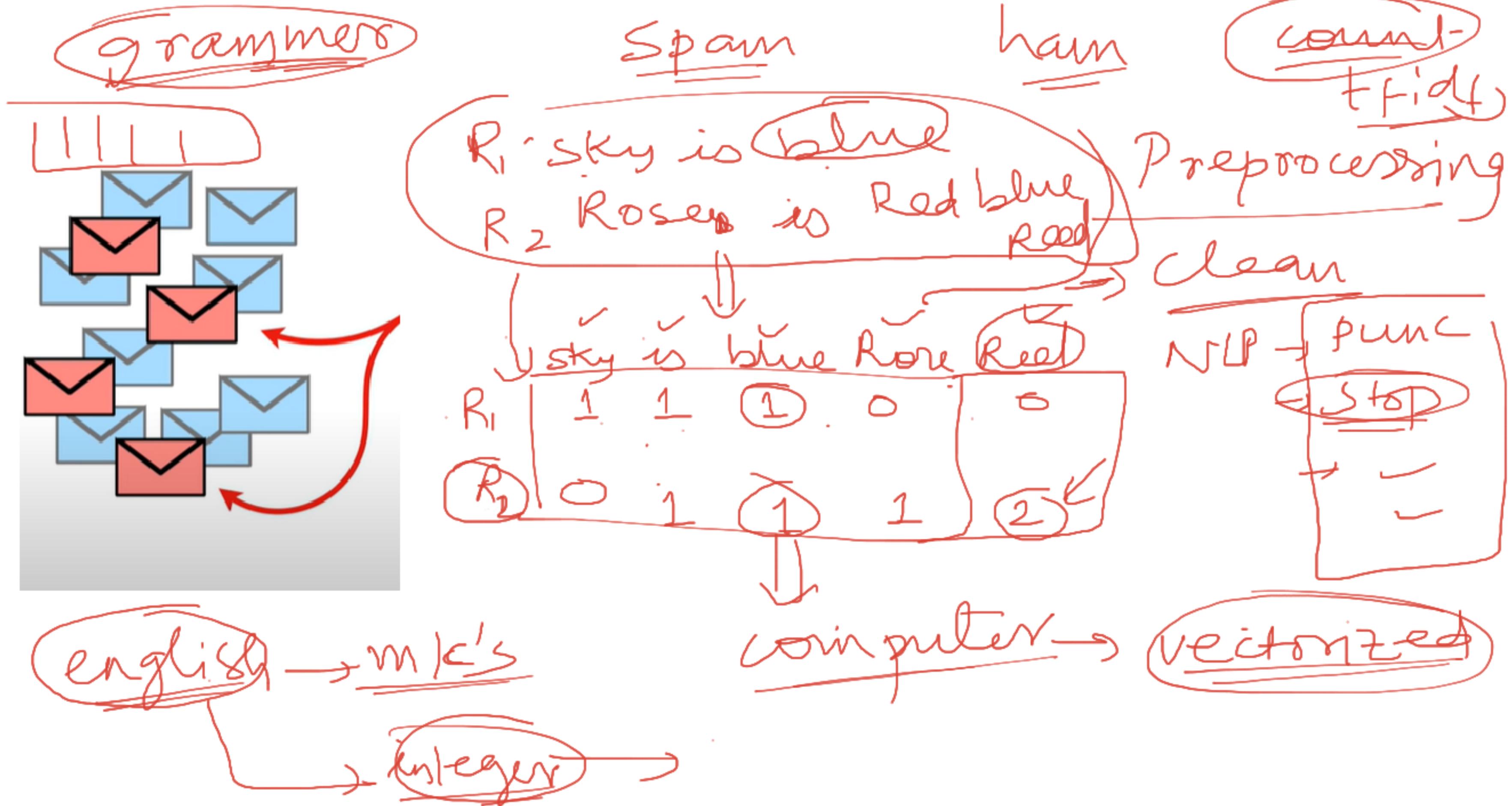
As it is fast, it can be used in real time predictions

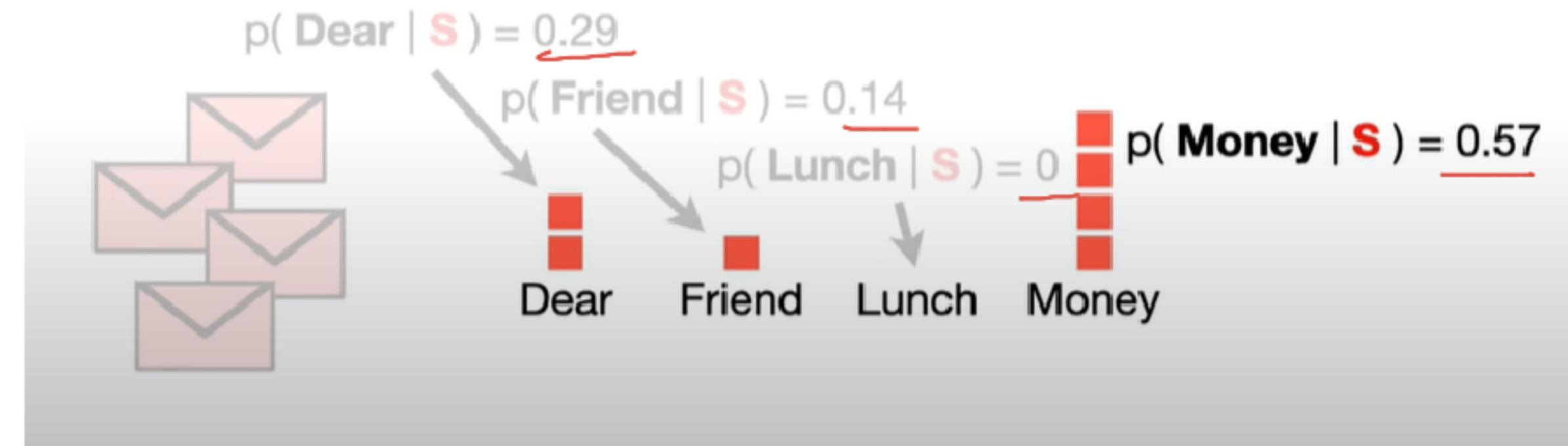
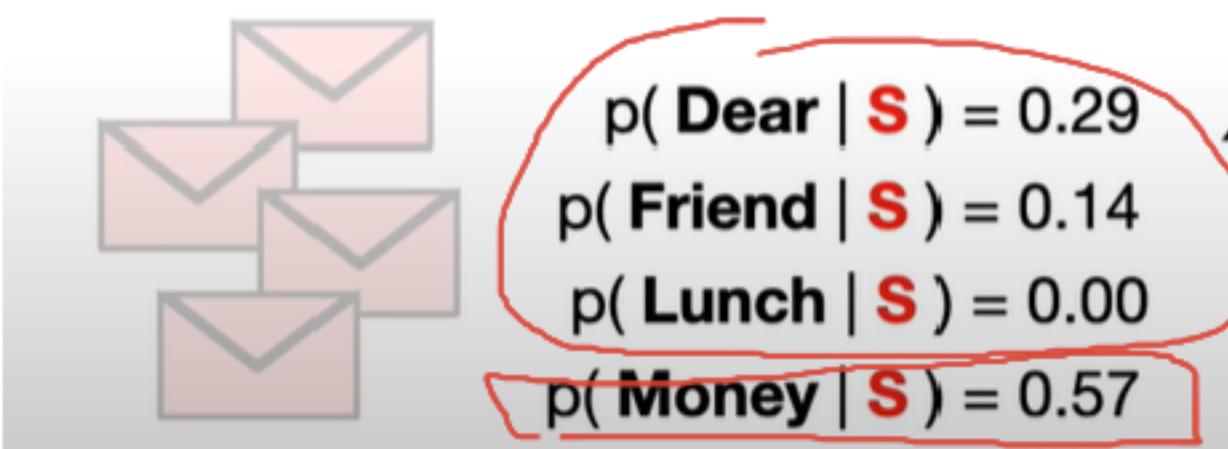
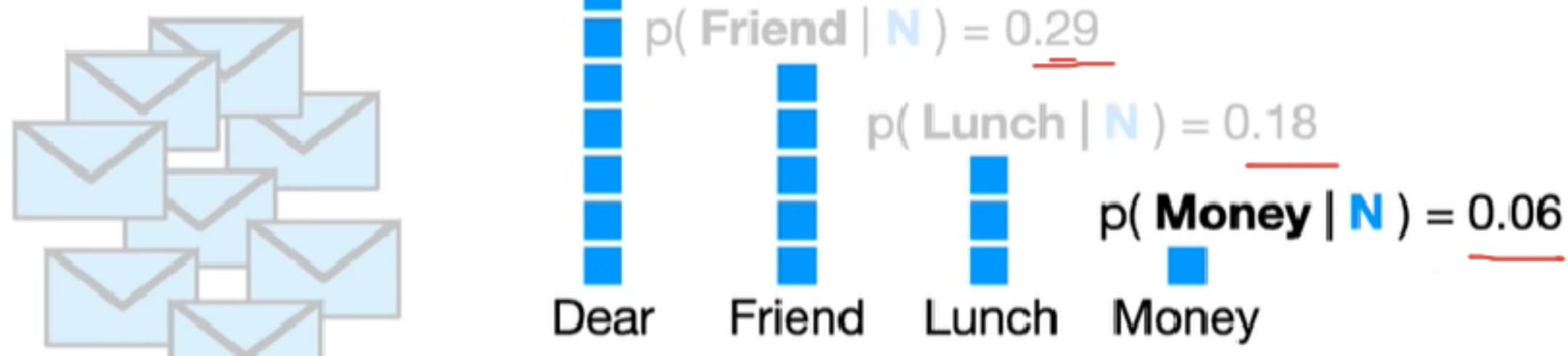
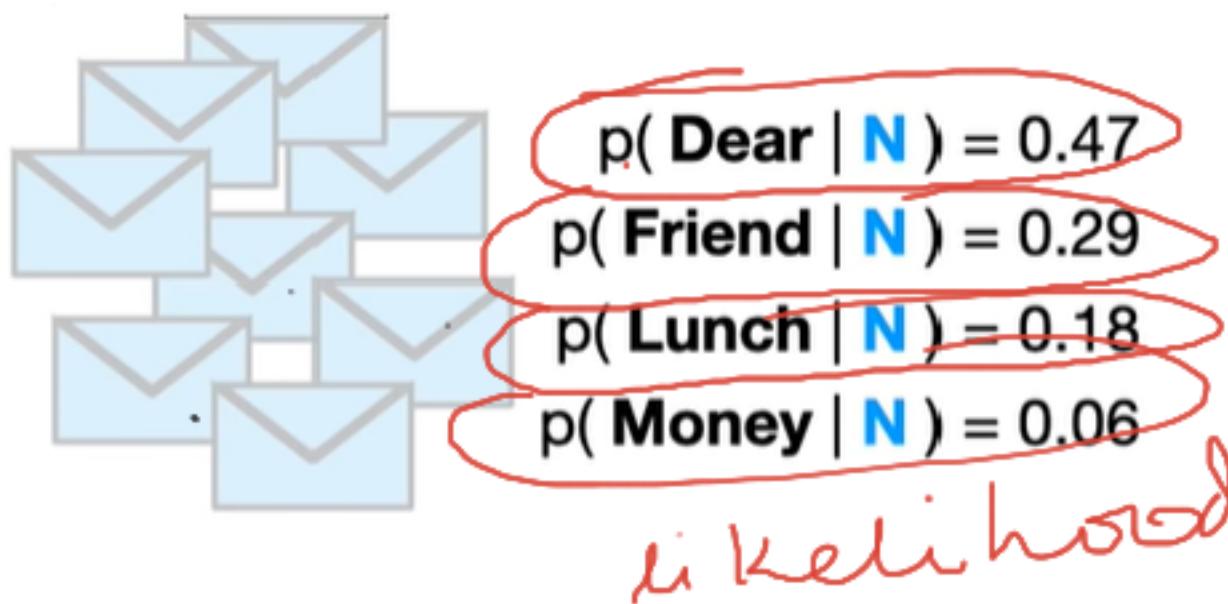
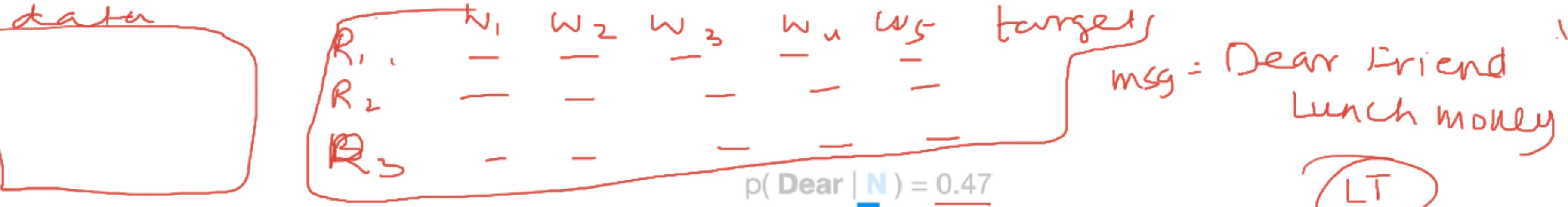
Highly scalable with number of predictors and data points

col rows,



Handles both continuous and discrete data





~~predict msg~~ ~~Dear | Friend~~

$$p(\text{N}) \times p(\text{Dear} | \text{N}) \times p(\text{Friend} | \text{N}) = 0.09$$

$$p(\text{S}) \times p(\text{Dear} | \text{S}) \times p(\text{Friend} | \text{S}) = 0.01$$

L_T

L_F_n



Lunch Money Money Money Money

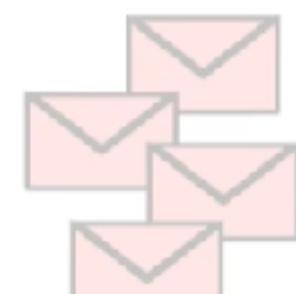


$p(\text{Dear} | \text{N}) = 0.47$
 $p(\text{Friend} | \text{N}) = 0.29$
 $p(\text{Lunch} | \text{N}) = 0.18$
 $p(\text{Money} | \text{N}) = 0.06$

$p(\text{N}) = 0.67$

...and that means we will always classify the messages with **Lunch** in them as **normal**, no matter how many times we see the word **Money**.

$$p(\text{N}) \times p(\text{Lunch} | \text{N}) \times p(\text{Money} | \text{N})^4 = 0.000002$$



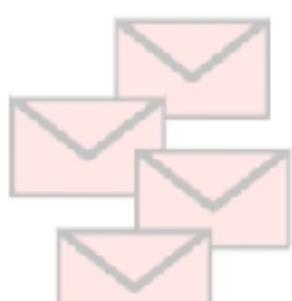
$p(\text{Dear} | \text{S}) = 0.29$
 $p(\text{Friend} | \text{S}) = 0.14$
 $p(\text{Lunch} | \text{S}) = 0.00$
 $p(\text{Money} | \text{S}) = 0.57$

$p(\text{S}) = 0.33$

$$p(\text{S}) \times p(\text{Lunch} | \text{S}) \times p(\text{Money} | \text{S})^4 = 0$$



$$p(\text{N}) = 0.67$$



$$p(\text{S}) = 0.33$$

Lunch Money Money Money Money

$$\begin{aligned} p(\text{Dear} \mid \text{N}) &= 0.43 \\ p(\text{Friend} \mid \text{N}) &= 0.29 \\ p(\text{Lunch} \mid \text{N}) &= 0.19 \\ p(\text{Money} \mid \text{N}) &= 0.10 \end{aligned}$$

$$p(\text{N}) \times p(\text{Lunch} \mid \text{N}) \times p(\text{Money} \mid \text{N})^4 = 0.00001$$

$$p(\text{S}) \times p(\text{Lunch} \mid \text{S}) \times p(\text{Money} \mid \text{S})^4 = 0.00122$$

$$\begin{aligned} p(\text{Dear} \mid \text{S}) &= 0.27 \\ p(\text{Friend} \mid \text{S}) &= 0.18 \\ p(\text{Lunch} \mid \text{S}) &= 0.09 \\ p(\text{Money} \mid \text{S}) &= 0.45 \end{aligned}$$

And since the value for **spam** is greater than the one for a **normal message**...

