



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

Introduction to Statistical Methods

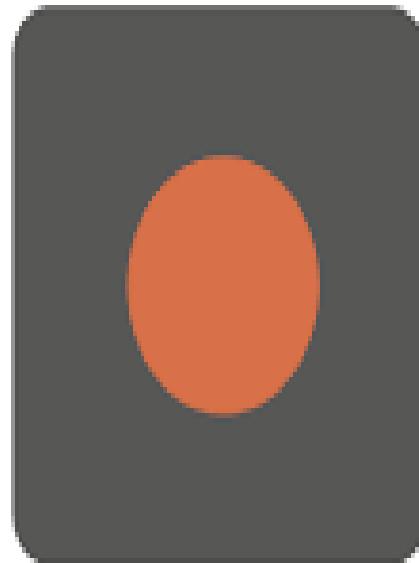
ISM Team



Session 4: Bayes theorem and Naïve Bayes theorem

(Session 4: 10th/11th June 2023)

IMP Note to Self



**Start
Recording**

RECAPTIULATION:

- Introduction to Probability
- Review of Set Theory
- Counting Principles
- Definition of Probability,
- Axioms of Probability,
- Addition Rule of Probability
- Conditional Probability,
- Independent events and,
- Total Probability

Contact Session 4

Contact Session	List of Topic Title	Reference
CS - 4	Bayes theorem(with proof),Introduction to Naïve Bayes concept.	T1 & T2
HW	Problems on Bayes theorem	T1 & T2
Lab	Bayes theorem & Naïve Bayes Concept	Lab 2

Agenda

- Bayes Theorem
- Introduction to Naïve Bayes concept

Text Books

No	Author(s), Title, Edition, Publishing House
T1	Statistics for Data Scientists, An introduction to probability, statistics and Data Analysis, Maurits Kaptein et al, Springer 2022
T2	Probability and Statistics for Engineering and Sciences, 8 th Edition, Jay L Devore, Cengage Learning

BAYES THEOREM:

Let $P = \{A_1, A_2, A_3, \dots, A_k\}$ be a set of exhaustive and mutually exclusive events of a sample space S with $P(A_i) \neq 0$. For each i . If B is any other event associated with A_i with $P(B) \neq 0$, then

$$P\left(\frac{A_i}{B}\right) = \frac{P(A_i)P\left(\frac{B}{A_i}\right)}{P(B)}$$
$$= \frac{P(A_i)P\left(\frac{B}{A_i}\right)}{\sum_{i=1}^n P(A_i)P\left(\frac{B}{A_i}\right)}$$



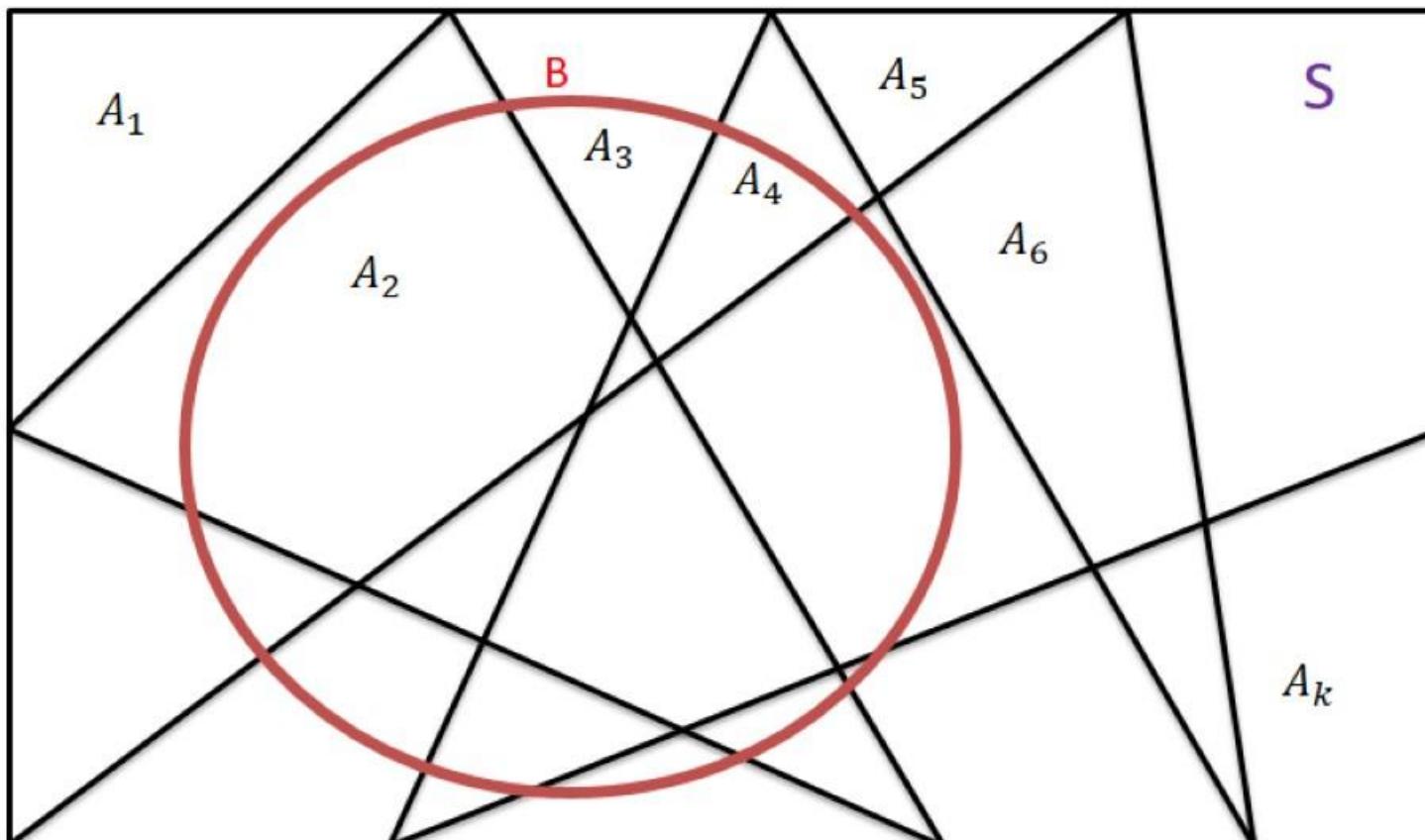
Geometrical Representation of Bayes theorem

innovate

achieve

lead

$$S = A_1 \cup A_2 \cup \dots \cup A_n \quad \text{and} \quad A_1 \cap A_2 \cap \dots \cap A_n = \emptyset$$



$$\therefore B = B \cap S = B \cap \{A_1 \cup A_2 \cup A_3, \dots, \cup A_n\}$$

Proof:

The conditional Probability of A_i for any i given B is given as

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

But we know that

$$P(A_i \cap B) = P(A_i) \cdot P(B | A_i) \text{ & } P(B) = \sum_{i=1}^{i=n} P(A_i) P(B | A_i)$$

$$P(A_i | B) = \frac{P(A_i) \cdot P(B | A_i)}{\sum_{i=1}^{i=n} P(A_i) P(B | A_i)}$$

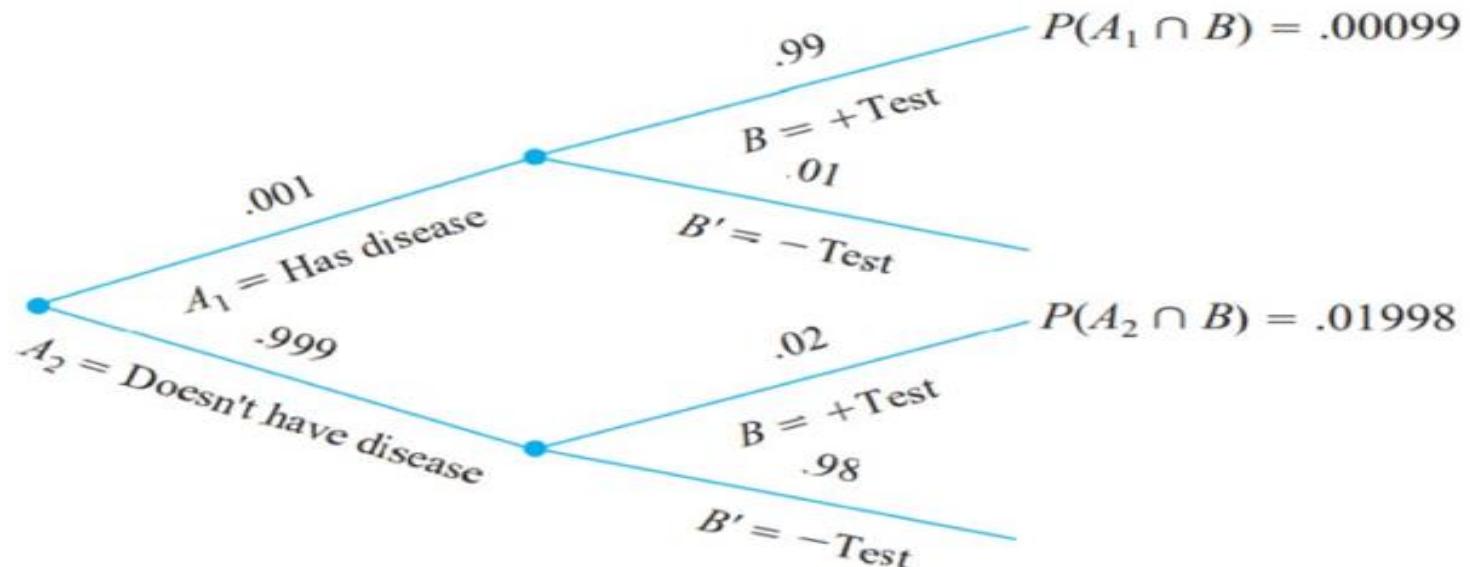
Hence the theorem proved

Example: 1

Only 1 in 1000 adults is afflicted with a rare disease for which a diagnostic test has been developed. The test is such that when an individual actually has the disease, a positive result will occur 99% of the time, whereas an individual without the disease will show a positive test result only 2% of the time. If a randomly selected individual is tested and the result is positive, what is the probability that the individual has the disease?

Solution:

To use Bayes' theorem, let A_1 = individual has the disease, A_2 = individual does not have the disease, and B = positive test result. Then $P(A_1) = .001$, $P(A_2) = .999$, $P(B|A_1) = .99$, and $P(B|A_2) = .02$. The tree diagram for this problem is in Figure



7/16

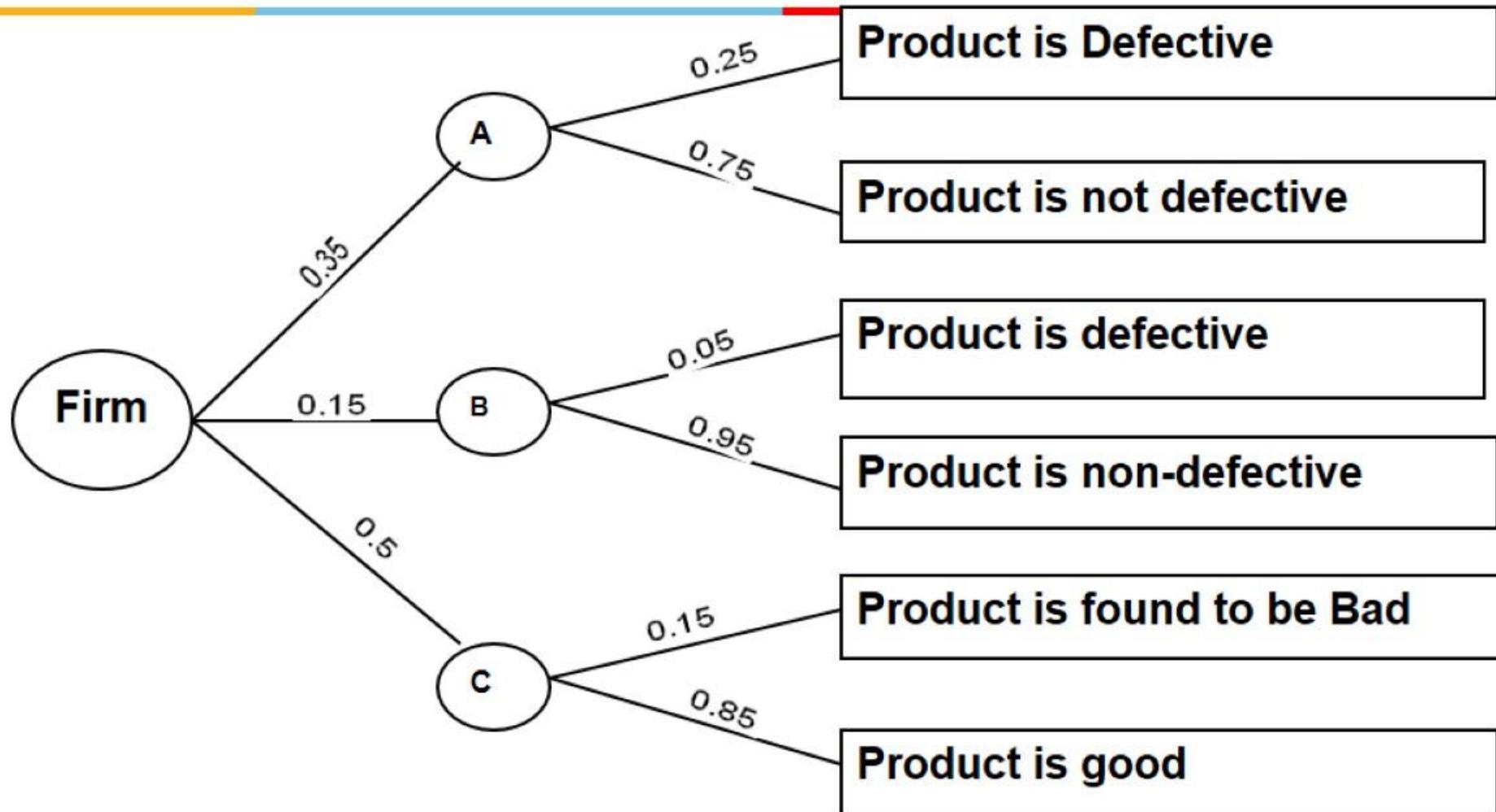
Next to each branch corresponding to a positive test result, the multiplication rule yields the recorded probabilities. Therefore, $P(B) = .00099 + .01998 = .02097$, from which we have

$$P(A_1 | B) = \frac{P(A_1 \cap B)}{P(B)} = \frac{.00099}{.02097} = .047$$

Example: 2

A certain firm has plants A, B, C producing, respectively 35 %, 15% and 50% of the total output. The probabilities of a non – defective product are, respectively, 0.75, 0.95 and 0.85. A Customer receives a bad product, what is the Chance that product came from the plant C?

Tree Diagram



Solution

Let X : “Customer receives a defective product”.

Clearly, $P(X) = P(A)P\left[\frac{X}{A}\right] + P(B)P\left[\frac{X}{B}\right] + P(C)P\left[\frac{X}{C}\right]$
 $= 0.17$

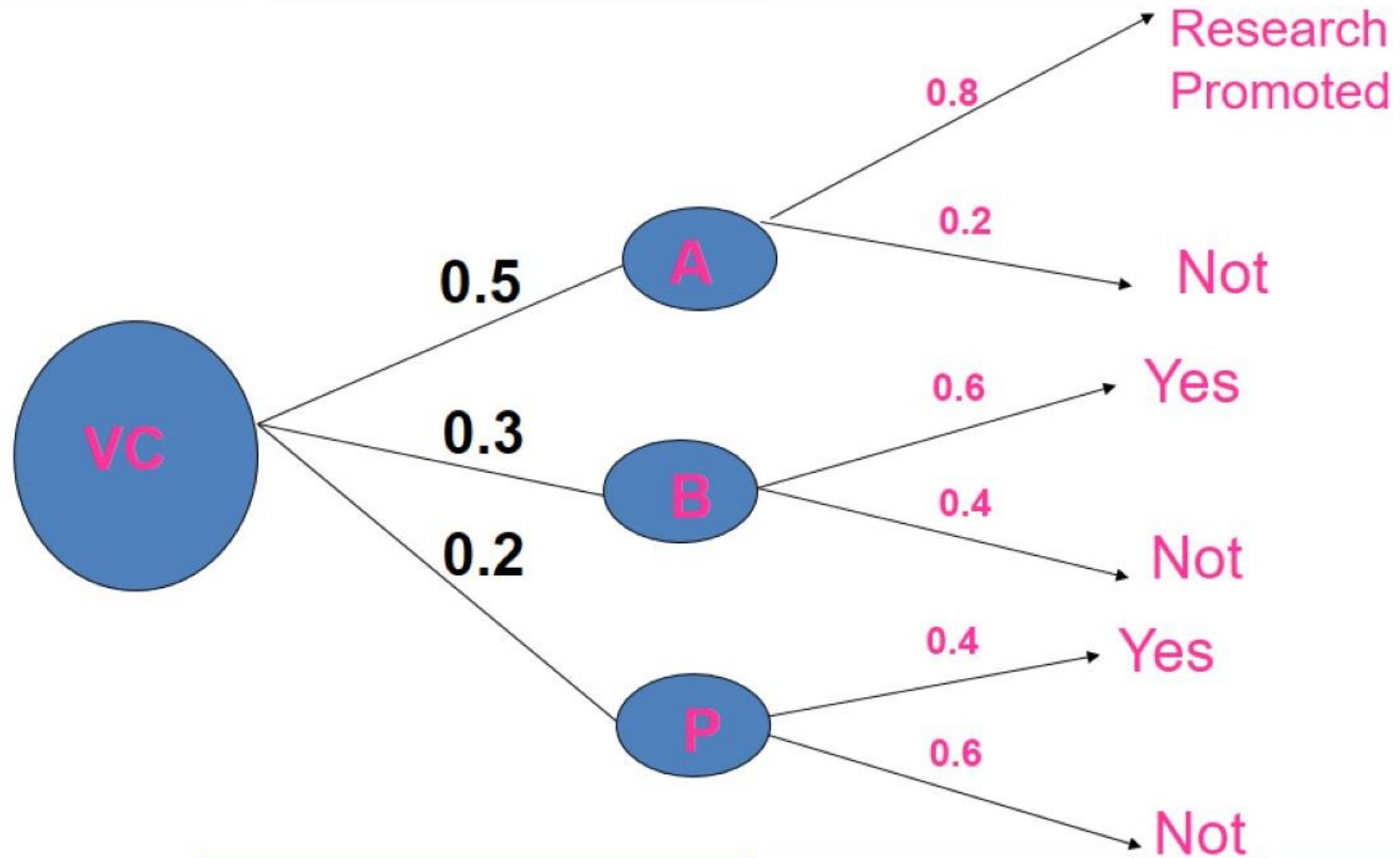
Therefore, the chance that product is manufactured by the plant C is

$$P(C | X) = \frac{P(C \cap X)}{P(X)} = \frac{0.5 \cdot 0.15}{0.17} = 0.4412$$

Example: 2

The chances that an academician, a business man and a politician becoming Vice Chancellor of an university are 0.5, 0.3 and 0.2 respectively. The probability that research work will be promoted in the university by these 3 gentlemen are respectively are 0.8, 0.6 and 0.4. It is found Research work has been promoted by the university. What is the chance that an academician has become the VC?

Tree Diagram



Let **X**: “Research work is promoted”

$$\text{Clearly, } P(X) = 0.5 \times 0.8 + 0.3 \times 0.6 + 0.2 \times 0.4 = 0.66$$

Now to find $P[\text{“An Academician is VC”} / \text{“Research work is promoted i.e. event } X\text{”}] = \frac{0.5 \times 0.8}{0.66} = 0.6061$



Example: A manufacturer of tablets receives its LED screens from three different suppliers, 60% from supplier B_1 , 30% from supplier B_2 , and 10% from supplier B_3 . In other words, the probabilities that any one LED screens received by the plant comes from these three suppliers are 0.60, 0.30, and 0.10. Also suppose that 95% of the LED screens from B_1 , 80% of those from B_2 , and 65% of those from B_3 perform according to specifications.

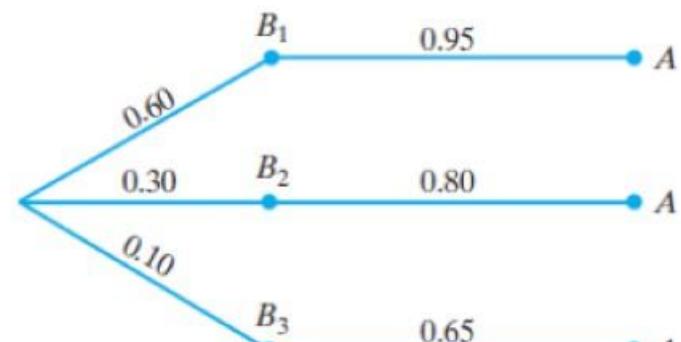
- 1) What is the probability that any one LED screen received by the plant will perform according to specifications?
- 2) Determine the probability that a particular LED screen, which is known to perform according to specifications, came from supplier B_3 .

If A denotes the event that a LED screen received by the plant performs according to specifications, and B_1 , B_2 , and B_3 are the events that it comes from the respective suppliers

$$\begin{aligned} A &= A \cap [B_1 \cup B_2 \cup B_3] \\ &= (A \cap B_1) \cup (A \cap B_2) \cup (A \cap B_3) \\ &= P(A \cap B_1) + P(A \cap B_2) + P(A \cap B_3) \end{aligned}$$

$$\begin{aligned} P(A) &= P(B_1) \cdot P(A | B_1) + \\ &\quad P(B_2) \cdot P(A | B_2) \\ &\quad + P(B_3) \cdot P(A | B_3) \end{aligned}$$

$$\begin{aligned} P(A) &= (0.60)(0.95) + (0.30)(0.80) + (0.10)(0.65) \\ &= 0.875 \end{aligned}$$



$$P(B_3 | A) = \frac{(0.10)(0.65)}{(0.60)(0.95) + (0.30)(0.80) + (0.10)(0.65)} = 0.074$$

Example: 3

Example: Two firms V and W consider bidding on a road-building job, which may or may not be awarded depending on the amounts of the bids. Firm V submits a bid and the probability is $\frac{3}{4}$ that it will get the job provided firm W does not bid. The probability is $\frac{3}{4}$ that W will bid, and if it does, the probability that V will get the job is only $\frac{1}{3}$.
(a) what is the probability that V will get the job? (b) If V gets the job, what is the probability that W did not bid?

Solution:

Answer: Given $P(V/W^1) = 3/4$, $P(W) = 1/3$, $P(V/W) = 1/3$ $P(W^1) = 2/3$

(a) $V = (V \cap W) \cup (V \cap W^1)$ $\Rightarrow P(V) = P(V \cap W) + P(V \cap W^1)$

$$P(V) = P(V/W) P(W) + P(V/W^1) P(W^1) = 11/18$$

(b) $P(W^1/V) = P(V/W^1) P(W^1) / P(V) = 9/11$

Suggested problems:

Example . An office has 4 secretaries handling respectively 20%, 60%, 15% and 5% of the files of all government reports. The probability that they misfile such reports are respectively 0.05, 0.1, 0.1 and 0.05. Find the probability that the misfiled report can be blamed on the first secretary.

Example . In a class 70% are boys and 30% are girls. 5% of boys and 3% of girls are irregular to the classes. What is the probability of a student selected at random is irregular to the classes and what is the probability that the irregular student is a girl?

Example:

Three machines A, B and C produce respectively 60%, 30% and 10% of the total number of items of a factory. The percentage of defective outputs of these machines are 2%, 3% and 4%. An item is selected at random and is found to be defective. (i) Find the probability that the item was produced by machine C? (ii) What is the probability that the item was produced by machine C or B?

Bayesian Learning

- Naive Bayes is a set of simple and efficient machine learning algorithms for solving a variety of classification and regression problems.
- Naive Bayes assumes conditional independence where Bayes theorem does not. This means the relationship between all input features are independent.
- Bayesian learning algorithms that calculate explicit probabilities for hypotheses, such as the naive Bayes classifier, are among the most practical approaches to certain types of learning problems
- For example: Problem of learning to classify text documents such as electronic news articles.
- For such learning tasks, the naive Bayes classifier is among the most effective algorithms known

Features of Bayesian learning

- Prior knowledge can be combined with observed data to determine the final probability of a hypothesis.
- Prior knowledge is provided by asserting
 - ❖ prior probability for each candidate hypothesis, and
 - ❖ probability distribution over observed data for each possible hypothesis.
- New instances can be classified by combining the predictions of multiple hypotheses, weighted by their probabilities.

Bayes Theorem

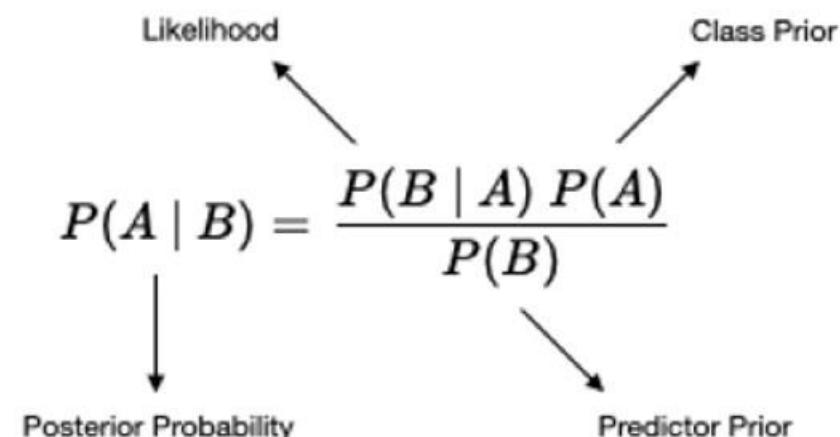
- $P(h)$ = prior probability of hypothesis h , before seeing the training data
- $P(D)$ = prior probability of training data D
- $P(h|D)$ = probability of h given D
- $P(D|h)$ = probability of D given h

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

Machine Learning

- Generative models
 - Build model to estimate the posterior probability $P(Y|X)$ by estimating
 - likelihood of data given target (hypothesis) $P(X|Y)$
 - Prior probabilities over target $P(Y)$
 - In general, for a specific class $Y=c_k$,

$$P(Y = c_k | X) = \frac{P(X|Y = c_k) * P(Y = c_k)}{P(X)}$$



Choosing Hypotheses

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

- Generally want the most probable hypothesis given the training data
- *Maximum a posteriori* hypothesis h_{MAP} :
$$\begin{aligned} h_{MAP} &= \arg \max_{h \in H} P(h|D) \\ &= \arg \max_{h \in H} \frac{P(D|h)P(h)}{P(D)} \\ &= \arg \max_{h \in H} P(D|h)P(h) \end{aligned}$$
- If assume $P(h_i) = P(h_j)$ then can further simplify, and choose the *Maximum likelihood* (ML) hypothesis

$$h_{ML} = \arg \max_{h_i \in H} P(D|h_i)$$

Brute Force MAP Hypothesis

- For each hypothesis h in H , calculate the posterior probability

$$P(h | D) = \frac{P(D | h)P(h)}{P(D)}$$

- Output the hypothesis h_{MAP} with the highest posterior probability

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(h|D)$$

MAP Hypothesis

- Using Bayes theorem, we compute the MAP hypothesis for all probable hypothesis (or all unique class labels)
- Identify the best hypothesis describing the data as

$$\begin{aligned} h_{MAP} &= \operatorname{argmax}_{h \in H} P(h|D) \\ &= \operatorname{argmax}_{h \in H} \frac{P(D|h)P(h)}{P(D)} \\ &= \operatorname{argmax}_{h \in H} P(D|h)P(h) \end{aligned}$$

H: set of all hypothesis

$P(D)$ is independent of h and is same for all hypothesis, therefore dropped

Maximum Likelihood Estimation



- When no prior information is available, all hypothesis are equally likely i.e $p(h_i) = p(h_j)$
- This is also true for a balanced class problem where all the classes are equally likely
- This is known as uniform prior
- MAP hypothesis further simplified to
- $h_{ML} = \operatorname{argmax} P(D|h)$ (where h belongs to H)

Conditional independence

- **Definition:** X is conditionally independent of Y given Z , if the probability distribution governing X is independent of the value of Y , given the value of Z

$$(\forall i, j, k) P(X = x_i | Y = y_j, Z = z_k) = P(X = x_i | Z_k)$$

$$P(X|Y, Z) = P(X|Z)$$

Example:

$$P(\text{Thunder}|\text{Rain, Lightning}) = P(\text{Thunder}|\text{Lightning})$$

Applying conditional independence

Naïve Bayes assumes X_i are conditionally independent given Y

e.g., $P(X_1|X_2, Y) = P(X_1|Y)$

$$\begin{aligned}P(X_1, X_2|Y) &= P(X_1|X_2, Y)P(X_2|Y) \\&= P(X_1|Y)P(X_2|Y)\end{aligned}$$

General form: $P(X_1, \dots, X_n|Y) = \prod_{j=1}^n P(X_j|Y)$

How many parameters to describe $P(X_1, \dots, X_n|Y)$? $P(Y)$?

Without conditional independence assumption?

With conditional independence assumption?

Naïve Bayes Independence assumption

Assumption:

$$P(X_1, \dots, X_n | Y) = \prod_{j=1}^n P(X_j | Y)$$

i.e., X_i and X_j are conditionally independent
given Y for $i \neq j$

Naïve Bayes classifier

- Bayes rule:

$$P(Y = y_k | X_1, \dots, X_n) = \frac{P(Y = y_k)P(X_1, \dots, X_n | Y = y_k)}{\sum_j P(Y = y_j)P(X_1, \dots, X_n | Y = y_j)}$$

- Assume conditional independence among X_i 's:

$$P(Y = y_k | X_1, \dots, X_n) = \frac{P(Y = y_k)\prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j)\prod_i P(X_i | Y = y_j)}$$

- Pick the most probable (MAP) Y

$$\hat{Y} \leftarrow \operatorname{argmax}_{y_k} P(Y = y_k)\prod_i P(X_i | Y = y_k)$$

↑
Prior
Probability ↑
MLE

NAÏVE BAYES CLASSIFIER

- Assume independence among attributes X_i when class is given:
 - ❖ $P(X_1, X_2, \dots, X_d | Y_j) = P(X_1 | Y_j) P(X_2 | Y_j) \dots P(X_d | Y_j)$
 - ❖ Now we can estimate $P(X_i | Y_j)$ for all X_i and Y_j combinations from the training data
 - ❖ New point is classified to Y_j if $P(Y_j) \prod P(X_i | Y_j)$ is maximal.

Example 1:

If the weather is sunny,
then the player will play
or not?

i.e. Play/ Sunny = Yes or No

Note if we know $P(\text{Yes/Sunny})$ and
 $P(\text{No/Sunny})$ then we can answer the
question asked

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No

Steps to Apply Bayes Theorem



Step 1- View or collect “raw” data.

Step 2 - Convert long data to a frequency table

weather	Play		Row Total
	no	yes	
Sunny	2	3	5
Overcast	0	4	4
Rainy	3	2	5
Column Total	5	9	14

Step 3 - Row and column sums to get probabilities

Weather probabilities

$$\text{sunny} = 5/14, \text{rainy} = 5/14$$

$$\text{Overcast} = 4/14$$

Play probabilities

$$\text{no} = 5/14$$

$$\text{yes} = 9/14$$

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No

Steps to Apply Bayes Theorem



		Play		
		no	yes	Row Total
weather				
Sunny	2	3	5	P(Sunny)= 5/14
Overcast	0	4	4	P(Overcast) = 4/14
Rainy	3	2	5	P(Rainy)=5/14
Column Total	5	9	14	
	P(no)=5/14 P(yes)=9/14			

Step 4 - Apply probabilities from frequency table to Bayes theorem

$$P(\text{yes} \mid \text{sunny}) = \frac{P(\text{sunny} \mid \text{yes}) P(\text{yes})}{P(\text{sunny})}$$

weather	no	yes	
Rainy	3	2	$\frac{5}{14} = 0.36$
sunny	2	3	$\frac{5}{14} = 0.36$
overcast	0	4	$\frac{4}{14} = 0.29$
Total	5	9	

$$\frac{5}{14} = 0.36 \quad \frac{9}{14} = 0.64$$

now $P(\text{yes} | \text{sunny}) = \frac{P(\text{sunny} | \text{yes}) P(\text{yes})}{P(\text{sunny})}$

$$= \frac{(3/9)(9/14)}{5/14} = \underline{\underline{0.60}} \quad \checkmark$$

$$P(\text{no} | \text{sunny}) = \frac{(2/5)(5/14)}{5/14} = \underline{\underline{0.40}}$$

Example 2:

If the features of
today = (Outlook is Sunny, Temp is Hot, Humidity is Normal, Windy is False),
 then the player will play or not?

S. No	Outlook	Temp	Humidity	Windy	Play Tennis
1	Rainy	Hot	High	False	No
2	Rainy	Hot	High	True	No
3	Overcast	Hot	High	False	Yes
4	Sunny	Mild	High	False	Yes
5	Sunny	Cool	Normal	False	Yes
6	Sunny	Cool	Normal	True	No
7	Overcast	Cool	Normal	True	Yes
8	Rainy	Mild	High	False	No
9	Rainy	Cool	Normal	False	Yes
10	Sunny	Mild	Normal	False	Yes
11	Rainy	Mild	Normal	True	Yes
12	Overcast	Mild	High	True	Yes
13	Overcast	Hot	Normal	False	Yes
14	Sunny	Mild	High	True	No

today = (Sunny, Hot, Normal, False)

$$\begin{aligned}
 P(Y|X) &= \frac{P(X|Y_{\text{yes}}) P(Y_{\text{yes}})}{P(X)} \\
 &= \frac{P(X_1, X_2, X_3, X_4 | Y_{\text{yes}}) P(Y_{\text{yes}})}{P(X_1, X_2, X_3, X_4)} = \frac{P(X_1 | Y) P(X_2 | Y) P(X_3 | Y) P(X_4 | Y)}{P(X_1) P(X_2) P(X_3) P(X_4)}
 \end{aligned}$$

Outlook

	Yes	No	P(yes)	P(no)
Sunny	2	3	2/9	3/5
Overcast	4	0	4/9	0/5
Rainy	3	2	3/9	2/5
Total	9	5	100%	100%

Temperature

	Yes	No	P(yes)	P(no)
Hot	2	2	2/9	2/5
Mild	4	2	4/9	2/5
Cool	3	1	3/9	1/5
Total	9	5	100%	100%

Humidity

	Yes	No	P(yes)	P(no)
High	3	4	3/9	4/5
Normal	6	1	6/9	1/5
Total	9	5	100%	100%

Wind

	Yes	No	P(yes)	P(no)
False	6	2	6/9	2/5
True	3	3	3/9	3/5
Total	9	5	100%	100%

Play		P(Yes)/P(No)
Yes	9	9/14
No	5	5/14
Total	14	100%

Example 3:

Magazine Promotion	Watch Promotion	Life Insurance Promotion	Credit Card Insurance	Sex
Yes	No	No	No	Male
Yes	Yes	Yes	Yes	Female
No	No	No	No	Male
Yes	Yes	Yes	Yes	Male
Yes	No	Yes	No	Female
No	No	No	No	Female
Yes	Yes	Yes	Yes	Male
No	No	No	No	Male
Yes	No	No	No	Male
Yes	Yes	Yes	No	Female

New Instance: Magazine Promotion = Yes , Watch Promotion = Yes,
 Life Insurance Promotion = No, Credit Card Insurance = No then Sex = ?

$D = \{ \text{magazine promotion, watch promotion, Life insurance promotion, credit card insurance} \}$

$h_i = \text{male or Female}$

$$P(\underline{\text{male}} / \text{Yes, yes, no, no}) = ?$$

$$P(\underline{\text{Female}} / \text{yes, Yes, no, no}) = ?$$

		magazine promotion		watch promotion		L.I promotion		credit card promotion	
		Male	Female	Male	Female	Male	Female	M	F
YES	Male	4	3	2	2	2	3	2	1
	Female	2	1	4	2	4	1	4	3
		Ratios (YES)		$\frac{4}{6}$		$\frac{3}{4}$		$\frac{2}{6}$	
				$\frac{2}{4}$		$\frac{2}{6}$		$\frac{3}{4}$	
		NO		$\frac{2}{6}$		$\frac{1}{4}$		$\frac{4}{6}$	
				$\frac{4}{6}$		$\frac{2}{4}$		$\frac{1}{4}$	
				$\frac{4}{6}$		$\frac{1}{4}$		$\frac{4}{6}$	
				$\frac{3}{4}$					

$$P(\text{male} | E)$$

$$= \frac{P(E | \text{male}) P(\text{male})}{P(E)}$$

$$= \frac{\left(\frac{4}{6} \cdot \frac{2}{6} \cdot \frac{4}{6} \cdot \frac{4}{6} \right) \left(\frac{3}{5} \right)}{P(E)}$$

$$= \frac{0.0593}{P(E)}$$

YES

YES

NO

NO

$$\frac{6}{10} \cdot \frac{3}{5}$$

$$P(\text{Female} | E)$$

$$= \frac{P(E | \text{Female}) P(\text{Female})}{P(E)}$$

$$= \frac{\left(\frac{3}{4}\right)\left(\frac{2}{4}\right)\left(\frac{1}{4}\right)\left(\frac{3}{4}\right) \cdot \frac{2}{5}}{P(E)}$$

$$= \frac{\left(\frac{9}{128}\right)\left(\frac{2}{5}\right)}{P(E)} = \frac{0.0281}{P(E)}$$

$$0.0593 > 0.0281$$

∴ male ✓

Example 4:

Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
whale	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
turtle	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
dolphin	yes	no	yes	no	mammals
eagle	no	yes	no	yes	non-mammals

Give Birth	Can Fly	Live in Water	Have Legs	Class
yes	no	yes	no	?

A: attributes

M: mammals

N: non-mammals

$$P(A | M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$$

$$P(A | N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$$

$$P(A | M)P(M) = 0.06 \times \frac{7}{20} = 0.021$$

$$P(A | N)P(N) = 0.004 \times \frac{13}{20} = 0.0027$$

$$P(A|M)P(M) > P(A|N)P(N)$$

=> Mammals

Issues with Naïve Bayes Classifier

Naïve Bayes Classifier:

$$P(\text{Refund} = \text{Yes} | \text{No}) = 3/7$$

$$P(\text{Refund} = \text{No} | \text{No}) = 4/7$$

$$P(\text{Refund} = \text{Yes} | \text{Yes}) = 0$$

$$P(\text{Refund} = \text{No} | \text{Yes}) = 1$$

$$P(\text{Marital Status} = \text{Single} | \text{No}) = 2/7$$

$$P(\text{Marital Status} = \text{Divorced} | \text{No}) = 1/7$$

$$P(\text{Marital Status} = \text{Married} | \text{No}) = 4/7$$

$$P(\text{Marital Status} = \text{Single} | \text{Yes}) = 2/3$$

$$P(\text{Marital Status} = \text{Divorced} | \text{Yes}) = 1/3$$

$$P(\text{Marital Status} = \text{Married} | \text{Yes}) = 0$$

$$\mid P(\text{Yes}) = 3/10$$

$$P(\text{No}) = 7/10$$

$$\mid P(\text{Yes} | \text{Married}) = 0 \times 3/10 / P(\text{Married})$$

$$P(\text{No} | \text{Married}) = 4/7 \times 7/10 / P(\text{Married})$$

For Taxable Income:

If class = No: sample mean = 110
sample variance = 2975

If class = Yes: sample mean = 90
sample variance = 25

Issues with Naïve Bayes Classifier

Consider the table with Tid = 7 deleted

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7				
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Naïve Bayes Classifier:

$$P(\text{Refund} = \text{Yes} | \text{No}) = 2/6$$

$$P(\text{Refund} = \text{No} | \text{No}) = 4/6$$

→ $P(\text{Refund} = \text{Yes} | \text{Yes}) = 0$

$$P(\text{Refund} = \text{No} | \text{Yes}) = 1$$

$$P(\text{Marital Status} = \text{Single} | \text{No}) = 2/6$$

→ $P(\text{Marital Status} = \text{Divorced} | \text{No}) = 0$

$$P(\text{Marital Status} = \text{Married} | \text{No}) = 4/6$$

$$P(\text{Marital Status} = \text{Single} | \text{Yes}) = 2/3$$

$$P(\text{Marital Status} = \text{Divorced} | \text{Yes}) = 1/3$$

$$P(\text{Marital Status} = \text{Married} | \text{Yes}) = 0/3$$

For Taxable Income:

If class = No: sample mean = 91

sample variance = 685

If class = Yes: sample mean = 90

sample variance = 25

Given $X = (\text{Refund} = \text{Yes}, \text{Divorced}, 120\text{K})$

$$P(X | \text{No}) = 2/6 \times 0 \times 0.0083 = 0$$

$$P(X | \text{Yes}) = 0 \times 1/3 \times 1.2 \times 10^{-9} = 0$$

Naïve Bayes will not be able to
classify X as Yes or No!

Naïve Bayes for Text Classification



- Naïve Bayes is commonly used for text classification
- For a document with k terms $d = (t_1, \dots, t_k)$

Fraction of documents in c

$$P(c|d) = P(c)P(d|c) = P(c) \prod_{t_i \in d} P(t_i|c)$$

- $P(t_i|c)$ = Fraction of terms from all documents in c that are t_i

Number of times t_i appears in some document in c

$$P(t_i|c) = \frac{N_{ic} + 1}{N_c + T}$$

Laplace Smoothing

Total number of terms in all documents in c

Number of unique words (vocabulary size)

- Easy to implement and works relatively well
- Limitation: Hard to incorporate additional features (beyond words).
 - E.g., number of adjectives used.

A Simple Example:

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

Which tag does the sentence *A very close game* belong to? i.e. $P(\text{sports} | \text{A very close game})$

Feature Engineering: Bag of words i.e use word frequencies without considering order

Using Bayes Theorem:

$$P(\text{sports} | \text{A very close game})$$

$$= P(\text{A very close game} | \text{sports}) P(\text{sports})$$

$$P(\text{A very close game})$$

We assume that every word in a sentence is **independent** of the other ones

"close" doesn't appear in sentences of sports tag, So $P(\text{close} | \text{sports}) = 0$, which makes product 0

A Simple Example

Draw | Naive Bayes
Text classification

Text	Tag	
"A great game"	Sports	Which tag does the sentence "A very close game" belong to? i.e. $P(\text{sports} \text{A very close game})$
"The election was over"	Not sports	Feature Engineering: Bag of words i.e use word frequencies without considering order
"Very clean match"	Sports	Using Bayes Theorem:
"A clean but forgettable game"	Sports	$P(\text{sports} \text{A very close game})$ $= P(\text{A very close game} \text{sports}) P(\text{sports})$ ----- $P(\text{A very close game})$
"It was a close election"	Not sports	

We assume that every word in a sentence is **independent** of the other ones

$$P(\text{A very close game}) = P(A) P(\text{very}) P(\text{close}) P(\text{game})$$

$$P(\text{A very close game} | \text{sports}) = P(\text{a} | \text{sports}) P(\text{very} | \text{sports})$$

$$P(\text{close} | \text{sports}) P(\text{game} | \text{sports})$$

"close" doesn't appear in sentences of sports tag, So $P(\text{close} | \text{sports}) = 0$, which makes product 0

$$P(\text{sports} | \text{A very close game}) = \frac{P(\text{A very close game} | \text{sports}) \cdot P(\text{sports})}{P(\text{A very close game})}$$

$$P(\text{not sports} | \text{A very close game})$$

Laplace smoothing

- Laplace smoothing: we add 1 or in general constant k to every count so it's never zero.
- To balance this, we add the number of possible words to the divisor, so the division will never be greater than 1
- In our case, the possible words are ['a', 'great', 'very', 'over', 'it', 'but', 'game', 'election', 'clean', 'close', 'the', 'was', 'forgettable', 'match'].

Apply Laplace Smoothing

Word	P(word Sports)	P(word Not Sports)
a	2+1 / 11+14	1+1 / 9+14
very	1+1 / 11+14	0+1 / 9+14
close	0+1 / 11+14	1+1 / 9+14
game	2+1 / 11+14	0+1 / 9+14

$$\begin{aligned}
 & P(a|Sports) \times P(very|Sports) \times P(close|Sports) \times P(game|Sports) \times \\
 & P(Sports) \\
 & = 2.76 \times 10^{-5} \\
 & = 0.0000276
 \end{aligned}$$

$$\begin{aligned}
 & P(a|Not\ Sports) \times P(very|Not\ Sports) \times P(close|Not\ Sports) \times \\
 & P(game|Not\ Sports) \times P(Not\ Sports) \\
 & = 0.572 \times 10^{-5} \\
 & = 0.00000572
 \end{aligned}$$

Example :

Doc No	Text
1	I LOVED THE MOVIE
2	I HATED THE MOVIE
3	A GREAT MOVIE ,GOOD MOVIE
4	POOR ACTING
5	GREAT ACTING , A GOOD MOVIE
NEW	I HATED THE POOR ACTING

Example :

Doc No	Text	
1	I LOVED THE MOVIE	_POSITIVE
2	I HATED THE MOVIE	_NEGATIVE
3	A GREAT MOVIE ,GOOD MOVIE	_POSITIVE
4	POOR ACTING	_NEGATIVE
5	GREAT ACTING , A GOOD MOVIE	_POSITIVE
NEW	I HATED THE POOR ACTING	_????

$$P(c/x)$$

$$P(+ / \text{I hated the poor acting}) =$$

$$P(- / \text{I hated the poor acting}) =$$

Based on these probabilities, we can decide the class which the new text belongs

$P(+ \mid \text{I hated the acting})$

i.e. $P(c_1 \mid x) = \frac{P(x \mid c_1) P(c_1)}{\cancel{P(x)}}$

$$= P(\text{I} \mid +) P(\text{hated} \mid +) P(\text{the} \mid +) P(\cancel{\text{acting}} \mid +) P(+)$$

$$\frac{P(\text{I}, +)}{P(+)}$$

$$\frac{P(\text{hated}, +)}{P(+)}$$

words	positive	negative
I	1	1
loved	1	0
the	1	1
movie	4	1
hated	0	1
a	2	0
great	2	0
Poor	0	1
acting	1	1
good	2	0

$$P(\pm | +)$$

$$= \frac{1+1}{14+10}$$

$$P(I | -)$$

$$= \frac{1+1}{6+10}$$

"I hated the poor
acting"

word

positive

negative

I

$$\frac{1+1}{14+10} = 0.0833$$

$$\frac{1+1}{6+10} = 0.125$$

hated

$$\frac{0+1}{14+10} = 0.0417$$

$$\frac{1+1}{6+10} = 0.125$$

the

$$\frac{1+1}{14+10} = 0.0833$$

$$\frac{1+1}{6+10} = 0.125$$

poor

$$\frac{0+1}{14+10} = 0.0417$$

$$\frac{1+1}{6+10} = 0.125$$

acting

$$\frac{1+1}{14+10} = 0.0833$$

$$\frac{1+1}{6+10} = 0.125$$

$x : I$ hate the Poor
acting

$$P(+|x)$$

$$= () () () () () * P(+)$$

\downarrow
 $3/5$

$$= 6.03 \times 10^{-7}$$

$$P(-|x)$$

$$= () () () () () () P(-)$$

\downarrow
 $2/8$

$$= 1.22 \times 10^{-5}$$

\therefore negative class

Example:

Suppose we got the new message with the words '**Dear Friend**', Decide whether this new message is a normal or spam message?

i.e. Normal/ Dear, Friend = Yes or No

Note if we know $P(\text{Normal} / \text{Dear, Friend})$ and $P(\text{Spam} / \text{Dear, Friend})$ then we can answer the question asked

Email word	Spam	Email word	Spam
Dear	Yes	Friend	No
Friend	No	Friend	Yes
Dear	No	Dear	No
Dear	No	Lunch	No
Dear	No	Friend	No
Friend	No	Dear	No
Lunch	No	Dear	No
Friend	No	Dear	No
Lunch	No	Dear	No
Dear	Yes	Money	Yes
Money	Yes	Money	No
Money	Yes	Money	Yes

Steps to Apply Bayes Theorem

Step 1- View or collect “raw” data.

Step 2 - Convert long data to a frequency table

word	Play			Row Total
	normal	spam		
Dear	8	2		10
Friend	5	1		6
Lunch	3	0		3
Money	1	4		5
Column Total	17	7		24

Step 3 - Row and column sums to get probabilities

Email word	Spam	Email word	Spam
Dear	Yes	Friend	No
Friend	No	Friend	Yes
Dear	No	Dear	No
Dear	No	Lunch	No
Dear	No	Friend	No
Friend	No	Dear	No
Lunch	No	Dear	No
Friend	No	Dear	No
Lunch	No	Dear	No
Dear	Yes	Money	Yes
Money	Yes	Money	No
Money	Yes	Money	Yes

Steps to Apply Bayes Theorem



Play

	normal	spam	Row Total
word			
Dear	8	2	10
Friend	5	1	6
Lunch	3	0	3
Money	1	4	5
Column Total	17	7	24

As $P(N/D, F) > P(S/D, F)$, we can decide that Dear Friend is Normal message.

Step 4 - Apply probabilities from frequency table to Bayes theorem

$$P(N/D, F) = \frac{P(D, F/N) \cdot P(N)}{P(D, F)} = \frac{P(D/N) \cdot P(F/N) \cdot P(N)}{P(D) \cdot P(F)} = \frac{\left(\frac{8}{17}\right) \cdot \left(\frac{5}{17}\right) \cdot \left(\frac{17}{24}\right)}{\left(\frac{10}{24}\right) \cdot \left(\frac{6}{24}\right)} = \frac{0.098}{0.104} = 0.9423$$

$$P(S/D, F) = \frac{P(D, F/S) \cdot P(S)}{P(D, F)} = \frac{P(D/S) \cdot P(F/S) \cdot P(S)}{P(D) \cdot P(F)} = \frac{\left(\frac{2}{7}\right) \cdot \left(\frac{1}{7}\right) \cdot \left(\frac{7}{24}\right)}{\left(\frac{10}{24}\right) \cdot \left(\frac{6}{24}\right)} = \frac{0.012}{0.104} = 0.1153$$

Example continued:

Suppose we got the new message contains the word '**Lunch Money Money Money Money**' , Decide whether this new message is a normal or spam message?

Email word	Spam	Email word	Spam
Dear	Yes	Friend	No
Friend	No	Friend	Yes
Dear	No	Dear	No
Dear	No	Lunch	No
Dear	No	Friend	No
Friend	No	Dear	No
Lunch	No	Dear	No
Friend	No	Dear	No
Lunch	No	Dear	No
Dear	Yes	Money	Yes
Money	Yes	Money	No
Money	Yes	Money	Yes

Steps to Apply Bayes Theorem



Play

	normal	spam	Row Total
word			
Dear	8	2	10
Friend	5	1	6
Lunch	3	0	3
Money	1	4	5
Column Total	17	7	24

We can observe that we have to classify any message with Lunch as Normal message, no matter how many times we see the word Money and that's the problem.

To work around this problem add 1 count to the frequency table to each word(Laplace smoothing)

Step 4 - Apply probabilities from frequency table

$$P(N) \cdot P(L/N) \cdot P(M/N)^4 = \left(\frac{17}{24}\right) \cdot \left(\frac{3}{17}\right) \cdot \left(\frac{1}{17}\right)^4 = 0.0000015$$

$$P(S) \cdot P(L/S) \cdot P(M/S)^4 = \left(\frac{7}{24}\right) \cdot \left(\frac{0}{7}\right) \cdot \left(\frac{4}{7}\right)^4 = 0$$

Steps to Apply Bayes Theorem



Play

	normal	spam	Row Total
word			
Dear	8+1	2+1	12
Friend	5+1	1+1	8
Lunch	3+1	0+1	5
Money	1+1	4+1	7
Column Total	21	11	32

As $P(S/L, M^4) > P(N/L, M^4)$, we can decide that
Lunch Money Money Money Money is Spam
message.

Step 4 - Apply probabilities from frequency table

$$P(N) \cdot P(L/N) \cdot P(M/N)^4 = \left(\frac{21}{32}\right) \cdot \left(\frac{4}{21}\right) \cdot \left(\frac{2}{21}\right)^4 = 0.00001$$

$$P(S) \cdot P(L/S) \cdot P(M/S)^4 = \left(\frac{11}{32}\right) \cdot \left(\frac{1}{11}\right) \cdot \left(\frac{5}{11}\right)^4 = 0.00133$$

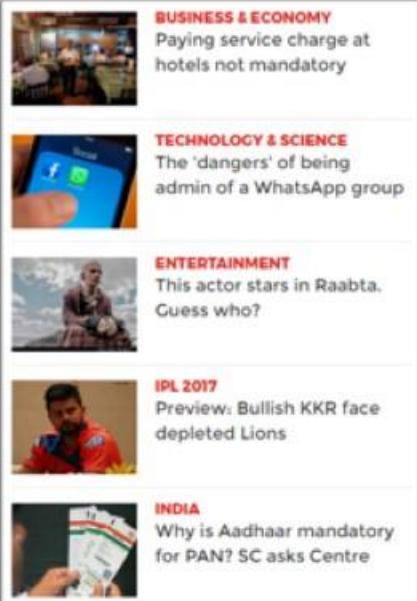
Naïve Bayes Classifier Applications

innovate

achieve

lead

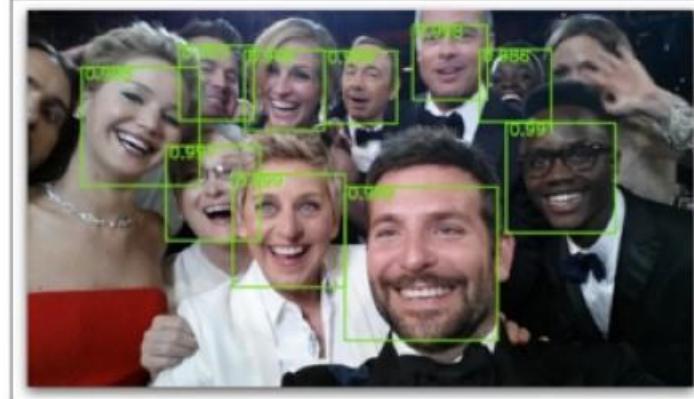
Categorizing News



Email Spam Detection



Face Recognition



Sentiment Analysis



Naive Bayes Classifier

- ✓ Along with decision trees, neural networks, one of the most practical learning methods.
- ✓ When to use
 - ✓ Moderate or large training set available
 - ✓ Attributes that describe instances are conditionally independent given classification
- ✓ Successful applications:
 - ✓ Diagnosis
 - ✓ Classifying text documents

Learning to Classify Text

➤ Why?

- ❖ Learn which news articles are of interest
- ❖ Learn to classify web pages by topic

➤ Naive Bayes is among most effective algorithms

➤ What attributes shall we use to represent text documents??

Baseline: Bag of Words Approach

the world of

TOTAL



all about the company

Our energy exploration, production, and distribution operations span the globe, with activities in more than 100 countries.

At TOTAL, we draw our greatest strength from our fast-growing oil and gas reserves. Our strategic emphasis on natural gas provides a strong position in a rapidly expanding market.

Our expanding refining and marketing operations in Asia and the Mediterranean Rim complement already solid positions in Europe, Africa, and the U.S.

Our growing specialty chemicals sector adds balance and profit to the core energy business.

All About The Company

- ▶ All About The Company
- Global Activities
- Corporate Structure
- TOTAL's Story
- Upstream Strategy
- Downstream Strategy
- Chemicals Strategy
- TOTAL Foundation
- Homepage



aardvark	0
about	2
all	2
Africa	1
apple	0
anxious	0
...	
gas	
...	
oil	
...	
Zaire	

- Require initial knowledge of many probabilities
 - Often estimated based on background knowledge, previously available data, and assumptions about the form of the underlying distributions.
- Significant computational cost required to determine the Bayes optimal hypothesis in the general case (linear in the number of candidate hypotheses)

HW: Exercise

Consider the car theft problem with attributes Color, Type, Origin, and the target, Stolen can be either Yes or No.

we need to classify whether the car is stolen, given the features of the car.

Given the Red color Domestic SUV car Find the probability of whether the car is stolen?

Color	Type	Origin	Stolen?
Red	Sports	Domestic	Yes
Red	Sports	Domestic	No
Red	Sports	Domestic	Yes
Yellow	Sports	Domestic	No
Yellow	Sports	Imported	Yes
Yellow	SUV	Imported	No
Yellow	SUV	Imported	Yes
Yellow	SUV	Domestic	No
Red	SUV	Imported	No
Red	Sports	Imported	Yes

HW: Exercise

If the weather is Snowy,
then the player will play
or not?

weather	Player play
Sunny	yes
Rainy	no
Cloudy	yes
Sunny	no
Sunny	yes
snowy	no
Rainy	yes
Cloudy	no
Cloudy	yes
Sunny	yes
snowy	no
Cloudy	yes
Rainy	no
snowy	no
snowy	yes





Thanks
