



**BITS Pilani**  
Pilani|Dubai|Goa|Hyderabad



# **M Tech(Data Science & Engineering) Introduction to Statistical Methods[ISM]**

**T Vamsidar**



## **Webinar Session No - 1**

**31<sup>st</sup> May -2022**

**Timimgs : 7.30 to 9.00 PM**

# Learning objectives



- Problems on Basic probability
- Conditional Probability
- Baye's theorem
- Random variables

# Problem - 1



Suppose that 55% of all adults regularly consume coffee, 45% regularly consume carbonated soda, and 70% regularly consume at least one of these two products.

a. What is the probability that a randomly selected adult regularly consumes both coffee and soda?  $\rightarrow P(A \cap B)$

b. What is the probability that a randomly selected adult doesn't regularly consume at least one of these two products?

$$\rightarrow P(\overline{A \cup B})$$

## Solution -1



Let **A** be the event that an adult consumes coffee.

Let **B** be the event that an adult consumes carbonated soda.

Given probabilities of the events.

$$\underline{P(A)} = \underline{0.55}, \underline{P(B)} = \underline{0.45}$$

$$and \underline{P(A \cup B)} = \underline{0.7}$$

## (a) Probability of consuming both coffee and soda

The case is intersection of event A and B,

$$\begin{aligned}P(A \cap B) &= P(A) + P(B) - P(A \cup B) \\&= 0.55 + 0.45 - 0.7 \\&= 0.3\end{aligned}$$

## Solution -1



(b) The given case represents compliment of the event when the adult consumes at least one of the drink.

$$\begin{aligned} P[(A \cup B)^c] &= 1 - P(A \cup B) \\ &= 1 - 0.7 = 0.3 \end{aligned}$$

## Problem - 2



Let  $\textcircled{A}$  denote the event that the next request for assistance from a statistical software consultant relates to the SPSS package, and let  $\textcircled{B}$  be the event that the next request is for help with SAS. Suppose that  $P(A) = 0.30$  and  $P(B) = 0.50$

- Calculate  $P(A')$ .
- Calculate  $P(A \cup B)$ .
- Calculate  $P(A' \cap B')$ .

$$P(A \cap B) = P(A) \cdot P(B)$$

## Solution -2

Let  $A$  be the event that the next request for assistance from a statistical software consultant relates to the SPSS package.

Let  $B$  be the event that the next request is for help with SAS.

$$P(A) = 0.30 \text{ and } P(B) = 0.50$$

---

i.  $P(A') = 1 - P(A)$   
=  $1 - 0.30$   
= **0.70**

ii.  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$   
=  $P(A) + P(B) - \{P(A) * P(B)\}$  [A and B are independent events]  
=  $0.3 + 0.5 - [0.3 * 0.5]$   
= **0.65**

✓ ✓

$$\text{iii. } P(\bar{A} \cap \bar{B}) = P[(A \cup B)^c]$$
$$= 1 - P(A \cup B)$$
$$= 1 - 0.65 = 0.35$$



# Problem - 3



The route used by a certain motorist in commuting to work contains two intersections with traffic signals. The probability that he must stop at the first signal is .4, the Analogous probability for the second signal is .5, and the probability that he must stop at least one of the two signals is .6.

What is the probability that he must stop

- a. At both signals?  $\rightarrow P(A \cap B)$
- b. At the first signal but not at the second one?  $\rightarrow P(A \cap B^c)$
- c. At exactly one signal?

## Solution -3



Let **A** be the event that he must stop at the first signal

Let **B** be the event that he must stop at the second signal

$$P(A) = 0.40 \text{ and } P(B) = 0.50$$

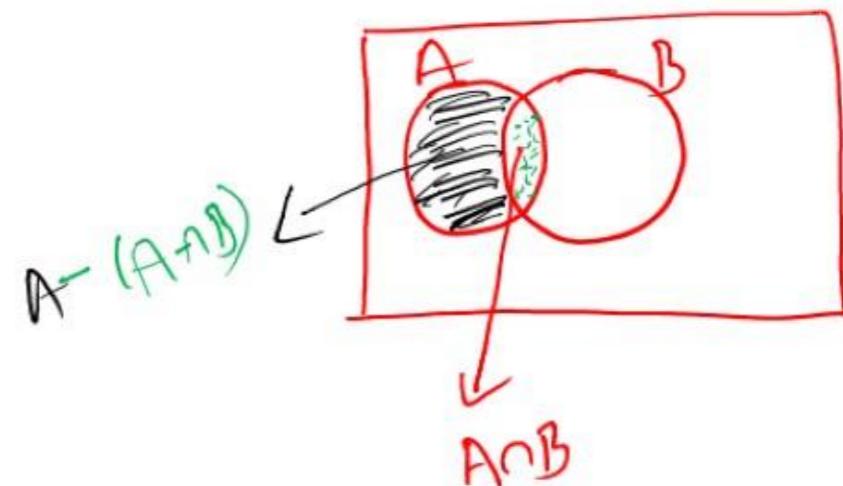
And given that

$$P(A \cup B) = 0.6$$

i.  $P(A \cap B) = P(A) + P(B) - P(A \cup B)$

$$= 0.4 + 0.5 - 0.6 = 0.3$$

$$\begin{aligned}\text{(ii)} \quad P(A \cap B^c) &= P(A) - P(A \cap B) \\ &= 0.4 - 0.3 \\ &= 0.1\end{aligned}$$



$$\text{ii. } P(A \cap B') = P(A) - P(A \cap B)$$

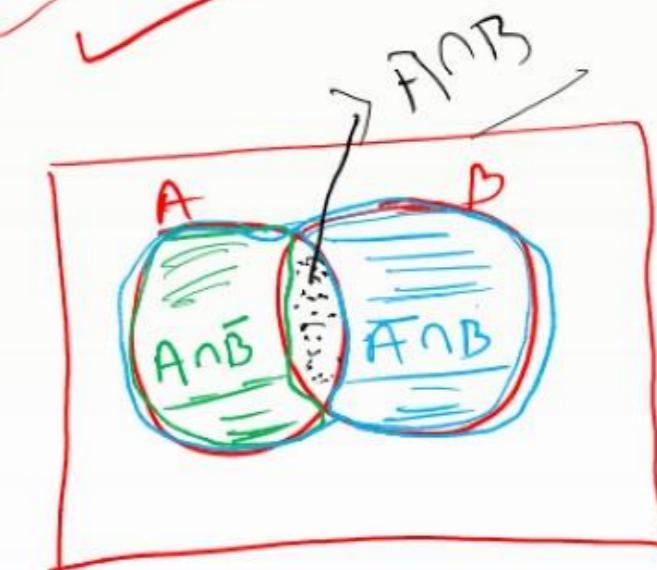
$$= 0.4 - 0.3 = 0.1$$

$$\begin{aligned} & P\{(A \cap B^c) \cup (\bar{A} \cap \bar{B})\} \\ &= P(A \cap B^c) + P(\bar{A} \cap \bar{B}) \\ &= P(A) - P(A \cap B) + P(\bar{B}) - P(A \cap B) \\ &= P(A \cup B) - P(A \cap B) \end{aligned}$$



$$\text{iii. } P(\text{At exactly one signal}) = P(A \cup B) - P(A \cap B)$$

$$= 0.6 - 0.3 = 0.3$$



## Problem - 4

The odds that person X speaks the truth are 3:2 and the odds that person Y speaks the truth are 5:3. In what percentage of cases are they likely to contradict each on an identical point?

# Solution -4



Let  $A = X$  speaks the truth

$\bar{A} = X$  tell a lie

Let  $B = Y$  speaks the truth

$\bar{B} = Y$  tell a lie

$$P(A) = \frac{3}{3+2}; P(\bar{A}) = \frac{2}{3+2}$$

$$P(B) = \frac{5}{5+3}; P(\bar{B}) = \frac{3}{5+3}$$

The event C that X and Y contradict each other on an-identical point.

That can happen in two ways

(i) X speaks the truth and Y tell a lie, i.e,  $A \cap \bar{B}$

(ii) X tell a lie and Y speaks the truth, i.e,  $\bar{A} \cap B$

# Solution -4



the events  $(A \cap \bar{B})$  and  $(\bar{A} \cap B)$  mutually exclusive events

$$\begin{aligned}P(C) &= P(A \cap \bar{B}) + P(\bar{A} \cap B) \\&= P(A) \cdot P(\bar{B}) + P(\bar{A}) \cdot P(B)\end{aligned}$$

[ A and B are independent events]

$$\begin{aligned}&= \frac{3}{5} \cdot \frac{3}{8} + \frac{2}{5} \cdot \frac{5}{8} \\&= 0.475\end{aligned}$$

47.5 % of cases are they likely to contradict each on an identical point.

## Problem - 5

Consider the following information about travellers on vacation (based partly on a recent Travelocity poll): **40%** check work email, **30%** use a cell phone to stay connected to work, **25%** bring a laptop with them, **23%** both check work email and use a cell phone to stay connected, and **51%** neither check work email nor use a cell phone to stay connected nor bring a laptop. In addition, **88** out of every **100** who bring a laptop also check work email, and **70** out of every **100** who use a cell phone to stay connected also bring a laptop.

A<sub>E</sub>  
B<sub>C</sub>  
C<sub>L</sub>

- a. What is the probability that a randomly selected traveller who checks work email also uses a cell phone to stay connected?
- b. What is the probability that someone who brings a laptop on vacation also uses a cell phone to stay connected?

$$\rightarrow P(A_C | A_L)$$

# Solution - 5

Let  $A_E$  = {Check work mail} ✓

Let  $A_c$  = {use a cell phone to stay connected to work}

Let  $A_l$  = {bring a laptop} ✓

We are given the following probabilities

$$\underline{P(A_E)} = \underline{0.4}$$

$$\underline{P(A_c)} = \underline{0.3} \text{ and } \underline{P(A_l)} = \underline{0.25}$$

$$\underline{P(A_E \cap A_c)} = \underline{0.23}$$

$$\underline{P(A'_E \cap A'_c \cap A'_l)} = \underline{0.51}$$

$$P\left(\frac{A_E}{A_l}\right) = \underline{0.88} = \underline{\underline{88\%}}_{100}$$

$$P\left(\frac{A_l}{A_c}\right) = \underline{0.7}$$

$$\begin{aligned}
 & A_C \wedge A_L \\
 \text{(ii)} \quad P(A_C/A_L) &= \frac{P(A_C \cap A_L)}{P(A_L)} \\
 &= P(A_L/A_C) \cdot \frac{P(A_C)}{P(A_C)}
 \end{aligned}$$

$$\begin{aligned}
 P(A \cap B) &= P(B) \cdot P(A|B) \\
 &= P(A) \cdot P(B|A)
 \end{aligned}$$

a. We need to find the probability that a randomly selected traveller who checks work email also uses a cell phone to stay connected.

Which is a conditional probability of  $A_c$  given  $A_E$  i.e,

$\underline{A_c / A_E}$

$$P(\underline{A_c / A_E}) = \frac{\underline{P(A_c \cap A_E)}}{\underline{A_E}}$$
$$= \frac{0.23}{0.4} = \underline{0.575}$$

b. The probability that someone who brings a laptop on vacation also uses a cell phone to stay connected.

Which is a conditional probability of  $A_c$  given  $A_l$  i.e,

$$\begin{aligned} P(A_c/A_l) &= \frac{P(A_c \cap A_l)}{P(A_l)} = \frac{P(A_l/A_c) \cdot P(A_c)}{P(A_l)} \\ &= \frac{(0.7)*(0.3)}{0.25} = 0.84 \end{aligned}$$

# Problem - 6

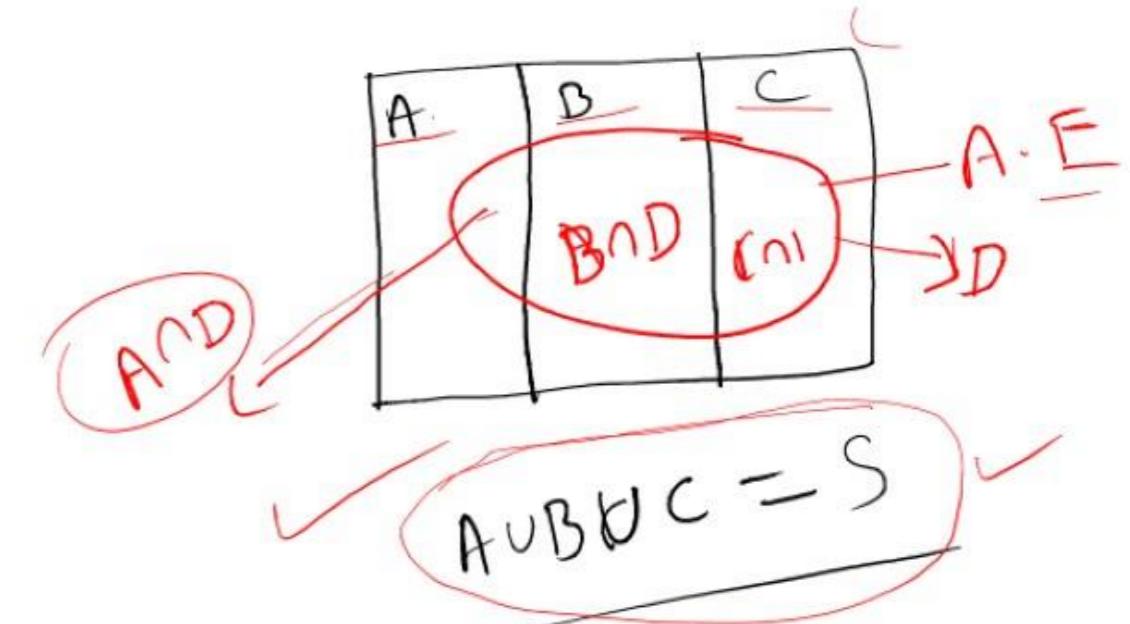
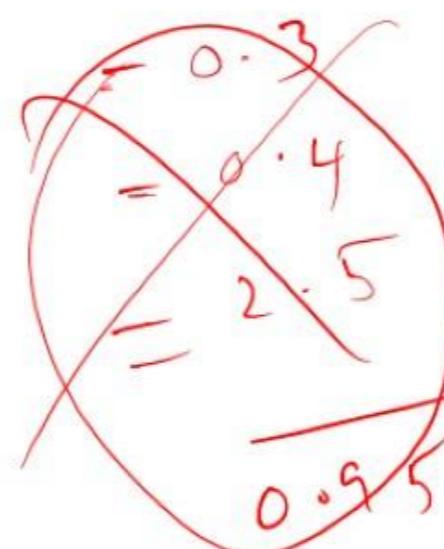
At a certain gas station, 40% of the customers use regular gas ( $A_1$ ), 35% use plus gas ( $A_2$ ), and 25% use premium ( $A_3$ ). Of those customers using regular gas, only 30% fill their tanks (event  $B$ ). Of those customers using plus, 60% fill their tanks, whereas of those using premium, 50% fill their tanks.

- What is the probability that the next customer will request plus gas and fill the tank?
- What is the probability that the next customer fills the tank?  $\rightarrow P(B)$
- If the next customer fills the tank, what is the probability that regular gas is requested? Plus? Premium?

$$P(A_1/B) \quad P(A_2/B) \quad P(A_3/B)$$

$$P(B/A_1) = 30\% \quad P(B/A_2) = 60\% \\ P(B/A_3) = 50\%$$

A



$$P(A_1) = 40\% = 0.4$$

$$P(A_2) = 35\% = 0.35$$

$$P(A_3) = 25\% = 0.25$$

$$= 1.0$$

# Solution - 6

Probabilities of customers using regular gas

$$P(A_1) = 40\% = 0.4$$

Probabilities of customers using plus gas

$$P(A_2) = 35\% = 0.35$$

Probabilities of customers using premium gas

$$P(A_3) = 25\% = 0.25$$

Also given with conditional probabilities of full gas tank

$$P(B/A_1) = 30\% = 0.3$$

$$P(B/A_2) = 60\% = 0.6$$

$$P(B/A_3) = 50\% = 0.5$$

$$P(B) = P(A_1 \cap B) + P(A_2 \cap B) + P(A_3 \cap B)$$

(a) The probability that next customer will require plus gas and fill the

$$P(A_2 \cap B) = P(A_2) * P(B/A_2)$$

$$= 0.35 * 0.6$$

$$= 0.21$$

$\checkmark A_2, B$

$$P(A \cup B) = P(A) + P(B)$$

$A, B \rightarrow \text{mutually F}$

$$P(B) = P(\bar{A}_1 \cap B) + P(A_2 \cap B) + P(\bar{A}_3 \cap B)$$

(b) The probability of next customer filling the tank is

$$\begin{aligned} P(B) &= P(A_1) P\left(B/A_1\right) + P(A_2) P\left(B/A_2\right) + P(A_3) P\left(B/A_3\right) \\ &= (0.4 \times 0.3) + (0.35 \times 0.6) + (0.25 \times 0.5) \\ &= 0.455 \end{aligned}$$

(c) If the next customer fills the tank, probability of requesting regular gas is

$$\begin{aligned} P(A_1/B) &= \frac{P(A_1) P\left(B/A_1\right)}{P(B)} \\ &= \frac{0.4 \times 0.3}{0.455} = 0.264 \end{aligned}$$

$\bar{A}_1$

$P(\bar{A}_1 \cap B)$

$P(B)$

If the next customer fills the tank, probability of requesting plus gas is

$$P(A_2/B) = \frac{P(A_2) P(B/A_2)}{P(B)}$$

$\frac{0.35*0.6}{0.455} = 0.462$

If the next customer fills the tank, probability of requesting premium gas is

$$P(A_3/B) = \frac{P(A_3) P(B/A_3)}{P(B)}$$

$\frac{0.25*0.5}{0.455} = 0.274$

## Problem - 7

“A company wants to find reasons for the dissatisfaction among employees because of which they are leaving the company. With this view they took the feedback among the employees and identified three reasons i.e. working conditions, pay hike, commuting. The probabilities of dissatisfaction with these three factors are 0.6, 0.3 and 0.1 respectively. And the probabilities that they are leaving the organization with these reasons are 0.3, 0.5 and 0.8 respectively.”

As a data scientist, use an appropriate statistical model / method to model this case and suggest the company the probable reasons on priority so that they can focus on it to retain the employees.

# Solution - 7

Let **W.C** = with the working conditions employees are dissatisfaction.

$$P(W.C) = 0.6$$

Let **P.H** = with the pay hike employees are dissatisfaction.

$$P(P.H) = 0.3$$

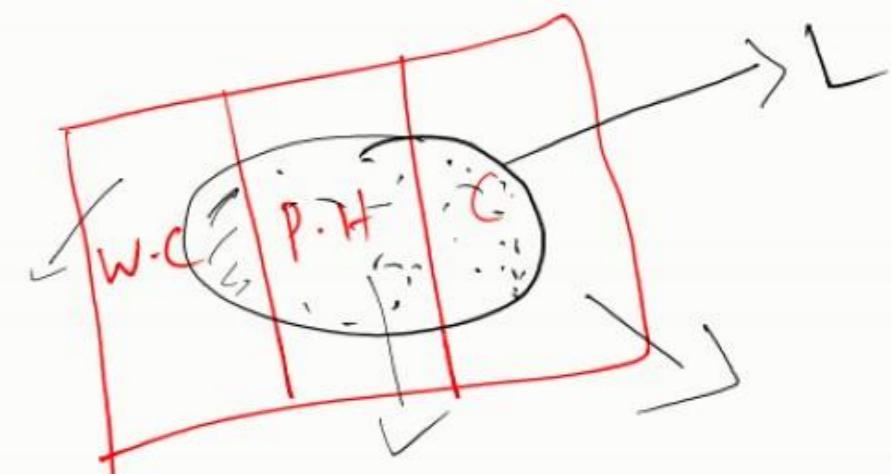
Let **CM** = with the commuting employees are dissatisfaction.

$$P(CM) = 0.1$$

$$\frac{1}{10}$$

And

Let **L** be event that they are leaving the company.



## Solution - 7

$$P(L/W.C) = 0.3, \quad P(L/P.H) = 0.5, \quad P(L/C.M) = 0.8$$

P[dissatisfaction among employees because of W.C and they are leaving the company]

$$P[(W.C) \cap L] = P(W.C)P(L/W.C) = (0.6)(0.3) = 0.18$$

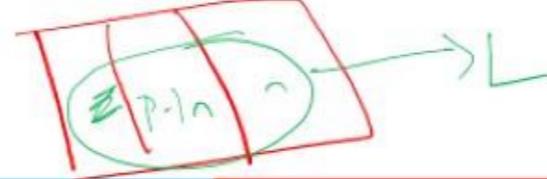
P[dissatisfaction among employees because of P.H and they are leaving the company]

$$P[(P.H) \cap L] = P(P.H)P(L/P.H) = (0.3)(0.5) = 0.15$$

P[dissatisfaction among employees because of CM and they are leaving the company]

$$P[(C.M) \cap L] = P(C.M)P(L/C.M) = (0.1)(0.8) = 0.08$$

# Solution - 7



$$\begin{aligned} P[\text{leaving the company}] &= P[L] = \underbrace{P[(W.C) \cap L] + P[(P.H) \cap L] + P[(CM) \cap L]}_{\text{Total probability}} \\ &= P(W.C)P(L/W.C) + P(P.H)P(L/P.H) + P(CM)P(L/CM) \\ &= (0.6)(0.3) + (0.3)(0.5) + (0.1)(0.8) \\ &= 0.41 \end{aligned}$$

$P[\text{they are leaving the organization, because of W.C}]$

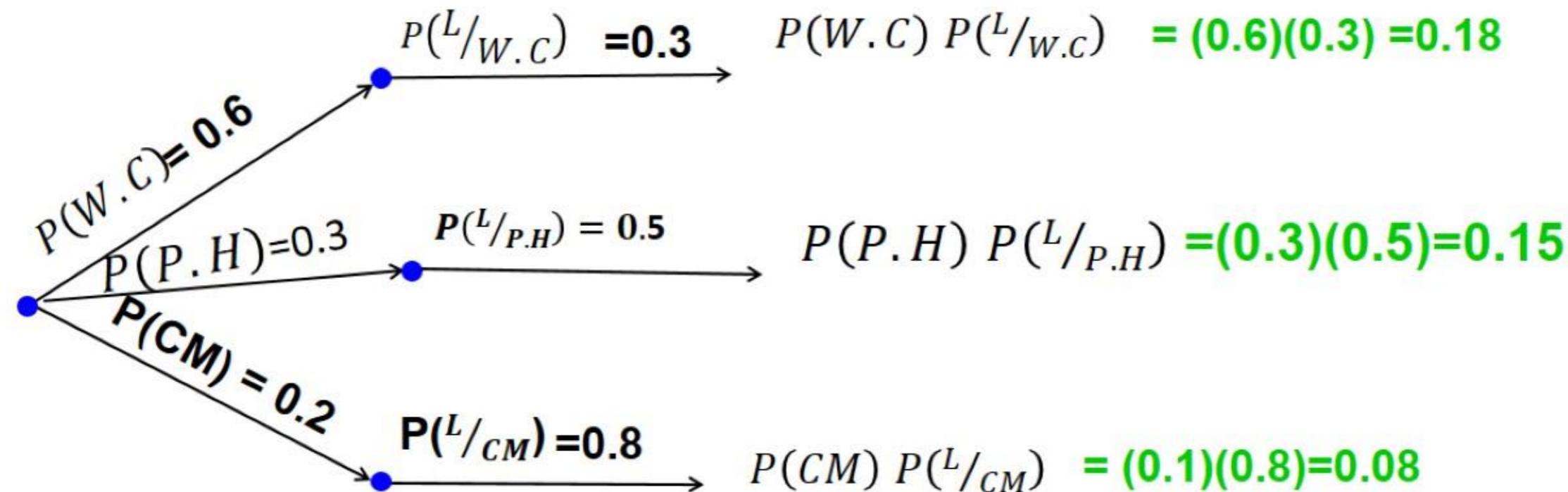
$$= P(W.C/L) = \frac{P[(W.C) \cap L]}{P[L]} = \frac{0.18}{0.41} = 0.4390$$

$$P(P.H/L) = \frac{P[(P.H) \cap L]}{P[L]} = \frac{0.15}{0.41} = 0.3658$$

$$P(CM/L) = \frac{P[(CM) \cap L]}{P[L]} = \frac{0.08}{0.41} = 0.1951$$

Most of the employees are leaving the organization with the reason working conditions.

# Probability Tree Diagram



---

$$\mathbf{P(L)} = P(W.C) P(L/W.C) + P(P.H) P(L/P.H) + P(CM) P(L/CM) = 0.41$$

---

## Problem - 8



A company that manufactures video camera as produces a basic model and a deluxe model. Over the years 40% of the cameras sold have been the basic model. Of those buying the basic model 30% purchase an extended warranty, whereas 50% of all deluxe purchasers do so. What is the probability that a randomly selected customer has an extended warranty, how likely is it that he or she has a basic model?

## Solution - 8



Let B be the event of a basic model.

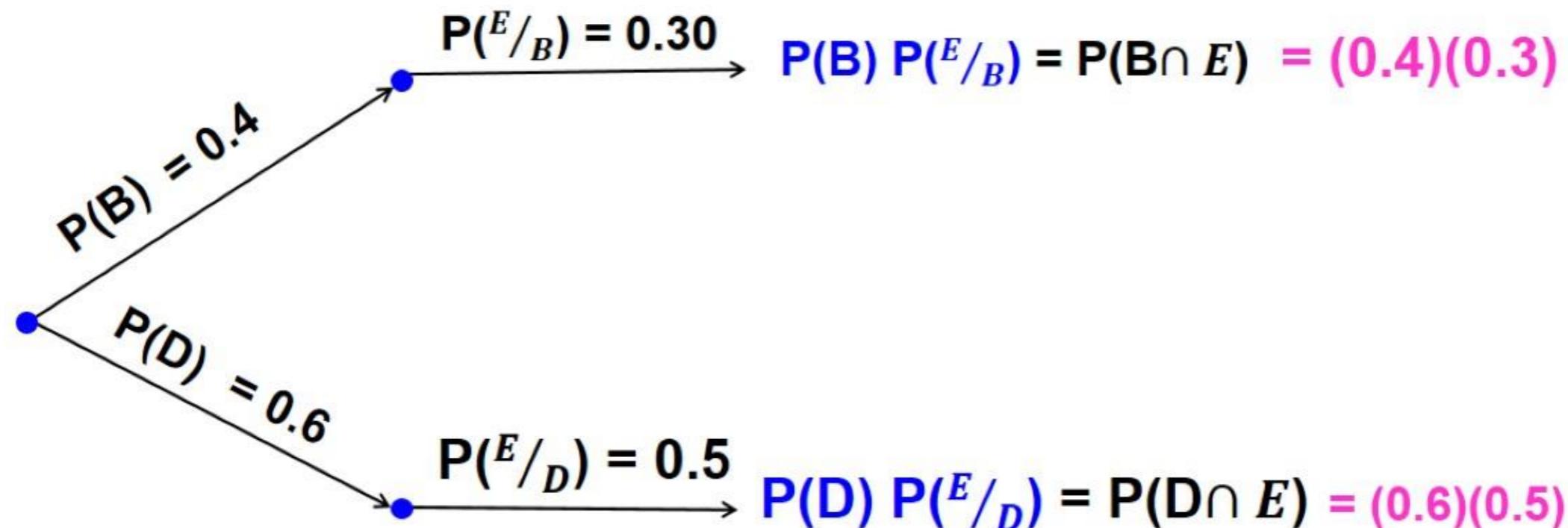
Let D be the event of a deluxe model.

Let E be event of extended warranty.

$$P(B) = 0.40, P(D) = 0.60,$$

$$P(E / B) = 0.30, \text{ and } P(E / D) = 0.50$$

# Probability Tree Diagram



---

$$P(E) = P(B) P(E|B) + P(D) P(E|D) = 0.42$$

---

## Solution - 8



the probability that a randomly selected customer has an extended warranty,  
how likely is it that he or she has a basic model is

$$P(B/E) = \frac{P(B \cap E)}{P(E)} = \frac{(0.4)(0.3)}{0.42} = 0.2857$$

## Problem - 9



The chance that doctor A will diagnose a disease X correctly is 60%. The chance that a patient will die by his treatment after correct diagnose is 40% and the chance of death by wrong diagnose is 70%. A patient of doctor A, who had disease X, died. What is the chance that his disease was correctly.

## Solution-9



Let  $E_1$  = event that disease X is diagnosed correctly by doctor A

Let  $E_2$  = event that a patient of doctor A who has disease X died

Then  $P(E_1) = \frac{60}{100} = 0.6$  and  $P(\bar{E}_1) = 1 - P(E_1) = 0.4$

$P(E_2/E_1) = \frac{40}{100} = 0.4$  and  $P(\bar{E}_2/\bar{E}_1) = \frac{70}{100} = 0.7$

## By Baye's theorem

$$P(E_1/E_2) = \frac{P(E_1)P(E_2/E_1)}{P(E_1)P(E_2/E_1) + P(\bar{E}_1)P(E_2/\bar{E}_1)}$$
$$= \frac{0.6 * 0.4}{0.6 * 0.4 + 0.4 * 0.7} = \frac{6}{13}$$

# Probability Mass function(Pmf)



## (Probability distribution function)

Let 'X' be a discrete random variable taking values  $x_1, x_2, \dots, x_n$  then the probability mass function  $P(X = x)$  is defined under the following conditions

- (i)  $P(X = x) \geq 0$
- (ii)  $\sum_{\forall x} P(X = x) = 1$

# Continuous Probability Distributions



A random variable is called continuous when it assumes values in a given interval.

The probability that a random variable  $X$  assumes different values  $x$  in a given interval, say  $[a, b]$ , is denoted by  $f(x) = P(a \leq X \leq b)$ , called **probability density function (pdf)**.

And

$$f(x) = P(a \leq X \leq b) = \int_a^b f(x) dx$$

# Properties of continuous probability distribution



A function  $f(x)$  to be a probability density function (pdf), if it satisfies the following conditions

1.  $f(x) \geq 0$
2.  $\int_{-\infty}^{\infty} f(x)dx = 1$
3. For any  $a \leq b$ ,

$$P(a \leq X \leq b) = \int_a^b f(x)dx$$

Let  $X$  be a continuous random variable.

Then 
$$\begin{aligned} P(a \leq X \leq b) &= P(a < X \leq b) \\ &= P(a \leq X < b) \\ &= P(a < X < b) \end{aligned}$$

**NOTE:** The Cumulative Distribution Function of continuous random variable  $X$  is denoted by  $F(X)$  and is defined as

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(x)dx$$

and  $f(x) = \frac{d}{dx}[F(x)]$

## Problem - 10



Verify that  $P(X) = \frac{x+3}{25}$  for  $x = 1,2,3,4,5$

Serve as Probability mass function?

# Solution:

Given that  $P(X) = \frac{x+3}{25}$  for  $x = 1, 2, 3, 4, 5$

$$P(X = 1) = \frac{1+3}{25} = \frac{4}{25}, \quad \checkmark$$

$$P(X = 2) = \frac{5}{25}, \quad \checkmark$$

$$P(X = 3) = \frac{6}{25}, \quad \checkmark$$

$$P(X = 4) = \frac{7}{25}, \quad \checkmark$$

$$P(X = 5) = \frac{8}{25}, \quad \checkmark$$

$$\begin{aligned} & P(A) \leq 1 \\ & (i) 0 \leq P(A) \leq 1 \quad \checkmark \end{aligned}$$

$$\begin{aligned}\sum_{x=1}^5 P(X = x) &= P(X = 1) + P(X = 2) + P(X = 3) + P(X = 4) \\&\quad + P(X = 5) \\&= \frac{4}{25} + \frac{5}{25} + \frac{6}{25} + \frac{7}{25} + \frac{8}{25} \\&= \frac{30}{25} > 1\end{aligned}$$

Given  $P(X)$  is not a Probability mass function  
[ we know that Total probability is 1 ]

# Problem - 11

For the discrete probability distribution

X=x	0	1	2	3	4	5	6	7
P(X=x)	0	k	2k	2k	3k	$k^2$	$2k^2$	$7k^2 + k$

- find (i) k (ii)  $P(X < 6)$  (iii)  $P(X \geq 6)$   
(iv) Mean and Variance

# Solution - 11

(i) We know that the Total probability is 1

$$\sum_{i=0}^7 P(X = x_i) = 1$$

$$0 + k + 2k + 2k + 3k + k^2 + 2k^2 + 7k^2 + k = 1$$

$$10k^2 + 9k - 1 = 0$$

$$k = -1$$

$$k = \frac{1}{10}$$

then  $\underline{k = \frac{1}{10}}$  [ since  $P(X) \geq 0$ , so  $k \neq -1$ ]

# Solution - 11



$$\begin{aligned}(ii) P(X < 6) &= \cancel{P(X = 0)} + \cancel{P(X = 1)} + \cancel{P(X = 2)} + \cancel{P(X = 3)} \\&\quad + \cancel{P(X = 4)} + \cancel{P(X = 5)} \\&= 0 + k + 2k + 2k + 3k + k^2 \\&= k^2 + 8k \\&= 0.81 [\cancel{k} = \frac{1}{10} = 0.1 ]\end{aligned}$$

$$(iii) \cancel{P(X \geq 6)} = \cancel{P(X = 6)} + \cancel{P(X = 7)}$$

or

$$\begin{aligned}P(X \geq 6) &= 1 - P(X < 6) = 1 - 0.81 \\&= 0.19\end{aligned}$$

# Solution - 11

$$\begin{aligned}\text{Mean} &= E[x] = \sum_{x=0}^7 xP(X=x) \\ &= 0 + k + 4k + 6k + 12k + 5k^2 + 12k^2 + 49k^2 + 7k \\ &= 66k^2 + 30k = 3.66\end{aligned}$$

$$var(x) = E(x^2) - [E(x)]^2$$

$$\begin{aligned}E(x^2) &= \sum_{x=0}^7 x^2 P(X=x) \\ &= 0 + k + 8k + 18k + 48k + 25k^2 + 72k^2 + 343k^2 + 49k \\ &= 440k^2 + 75k = 16.8\end{aligned}$$

$$\begin{aligned}\text{Then } var(x) &= E(x^2) - [E(x)]^2 = 16.8 - (3.66)^2 \\ &= 3.4044\end{aligned}$$

X=x	0	1	2	3	4	5	6	7
P(X=x)	0	k	2k	2k	3k	$k^2$	$2k^2$	$7k^2 + k$

# Thanks