



शिक्षा मंत्रालय
MINISTRY OF
EDUCATION

INDIAN INSTITUTE OF TECHNOLOGY
JODHPUR



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING



P M R F

Prime Minister's Research Fellowship

Week 4 - Live Session

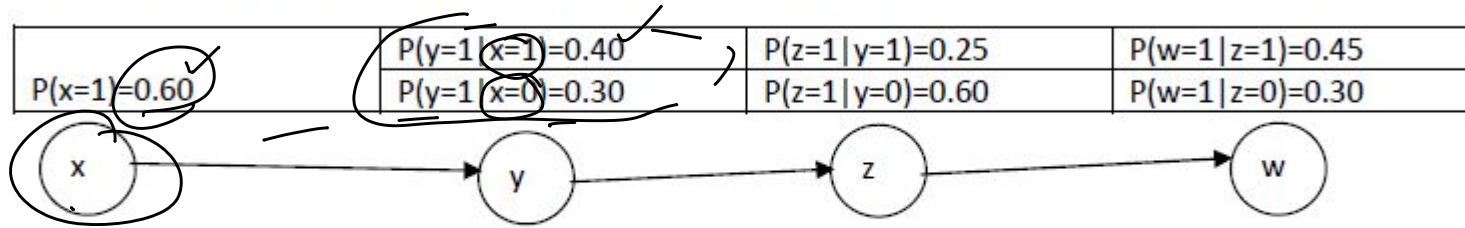
Data Mining

Swapnil S. Mane

mane.1@iitj.ac.in

PMRF Research Scholar

Q1-3 are based on a simple Bayesian Network shown below:



The Bayesian Network is fully specified by the marginal probabilities of the root node(x) and the conditional probabilities.

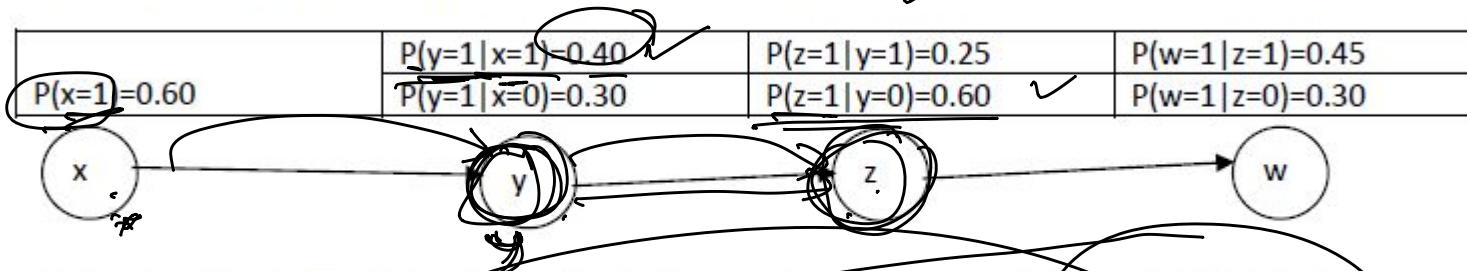
ca

Q1. $P(y=0)$ is:

- a) 0.70
- b) 0.12
- ✓ c) 0.64
- d) 0.36

$$\begin{aligned}
 P(y=0) &= 1 - P(y=1) \\
 P(y=1) &= P(y=1|x=1) \cdot P(x=1) + P(y=1|x=0) \cdot P(x=0) \\
 &= (0.40 \cdot 0.60) + (0.30 \cdot 0.40) \\
 P(y=1) &= 0.36 \\
 P(y=0) &= 1 - 0.36 \\
 &= 0.64
 \end{aligned}$$

Q1-3 are based on a simple Bayesian Network shown below:



The Bayesian Network is fully specified by the marginal probabilities of the root node(x) and the conditional probabilities.

$$p(y=1|x=1) + p(y=0|x=1) = 1$$

Q2. $P(z=1|x=1)$ is:

- a) 0.50
- b) 0.60
- ✓ c) 0.46
- d) 0

$$\begin{aligned}
 P(z=1|x=1) &= P(z=1|y=1) * P(y=1|x=1) \\
 &\quad + P(z=1|y=0) * P(y=0|x=1) \\
 &= 0.46
 \end{aligned}$$

Q1-3 are based on a simple Bayesian Network shown below:

$P(x=1)=0.60$	$P(y=1 x=1)=0.40$	$P(z=1 y=1)=0.25$	$P(w=1 z=1)=0.45$
	$P(y=1 x=0)=0.30$	$P(z=1 y=0)=0.60$	$P(w=1 z=0)=0.30$



The Bayesian Network is fully specified by the marginal probabilities of the root node(x) and the conditional probabilities.

$$P(z=1|x=1) = 0.46$$

Q3. $P(w=0|x=1)$ is:

- a) 0.37
- b) 0.63
- c) 1
- d) None of the above

$$\begin{aligned}
 P(w=0|x=1) &= P(w=0|z=1) \cdot P(z=1|x=1) \\
 &\quad + P(w=0|z=0) \cdot P(z=0|x=1) \\
 &= (0.55 \cdot 0.46) + (0.70 \cdot 0.54) \\
 &= 0.63
 \end{aligned}$$

Q4. Consider a binary classification problem with two classes C1 and C2. Class labels of ten other training set instances sorted in increasing order of their distance to an instance x is as follows: ~~{C1, C2, C1, C2, C2, C2, C1, C2, C1, C2}~~. How will a $K=5$ nearest neighbor classifier classify x ?

a) There will be a tie

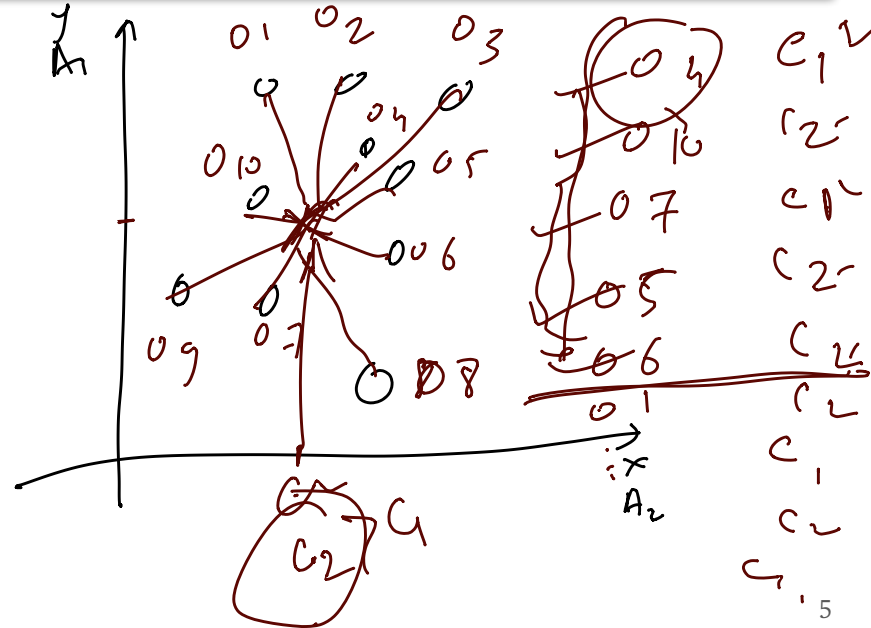
b) C1

☒ c) C2

d) Not enough information to classify

training set

	A_1	A_2	C
o_1			C_1
o_2			C_2
o_3			C_2



Consider the following data for questions 5-6.

You are given the following set of training examples. Each attribute can take value either 0 or 1.

A1	A2	A3	Class
0	0	1	C1
0	1	0	C1
0	1	1	C1
1	0	0	C2
1	1	0	C1
1	1	1	C2

Q5. How would a 3-NN classify the example (A1 = 1, A2 = 0, A3 = 1) if the distance metric is Euclidean distance?

a) C1

b) C2

c) There will be a tie

d) Not enough information to classify

$$k=3$$

$$(1) \text{dis}(x, o_1) = \sqrt{1+0+0} = 1$$

$$(5) \text{dis}(x, o_2) = \sqrt{1+1+1} = 1.73$$

$$\text{dis}(x, o_3) = \sqrt{1+1+0} = 1.41$$

$$\text{dis}(o_1, x) = \sqrt{\sum_{i=1}^3 (x_i - o_{1i})^2}$$

$$\text{dis}(x, o_4) = \sqrt{0+0+1} = 1$$

$$\text{dis}(x, o_5) = \sqrt{0+1+1} = 1.41$$

$$\text{dis}(x, o_6) = \sqrt{0+1+0} = 1$$

Consider the following data for questions 5-6.

You are given the following set of training examples. Each attribute can take value either 0 or 1.

A1	A2	A3	Class
0	0	1	C1
0	1	0	C1
0	1	1	C1
1	0	0	C2
1	1	0	C1
1	1	1	C2

Q6. How would a 3-NN classify the example (A1 = 0, A2 = 0, A3 = 0) if the distance metric is Euclidean distance?

a) C1

b) C2

c) There will be a tie

d) Not enough information to classify

$$dis(x, o_1) = \sqrt{0 + 0 + 1} = 1$$

$$dis(x, o_2) = \sqrt{0 + 1 + 0} = 1$$

$$dis(x, o_3) = \sqrt{0 + 1 + 1} = 1.4$$

$$dis(x, o_4) = \sqrt{1 + 0 + 0} = 1$$

$$dis(x, o_5) = \sqrt{1 + 1 + 0} = 1.4$$

$$dis(x, o_6) = \sqrt{1 + 1 + 1} = 1.73$$

$$A(1, 1, 100)$$

$$B(1, 2, 100)$$

$$C(9, 9, 2)$$

$$AB = \sqrt{0 + 1 + 0} = 1 \quad \checkmark$$

$$AC = \sqrt{8^2 + 8^2 + 998} \approx 999$$

Q7. Issues with Euclidean measure are:

- a) High dimensional data.
- b) Can produce counter-intuitive results.
- c) Shrinking density – sparsification effect
- d) All of the above.

$$\begin{aligned} A^1 A^2 \\ A(2,3) \\ B(4,1) &= \sqrt{(2-4)^2 + (3-1)^2} \\ &= \sqrt{4+4} = \sqrt{8} \approx 2.8 \end{aligned}$$

$$\begin{aligned} A(2,3,0.01,0.01) \\ B(4.5,0.01,0.01) &= \sqrt{(2-4.5)^2 + (3-0.01)^2 + (0.01-0.01)^2 + (0.01-0.01)^2} \\ &= \sqrt{6.25 + 8.9801} = \sqrt{15.2301} \approx 3.9 \end{aligned}$$

$$\begin{aligned} A(1,1) \\ B(1,2) \\ C(9,9) \\ AB &= \sqrt{0+1} = 1 \\ BC &= \sqrt{(8)^2 + (-7)^2} = 11.3 \end{aligned}$$

A(1,1)
B(1,2)
C(9,9)