

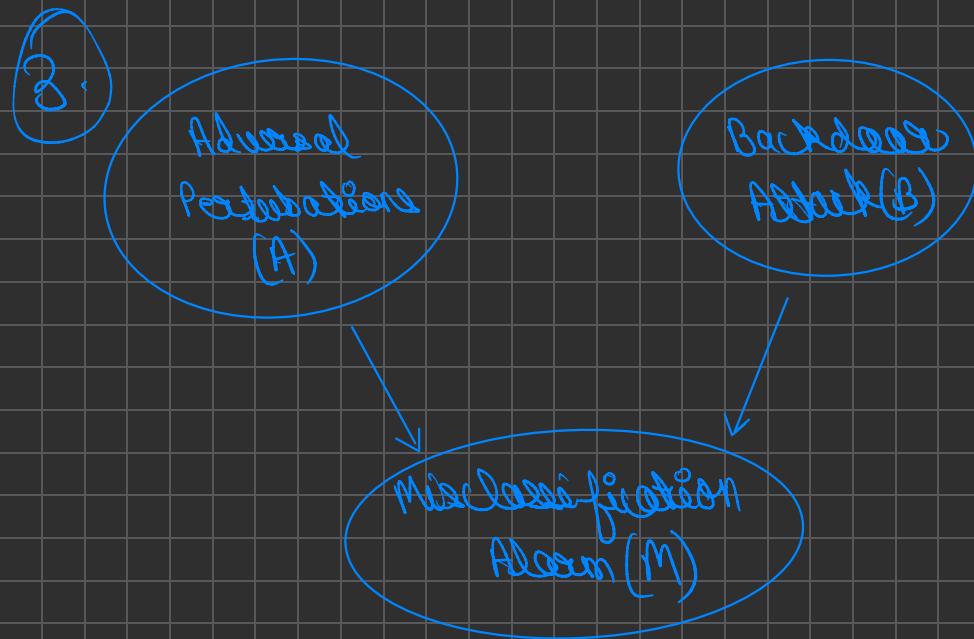
$$= \frac{0.1312}{0.416} = 0.32$$

$$P(J \wedge C | R) = P(J | R) P(C | R)$$

$$P(J | R) \times P(C | R) = 0.6659 \times 0.32 \\ = 0.213$$

$$P(J \wedge C | R) = 0.213 = 0.213$$

J and C are conditionally independent given R



Initial Observation , A and B are independent

a)

$$\text{Posterior Probability} = \frac{P(M|A) \cdot P(A)}{P(M)}$$

$$P(M) = P(M|A, B) \cdot P(A \cap B) + P(M|A, \neg B) \\ P(A \cap \neg B) \\ + P(M|\neg A, B) \cdot P(\neg A \cap B)$$

b) Prior Probabilities

$P(A), P(B) \rightarrow \text{independent}$

$P(A)$: initial prior belief of adversarial perturbation

$P(B)$: initial prior belief of backdoor attack

Likelihood Probabilities are:

$P(M|A)$: Prob. of observing misclassification given that adversarial perturbation has occurred

$P(M|B)$: Prob. of observing misclassification given that backdoor attack has occurred

Posterior Probability

$P(A|M)$: Post. belief OR likelihood of adversarial perturbation given that a missclassification has occurred

$P(B|M)$: Updated belief in backdoor attacks after observing M

C.

(i) Initial Independence: Before the report Adversarial perturbations (A) and backdoor attacks (B) were independent causes of M.

Observing M increased the likelihood of both A and B

(ii) Impact Of New Information: The report increases $P(B)$, the prior probability of backdoor attacks, making B a more likely explanation of M

(ii) Explaining Away Effect: As B becomes more probable, it explains away M , reducing the posterior prob. of A being the cause:

$$P(A|M, \text{high } P(B)) < P(A|M, \text{original } P(B))$$

Conclusion: The increased likelihood of background attacks shifts the explanation for M toward B , decreasing the belief that A caused the misclassification. This is a result of the explaining-away effect.