

Q.3

- (a) A : adversal perturbations
B : Backdoor attacks
m : Misclassification alarm

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

→ likelihood of misclassification given adversal perturbation occurred
→ Prior probability of adversal perturbation
→ Total probability of misclassification

here $P(m) = P(m|A) \cdot P(A) + P(m|B) \cdot P(B)$

(b) Prior Probabilities are:

- $P(A)$: initial / Prior belief of adversal perturbations
- $P(B)$: initial / Prior belief of backdoor attacks

Likelihood Probabilities are:

- $P(m|A)$: Probability of observing misclassification given that adversal perturbation has occurred

- $P(m|B)$: Probability of observing misclassification given that backdoor attack has occurred

Posterior Probability:

- $P(A|m)$: Probability / belief of likelihood of adversal perturbation given that a misclassification has occurred

- (c) • it is given that there is an increase in backdoor attacks, which results in updation of prior $P(B)$ to increase

- since $P(m)$ is directly positively related to $P(B)$.

using part (a) : $P(M) = P(M|A) \cdot P(A) + P(M|B) \cdot P(B)$

increase in probability of backdoor attack will result in an increase in probability of a misclassification alarm

$$\bullet P(A|M) = \frac{P(M|A) \cdot P(A)}{P(M|A) \cdot P(A) + P(M|B) \cdot P(B)}$$

An increase in $P(B)$ will result in reduction of our belief regarding the likelihood of adversarial perturbation causing the misclassification due to inverse positive relationship.