

Conditional Probability:

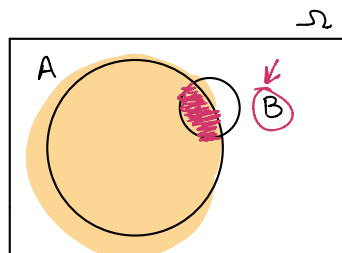
$$\underline{P(A|B)} = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B) \neq P(B|A)$$

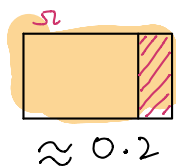
Interpretation

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

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

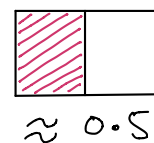


$$P(B|A)$$



≈ 0.2

$$P(A|B)$$



≈ 0.5

Problem:

- A disease is prevalent in 0.2% of a population.
- We have a test that, given to a sick person, gives a +ve result 85% of the time.
- Of all the people ever tested, 8% were positive.

Q: If Nazo is tested and test comes back positive, what are the chances that she actually has the disease?

☐ 85%

☐ 77%

☐ 21%

☒ 2%

$$P(\text{Disease}) = 0.002 \quad 0.2\%$$

(*) Event "Disease" is difficult to measure directly!

$$P(\text{Pos} | \text{Disease}) = 0.85$$

Event "Pos" is relatively easier to measure

$$P(\text{Pos}) = 0.08$$

$$P(\text{Disease} | \text{Pos}) = ?$$

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

$$P(B|A) = \frac{P(A|B) P(B)}{P(A)}$$

$$P(\overset{B}{\text{Disease}} | \overset{A}{\text{Pos}}) = \frac{P(\overset{A}{\text{Pos}} | \overset{B}{\text{Disease}}) P(\overset{B}{\text{Disease}})}{P(\overset{A}{\text{Pos}})}$$

$$= \frac{(0.85)(0.002)}{0.08}$$

$$= 0.021$$

$$\underline{\underline{2.1\%}}$$

$$\begin{array}{c}
 \text{Posterior} \\
 \hline
 P(\text{Disease} \mid \text{Pos}) = \frac{\overbrace{P(\text{Pos} \mid \text{Disease})}^{\text{Likelihood}} \overbrace{P(\text{Disease})}^{\text{Prior}}}{\underbrace{P(\text{Pos})}_{\text{normalizing factor}}}
 \end{array}
 \quad \leftarrow \begin{array}{c} \text{before} \end{array}$$

\uparrow
 after

$$P(\text{Disease}) = 0.002$$

— Prior belief

$$P(\text{Pos} \mid \text{Disease}) = 0.85$$

— Result of experiment

$$P(\text{Pos}) = 0.08$$

$$P(\text{Disease} \mid \text{Pos}) = 0.21$$

— Updated belief

⊗ You start off with some belief and update it based on some experiment!

This is the "Bayes' Rule" of inference.

$$P(B \mid A) = \frac{P(A \mid B) P(B)}{P(A)}$$

Classical Statistics \longrightarrow

\longrightarrow Bayesian " \longrightarrow

Applying Bayes' Rule to Spam Detection.

- "You have inherited a million dollars." ← Spam
- "There will be a meeting at noon." ← not-spam
- Assumption: We have a dataset of spam emails. ← test dataset
- Need to find whether a piece of text is spam.
- Let's first consider a single word.

$P(\text{Spam} | w) = \frac{P(w | \text{spam}) P(\text{spam})}{P(w)}$

"how frequently does this word appear in spam?" "How much spam is there in the world?"

"Given that this word appears, how likely is it that the message is spam?" from dataset. "How frequent is this word?"

$P(\text{spam}) = \frac{\text{\# of spam messages}}{\text{\# of all messages}} = \frac{100}{500}$

$P(w | \text{spam}) = \frac{\text{\# of times this word appears in spam}}{\text{\# of spam messages}}$

$P(w) = \frac{\text{\# of times this word appears}}{\text{\# of total messages}}$

- Now, do this for all words

- $P(\text{spam} | \text{words}) = P(\text{spam} | w_1) * P(\text{spam} | w_2) * \dots * P(\text{spam} | w_n)$

"independence"

$$P(\text{spam} | \text{words}) = \prod_{i=1}^{|\text{words}|} P(\text{spam} | w_i)$$

→ words → independent
"Naïve Bayes Model"

$$\prod_{i=1}^n$$

$$\sum_{i=1}^n$$