BAYES' THEOREM

# BAYES THEOREM

- one of the most famous equations in statistics and probability

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

# A FRAMEWORK FOR UPDATING OUR BELIEFS

- What's the point of probability ? **=> decision making under uncertainty**

- Our knowledge about the world is never totally exact, so how do we decide whether to go ahead with a decision or not ?

- **Bayes' Theorem gives us a quantitative framework for updating our beliefs as the facts around us change…**

# THE INTUITION WITH AN EXAMPLE

- It's 9AM on Monday morning, and you receive an email from your boss. You notice that it seems a little different from her usual notes: the message contains several grammatical errors, and ends by asking you to provide your social security number. Though you first assumed it was a legitimate email, the grammar mistakes and suspicious request convince you to send it right to the spam folder.

(https://medium.com/opex-analytics/bayes-theorem-101-6a9a1ea5d4a6)

- When making that quick decision to ignore the email from your "boss," you unconsciously estimated several different probabilities.
  - First, you judged the likelihood of a work email's legitimacy to be fairly high.
  - But then you assessed the probability that such a weird email could come from your boss to be low. You also have some general sense that phishing emails tend to be weird in a few specific ways, and you know that phishing scams are common enough that this particular email could plausibly be harmful.

- With all this information swirling around in your head, you decide that the email is most likely spam.

- That's pretty much all bayes theorem is: **updating** our prior beliefs given some particular piece of information.

# TAKING A CLOSER LOOK AT THE FORMULA

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

- **P(A|B)** – is the probability of A given that B has already happened.
- **P(B|A)** – is the probability of B given that A has already happened. It looks circular and arbitrary for now…
- **P(A)** – is the unconditional probability of A occurring.
- **P(B)** – is the unconditional probability of B occurring.

# TAKING A CLOSER LOOK AT THE FORMULA

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

P(A|B) is a conditional probability – one that measures probability over only certain states of the world (states where B has occurred).

P(A) is an example of an unconditional probability and is measured over all states of the world.

# VIDEO EXPLANATION

- https://www.youtube.com/watch?v=HZGCoVF3YvM&t=542s

# PRACTICE PROBLEM 1

**Using Bayes for NLP to predict spams based on the content of an email.**

- Assume that the word *offer* occurs in 80% of the spam messages

- Also assume *offer* occurs in 10% of desired e-mails (hams)

- If 30% of the received e-mails are considered to be spam, and I receive a new message which contains *offer*, what is the probability that this new email is a spam?

- Draw a tree diagram to help you and consider the case where you have a sample of 100 emails : you can first find the solution by counting, and then try and find it using the theorem.

# PRACTICE PROBLEM 1

- P( contains *offer*|spam) = 0.8 (given in the question)

- P(spam) = 0.3 (given in the question)

- P(contains *offer*) = 0.3*0.8 + 0.7*0.1 = 0.31

$$P(spam|contains\ offer) = \frac{P(contains\ offer|spam) * P(spam)}{P(contains\ offer)}$$

# PRACTICE PROBLEM 1

- Both results should be the same :

$$P(spam|contains\ offer) = \frac{0.8 * 0.3}{0.31} = 0.774$$

# PRACTICE PROBLEM 2

- Covid-19 tests are common nowadays, but some test results can be wrong...

- Let's assume:
    - a diagnostic test has 99% accuracy
    - and 60% of all people have Covid-19.

- If a patient tests positive, what is the probability that they actually have the disease?

- Same as previously: take a sample of 100 patients first and find the probability using counts and then use the theorem.

$$P(covid19|positive) = \frac{P(positive|covid19) * P(covid19)}{P(positive)}$$

- P(positive|covid19) = 0.99

- P(covid19) = 0.6

- P(positive) = 0.6*0.99+0.4*0.01=0.598

$$P(covid19|positive) = \frac{0.99 * 0.6}{0.598} = 0.993$$

# PRACTICE PROBLEM 3 (MONTY HALL IS BACK)

- You're on a gameshow called "**Let's Make a Deal**". There are 3 closed doors in front of you.

- Behind each door is a prize. One door has a **car**, one door has **breath mints**, and one door has a **bar of soap**. You'll get the prize behind the door you pick, but you don't know which prize is behind which door. Obviously you want the car!

- Imagine you pick **door A**.

- After picking **door A**, the host of the show, Monty Hall, now opens **door B,** revealing a bar of soap. He then asks you if you'd like to change your guess. Should you?

- By working through Bayes Theorem, we can calculate the actual odds of winning the car if we stick with **door A**, or switch to **door C**.

# PRACTICE PROBLEM 3

- The posteriors we want to compute :

  1.P(prize=A|opened=B) vs. 2.P(prize=C|opened=B)

# PRACTICE PROBLEM 3

- **Priors**
  - The probability of any door being correct before we pick a door is 1/3. Prizes are randomly arranged behind doors and we have no other information. So the **prior**, P(A), of any door being correct is **1/3**.

  1. $P(prize = A)$, the prior probability that door A contains the car = 1/3
  2. $P(prize = C)$, the prior probability that door C contains the car = 1/3

# PRACTICE PROBLEM 3

- **Likelihood**
  - If the car is behind door A, then Monty can open door B or C. So the probability of opening either is 50%.

    1. $P(opens = B | prize = A) = \frac{1}{2}$, the likelihood Monty opened door B if door A is correct

  - If the car is in fact behind door C then Monty can only open door B. He cannot open A, the door we picked. He also cannot open door C because it has the car behind it.

    2. $P(opens = B | prize = C) = 1$, the likelihood Monty opened door B if door C is correct

# PRACTICE PROBLEM 3

- **Numerator: P(A) x P(B|A)**

- $P(prize = A) \times P(opens = B|prize = A) = 1/3 \times 1/2 = 1/6$

- $P(prize = C) \times P(opens = B|prize = C) = 1/3 \times 1 = 1/3$

# PRACTICE PROBLEM 3

- **Normalize**

- This is the marginal probability P(opens=B) which is the total probability, removing dependence from any event:
  - In this case:

$$\sum P(opens = B | prize = A)P(prize = A), P(opens = B | prize = C)P(prize = C)$$

$$P(opens = B) = \frac{1}{6} + \frac{1}{3} = \frac{1}{2}$$

- Putting everything together:

1. $P(prize = A | opens = B) = \dfrac{\frac{1}{6}}{\frac{1}{2}} = \dfrac{1}{3}$

2. $P(prize = C | opens = B) = \dfrac{\frac{1}{3}}{\frac{1}{2}} = \dfrac{2}{3}$

=> the prize is more likely to be hidden behind door C, so we should switch !