

# DSC 40A

*Theoretical Foundations of Data Science I*

---

# Question

Answer at [q.dsc40a.com](http://q.dsc40a.com)

Remember, you can always ask questions at  
[q.dsc40a.com](http://q.dsc40a.com)!

If the direct link doesn't work, click the "Lecture Questions" link in the top right corner of [dsc40a.com](http://dsc40a.com).

# Agenda

- Bayes Theorem
- Naïve Bayes Classifier

MODIFIED BAYES' THEOREM:

$$P(H|X) = P(H) \times \left( 1 + P(C) \times \left( \frac{P(x|H)}{P(x)} - 1 \right) \right)$$

H: HYPOTHESIS

x: OBSERVATION

P(H): PRIOR PROBABILITY THAT H IS TRUE

P(x): PRIOR PROBABILITY OF OBSERVING x

P(C): PROBABILITY THAT YOU'RE USING  
BAYESIAN STATISTICS CORRECTLY

Source: xkcd

# Bayes Theorem

The background of the slide features abstract, overlapping green geometric shapes, primarily triangles and polygons, in various shades of green. These shapes are concentrated on the right side of the slide, with some extending towards the left. The overall effect is a modern, minimalist design.



# Last Week

- We defined Bayes' Theorem:

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

- Bayes' Theorem describes how to update the probability of one event given that another has occurred.

# Bayes' Theorem

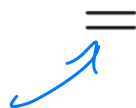
Bayes' Theorem follows from the multiplication rule, or conditional probability.

$$P(A) * P(B|A) = P(A \text{ and } B) = P(B) * P(A|B)$$

Bayes' Theorem:

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

law of total  
probability



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

not  
B



# Bayes' Theorem

For hypothesis  $H$  and evidence (data)  $E$

$$\underline{P(H | E)} = \frac{\underline{P(E|H)} \cdot P(H)}{\underline{P(E)}} = \frac{P(E|H) \cdot P(H)}{P(E|H) \cdot P(H) + \underline{P(E|\bar{H}) \cdot P(\bar{H})}}$$

- $P(H)$  - prior, initial probability before  $E$  is observed
- $P(H|E)$  - posterior, probability of  $H$  after  $E$  is observed
- $P(E|H)$  - likelihood, probability of  $E$  if the hypothesis is true
- $P(E)$  - marginal, probability of  $E$  regardless of  $H$

$$\underline{P(E|\bar{H})} \neq P(\bar{E} | H)$$

The likelihood function is a function of  $E$ , while the posterior probability is a function of  $H$ .

# Bayes' Theorem: Example

$$P(H|E) = \frac{P(E|H)P(H)}{P(E|H)P(H) + P(E|\sim H)P(\sim H)}$$

A manufacturer claims that its drug test will **detect steroid use 95% of the time**. What the company does not tell you is that 15% of all steroid-free individuals also test positive (the false positive rate). 10% of the Tour de France bike racers use steroids. Your favorite cyclist just tested positive. What's the probability that he used steroids?

What is your first guess?

- A. Close to 95%
- ~~70%~~ B. Close to 85%
- C. Close to 40%
- D. Close to 15%

# Bayes' Theorem: Example

$$P(H|E) = \frac{P(E|H)P(H)}{P(E|H)P(H) + P(E|\sim H)P(\sim H)}$$

$\bar{H}$ , not  $H$

A manufacturer claims that its drug test will **detect steroid use 95% of the time**. What the company does not tell you is that 15% of all steroid-free individuals also test positive (the false positive rate). 10% of the Tour de France bike racers use steroids. Your favorite cyclist just tested positive. What's the probability that he used steroids?

Now, calculate it and choose the best answer.

- A. Close to 95%
- B. Close to 85%
- C. Close to 40%
- D. Close to 15%

$E$  - positive drug test

$H$  - uses steroids

$\bar{H}$  - doesn't use steroids

# Bayes' Theorem: Example

$$P(H|E) = \frac{P(E|H)P(H)}{P(E|H)P(H) + P(E|\sim H)P(\sim H)} = \frac{0.95 \cdot 0.1}{0.95 \cdot 0.1 + 0.5 \cdot 0.9} \approx 0.41$$

A manufacturer claims that its drug test will **detect steroid use 95% of the time**. What the company does not tell you is that **15% of all steroid-free individuals also test positive** (the false positive rate). **10% of the Tour de France bike racers use steroids**. Your favorite cyclist just tested positive. What's the probability that he used steroids?

**Solution:**

**H: used steroids**

**E: tested positive**

$$P(E|H) = 95\% \text{ (TP)}$$

$$P(E|\bar{H}) = 15\% \text{ (FP)}$$

$$P(H) = 10\%$$

$$P(\bar{H}) = 90\%$$

$$\Rightarrow P(H|E) = 41\%$$

# Bayes' Theorem: Example

$$P(H|E) = \frac{P(E|H)P(H)}{P(E|H)P(H) + P(E|\sim H)P(\sim H)}$$

A manufacturer claims that its drug test will **detect steroid use 95% of the time**. What the company does not tell you is that 15% of all steroid-free individuals also test positive (the false positive rate). 10% of the Tour de France bike racers use steroids. Your favorite cyclist just tested positive. What's the probability that he used steroids?

**Solution:**

**H: used steroids**

**E: tested positive**

Despite manufacturer's claims, only **41% chance** that cyclist used steroids.

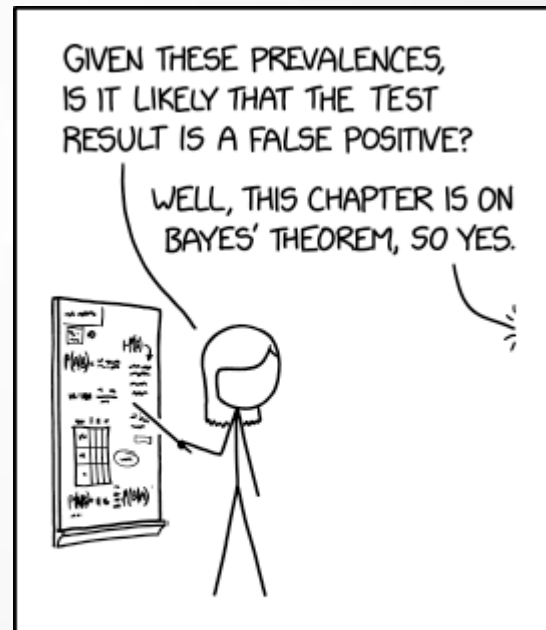
# Bayes' Theorem: Example

## Example

- 1% of people have a certain genetic defect
- 90% of tests accurately detect the gene (true positives).
- 7% of the tests are false positives.

$H$  - has genetic disorder  
 $E$  - positive test result

If Olaf gets a positive test result, what are the odds he actually has the genetic defect?



SOMETIMES, IF YOU UNDERSTAND  
BAYES' THEOREM WELL ENOUGH,  
YOU DON'T NEED IT.



# Bayes' Theorem: Example

- Hypothesis: Olaf has the gene,  $P(H) = 0.01$
- Evidence: Olaf got a positive test result,  $P(E)$
- True positive: Probability of positive test result if someone has the gene  $P(E|H) = 0.9$
- False positive: Probability of positive test result if someone doesn't have the gene  $P(E|\bar{H}) = 0.07$

# Bayes' Theorem: Example

Calculate

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)} = \frac{P(E|H)P(H)}{P(E|H) \cdot P(H) + P(E|\bar{H}) \cdot P(\bar{H})}$$
$$= \frac{0.9 \cdot 0.01}{0.9 \cdot 0.01 + 0.07 \cdot 0.99} = \dots = 0.115$$

↖  $1 - 0.01$

The probability that Olaf has the gene is only 11.5% despite the positive test result!

# Bayes' Theorem: Example

What happens if there are less false positives?

Consider  $P(E|\bar{H}) = 0.02$ :

$$P(H|E) = \frac{0.9 \cdot 0.01}{0.9 \cdot 0.61 + 0.02 \cdot 0.99} = \dots \approx 31\%$$

The probability that Olaf has the gene is now 31 %.

# Bayes' Theorem: Example

What happens if there are more true positives?

Consider  $P(E|H) = 0.95$ :

$$P(H|E) = \frac{0.95 \cdot 0.01}{0.95 \cdot 0.11 + 0.07 \cdot 0.99} = \dots = 0.12$$

Improving the accuracy of true positives raised the probability that Olaf has the gene to 12%.

# The Monty Hall Problem



# The Monty Hall Problem

You're in a game show.  
Behind one door is a car.  
Behind the others, goats.

You pick one of three doors,  
say #1.



The host, Monty Hall, opens one door, revealing...a goat!

# The Monty Hall Problem

He then gives you the option of switching doors or sticking with your original choice.



The question is: **should you switch?**

# The Monty Hall Problem

$\{1, 2, 3\}$  - door numbers

$C$  - door number with a car

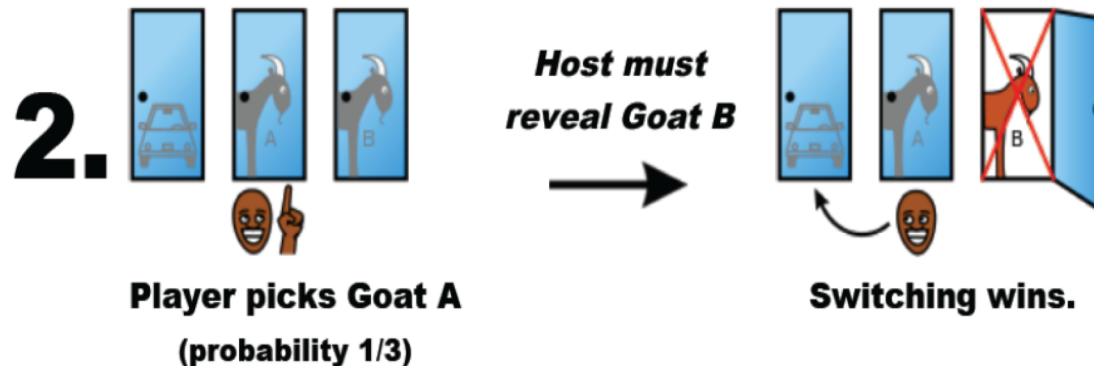
$S$  - door number player selected

$H$  - door number host opened



# The Monty Hall Problem

## The Monty Hall game



# The Monty Hall Problem

Selected  
door #1

Car location:	Host opens:	Total probability:	Stay:	Switch:
$\frac{1}{3}$ Door 1	$\frac{1}{2}$ Door 2	$\frac{1}{6}$	Car	Goat
$\frac{1}{3}$ Door 1	$\frac{1}{2}$ Door 3	$\frac{1}{6}$	Car	Goat
$\frac{1}{3}$ Door 2	1 Door 3	$\frac{1}{3}$	Goat	Car
$\frac{1}{3}$ Door 3	1 Door 2	$\frac{1}{3}$	Goat	Car

**Tree showing the probability of every possible outcome if the player initially picks Door 1**

# The Monty Hall Problem

$$\underline{S=1} \quad H=2$$

$$P(C=1 | H=2, \underline{S=1})$$

$$P(C=3 | H=2, \underline{S=1})$$

$$P(C=1 | H=2) = \frac{P(C=1 \cap H=2)}{P(H=2)}$$

$$P(C=1, H=2) = P(H=2 | C=1) \cdot P(C=1) = \frac{1}{2} \cdot \frac{1}{3} = \underline{\frac{1}{6}}$$

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

$$\rightarrow P(A \cap B) = P(A|B)P(B) = P(B|A) \cdot P(A)$$

can open 2 or 3

# The Monty Hall Problem

$$P(C=3 | H=2) = \frac{P(C=3 \cap H=2)}{P(H=2)}$$

$$P(C=3 \cap H=2) = P(H=2 | C=3) P(C=3) = 1 \cdot \frac{1}{3} = \underline{\frac{1}{3}}$$

has to open door 2

$$P(H=2) = P(H=2, C=1) + \underline{P(H=2, C=2)} + P(H=2, C=3)$$
$$= \frac{1}{6} + 0 + \frac{1}{3} = \underline{\underline{\frac{1}{2}}}$$

$$P(C=1 | H=2) = \frac{1/6}{1/2} = \frac{1}{3} < P(C=3 | H=2) = \frac{1/3}{1/2} = \frac{2}{3}$$

conditioned on  $S=1$