# DSC 40A

Theoretical Foundations of Data Science I

# Question

Remember, you can always ask questions at
q.dsc40a.com!
If the direct link doesn't work, click the "Lecture
Questions" link in the top right corner of dsc40a.com.

**Example 15.** What is the probability that your five-card poker hand is a straight?

52 cards

4 suits

A 2 3 4 5 6 7 8 9 10 J Q K, A

♡ ◇ ♣ ♠

5 cards

$\_\_ \_\_ \_\_ \_\_ \_\_$

4 possible suits

$S$ = set of 5 cards (equally likely)

card = value + suit

straight   A 2 3 4 5   $\longrightarrow$   10 J Q K A

$$\text{prob}(\text{straight}) = \frac{|E|}{|S|} = \frac{\#\text{straight hands}}{\#\text{sets of 5 cards}} = \frac{10 \cdot 4^5}{C(52,5)}$$

**Example 16.** Suppose you look at your first card as it is dealt, and you see that it is a Queen. What is the probability that your five-card hand is a straight?

$S = $ set of 4 cards (not including Q)

$$\text{prob}\binom{\text{straight}}{\text{w/ Q}} = \frac{\text{\# set of 4 cards that make straight w/ Q}}{\text{\# sets of 4 cards (out of 51)}}$$

$$= \frac{3 \cdot 4^4}{C(51,4)}$$

10 J Q K A

9 to J Q u

8 9 to J Q

# Agenda

- Law of total probability.

- Bayes theorem.

# Getting to Campus

- You conduct a survey:
  - How did you get to campus today? Walk, bike, or drive?
  - Were you late?

*Walk & late*

*bike and not late*

|        | Late | Not Late |
|--------|------|----------|
| Walk   | 6%   | 24%      |
| Bike   | 3%   | 7%       |
| Drive  | 36%  | 24%      |

*all sum to 100%*

# Getting to Campus

|        | Late | Not Late |
|--------|------|----------|
| Walk   | 6%   | 24%      |
| Bike   | 3%   | 7%       |
| Drive  | 36%  | 24%      |

Prob(bike)

What is the probability that a randomly selected person is late?
- A. 24%
- B. 30%
- C. 45%
- D. 50%

$6\% + 3\% + 36\% = 45\%.$

# Getting to Campus

|  | Late $B$ | Not Late |
|---|---|---|
| Walk $A_1$ | 6% | 24% |
| Bike $A_2$ | 3% | 7% |
| Drive $A_3$ | 36% | 24% |

- Since everyone either walks, bikes, or drives,

P(Late) = P(Late AND Walk) + P(Late AND Bike) + P(Late AND Drive)

$$P(B) = P(B \cap A_1) + P(B \cap A_2) + P(B \cap A_3)$$

- This is called the **Law of Total Probability**.

# Getting to Campus

|  | Late | Not Late |
|---|---|---|
| Walk | 6% | 24% |
| Bike | 3% | 7% |
| Drive | 36% | 24% |

$\cap = $ And $= +$

Suppose someone tells you that they walked. What is the probability that they were late?
- A. 6%
- B. 20%
- C. 25%
- D. 45%

Hint: conditional probability!

$P\left(\text{late} \mid \text{walk}\right) = \dfrac{P(\text{late} \cap \text{walk})}{P(\text{walk})_{\to >0}} = \dfrac{6\%}{30\%} = \dfrac{1}{5} = 20\%.$

Need $P(\text{walk}) = P(\text{walk} + \text{late}) + P(\text{walk} + \text{not late})$
$= 6\% + 24\% = 30\%$

multiplication rule
$P(\text{late} \cap \text{walk}) = P(\text{late} \mid \text{walk}) \, P(\text{walk})$

# Getting to Campus

|       | Late | Not Late |
|-------|------|----------|
| Walk  | 6%   | 24%      |
| Bike  | 3%   | 7%       |
| Drive | 36%  | 24%      |

- Since everyone either walks, bikes, or drives,

P(Late) = P(Late AND Walk) + P(Late AND Bike) + P(Late AND Drive)

P(Late)  = P(Late|Walk)*P(Walk) + P(Late|Bike)*P(Bike)
          +P(Late|Drive)*P(Drive)
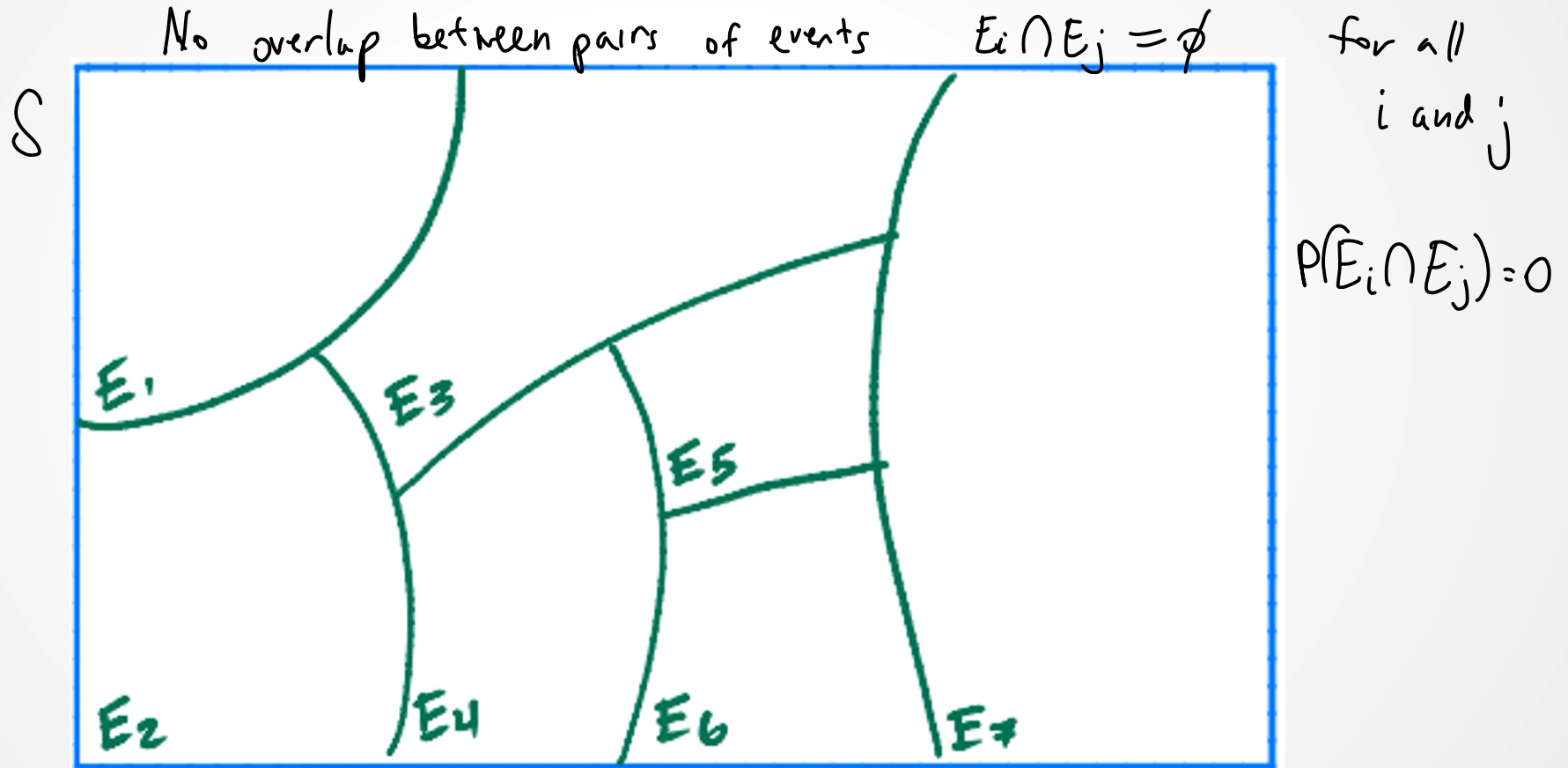
# Partitions

- A set of events $E_1$, $E_2$, ..., $E_k$ is a **partition** of $S$ if
  - $P(E_i \cap E_j) = 0$ for all $i,j$ $\longleftarrow$ mutually exclusive
  - $P(E_1) + P(E_2) + ... + P(E_k) = 1 = \sum_k P(E_k) = P(S)$

Every outcome in $S$ belongs to only one of the $E_k$'s

Walk

| | 20% |
|---|---|
| bike | 10% |
| drive | 60% |

Late    Not late

| Late | Not late |
|---|---|
| 45% | 55% |

either late or not late
cannot be both

# Partitions

No overlap between pairs of events $E_i \cap E_j = \phi$ for all $i$ and $j$

$P(E_i \cap E_j) = 0$

$S$

$E_1$

$E_3$

$E_5$

$E_2$

$E_4$

$E_6$

$E_7$

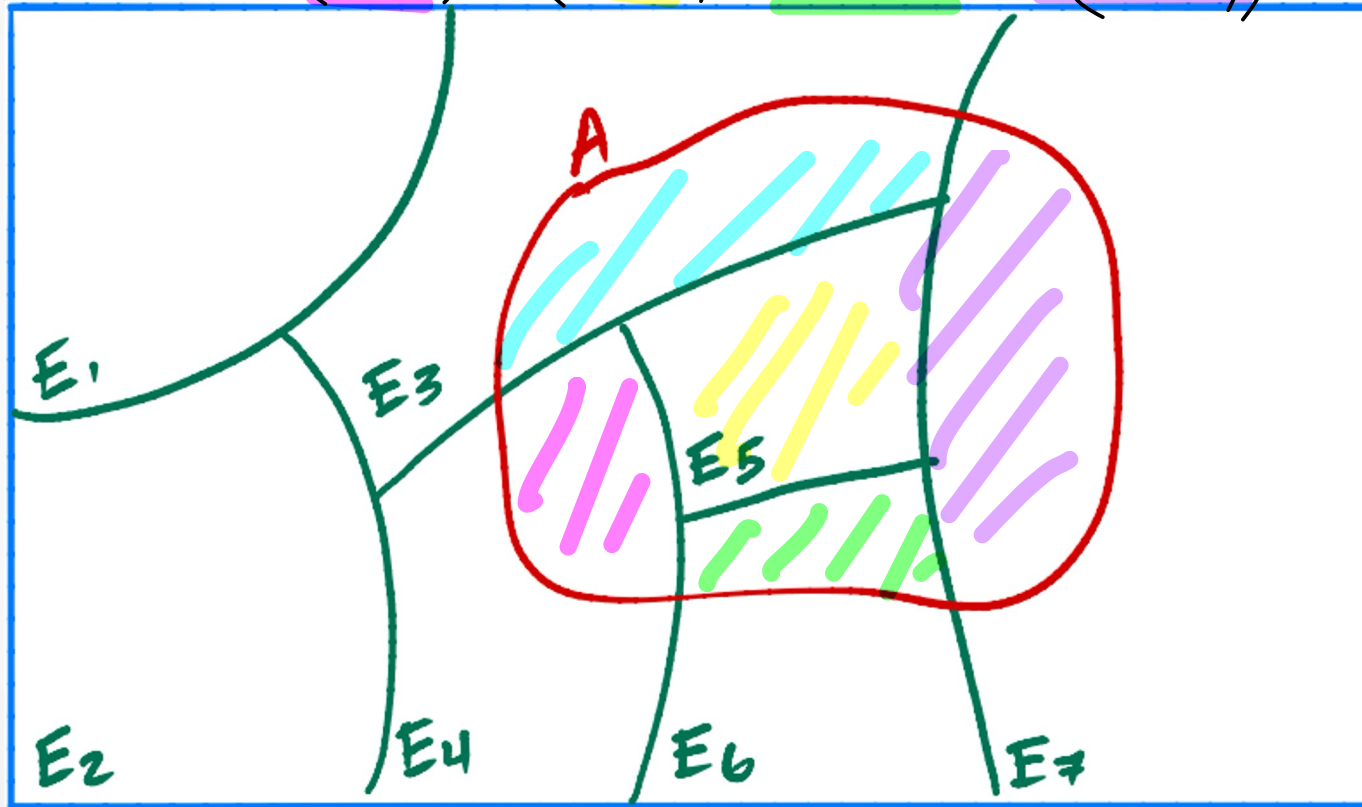# Law of Total Probability

(mutually exclusive)

- If $A$ is an event and $E_1, E_2, \ldots, E_k$ is a **partition** of $S$, then

$$P(A) = P(A \cap E_1) + P(A \cap E_2) + \ldots + P(A \cap E_k)$$

$$= \sum_{i=1}^{k} P(A \cap E_i)$$

# Law of Total Probability

$$P(A) = P(A \cap E_3) + P(A \cap E_4) + P(A \cap E_5) + P(A \cap E_6) + P(A \cap E_7)$$

$A \subset S$

$A \cap E_1 = \emptyset$

$A \cap E_2 = \emptyset$

$P(A \cap E_1) = 0$

$= P(A \cap E_2)$

$E_1$

$E_3$

$A$

$E_5$

$E_2$

$E_4$

$E_6$

$E_7$

# Law of Total Probability

- If $A$ is an event and $E_1$, $E_2$, …, $E_k$ is a **partition** of $S$, then

$$P(A) = P(A \cap E_1) + P(A \cap E_2) + ... + P(A \cap E_k)$$

$$= \sum_{i=1}^{k} P(A \cap E_i)$$

- Written another way,

$$P(A) = P(A \mid E_1) \cdot P(E_1) + ... + P(A \mid E_k) \cdot P(E_k)$$

$$\underbrace{P(A \cap E_1)}_{} \longrightarrow \text{multiplication rule}$$

$$= \sum_{i=1}^{k} P(A \mid E_i) \cdot P(E_i)$$

# Getting to Campus

|  | Late | Not Late |
|---|---|---|
| Walk | 6% | 24% |
| Bike | 3% | 7% |
| Drive | 36% | 24% |

= 45%

Suppose someone is late. What is the probability that they walked? Choose the best answer.
- A. Close to 5%
- B. Close to 15%
- C. Close to 30%
- D. Close to 40%

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

$$P(walk \mid late) = \frac{P(walk \cap late)}{P(late)} =$$

$$= \frac{6\%}{45\%} \approx 13\%$$

$$P(walk \mid late) \neq P(late \mid walk)$$

# Getting to Campus

- Suppose all you know is
  - P(Late) = 45%
  - P(Walk) = 30%
  - P(Late|Walk) = 20%
- Can you still find P(Walk|Late)?

$$P(\text{walk}|\text{late}) = \frac{P(\text{walk} \cap \text{late})}{P(\text{late})} = \frac{P(\text{late}|\text{walk}) \cdot \text{Prob}(\text{walk})}{P(\text{late})} = \frac{0.2 \cdot 0.3}{0.45} \approx 13\%.$$

# 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:

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

$> 0$

can calculate $P(A)$ using law of total prob.

# 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:

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

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

not B

# Bayes' Theorem

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

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

- $P(H)$ - 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 $H$

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

# Bayes' Theorem: Example

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

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%

# 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**

# 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.

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

# Bayes' Theorem: Example

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

# Bayes' Theorem: Example

Calculate

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

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

# Bayes' Theorem: Example

What happens if there are less false positives?
Consider $P(E|\bar{H}) = 0.02$:

The probability that Olaf has the gene is now _____.

What happens if there are more true positives?
Consider $P(E|H) = 0.95$:

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

# Preview: Bayes' Theorem for Classification

Bayes' Theorem is very useful for classification problems, where we want to predict a class based on some features.

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

A = having certain features

$$P(\text{class}|\text{features}) = \frac{P(\text{features}|\text{class}) * P(\text{class})}{P(\text{features})}$$

# Summary

- When a set of events partitions the sample space, the law of total probability applies.

$$P(A) = P(A \cap E_1) + P(A \cap E_2) + ... + P(A \cap E_k)$$

$$= \sum_{i=1}^{k} P(A \cap E_i)$$

- Bayes Theorem says how to express $P(B|A)$ in terms of $P(A|B)$.

- **Next time:** independence and conditional independence