

DSC 40A

Theoretical Foundations of Data Science I

Random Sampling

Announcements

- Groupwork due tonight
- Homework 5 due Friday
- Upcoming homework schedule:
Homework 6 released Monday 11/18 and due 11/25

* New FAQ on Convexity

Agenda

- Conditional probability continued
- Sampling with and without replacement

Question

Answer at q.dsc40a.com

Remember, you can always ask questions at
q.dsc40a.com!

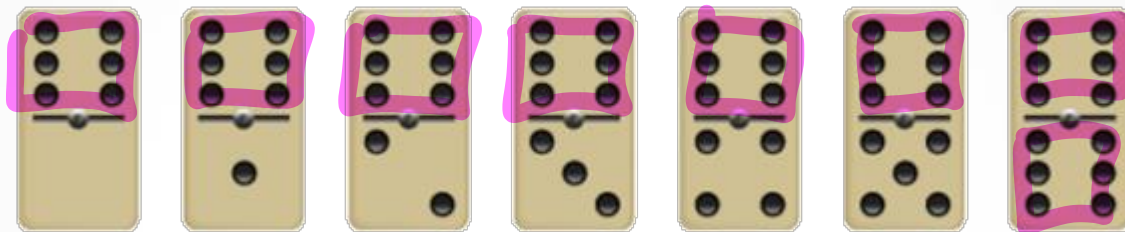
If the direct link doesn't work, click the "Lecture Questions" link in the top right corner of dsc40a.com.

Conditional probability continued



Dominoes


Question 3: Now you pick a random tile from the set and uncover only one side, revealing that it has 6 dots. What is the probability that this tile is a double, with 6 on both sides?




$S =$ Dominoes with side that has a 6

$E =$ Domino that is a double 6

$$P(E) = \frac{\# \text{outcomes in } E}{\# \text{outcomes in } S} = \frac{2}{8} = \frac{1}{4} > \frac{1}{7}$$

$S = 28 \cdot 2$ possible layouts


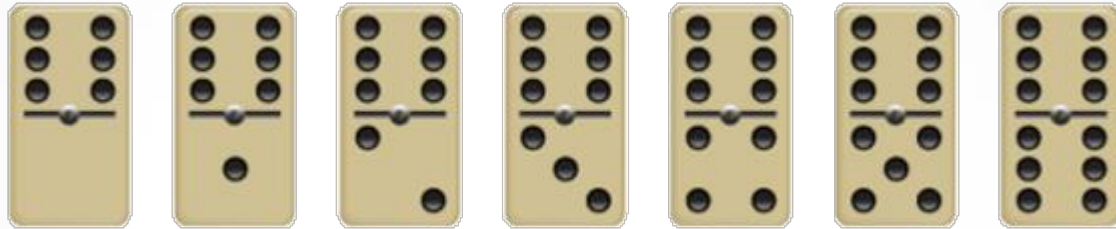
$E =$ both sides are the same: $7 \cdot 2 = 14$



$F =$ layout in which 6 was observed: 8 layouts

Dominoes

Question 3: Now you pick a random tile from the set and uncover only one side, revealing that it has 6 dots. What is the probability that this tile is a double, with 6 on both sides?



$$P(E|F) = \frac{P(E \cap F)}{P(F)} = \frac{2/56}{8/56} = \frac{1}{4}$$

Try it out in [code!](#)

Conditional probabilities: Simpson's Paradox

	Treatment A	Treatment B
Small kidney stones	81 successes / 87 (93%)	234 successes / 270 (87%)
Large kidney stones	192 successes / 263 (73%)	55 successes / 80 (69%)
Combined	273 successes / 350 (78%)	289 successes / 350 (83%)

Which treatment is better?

~19%. A. Treatment A for all cases.

~31%. B. Treatment B for all cases.

6%. C. A for small and B for large.

43%. D. A for large and B for small.

Conditional probabilities: Simpson's Paradox

	Treatment A	Treatment B
Small kidney stones	81 successes / 87 (93%)	234 successes / 270 (87%)
Large kidney stones	192 successes / 263 (73%)	55 successes / 80 (69%)
Combined	273 successes / 350 (78%)	289 successes / 350 (83%)

Simpson's Paradox

"When the less effective treatment is applied more frequently to easier cases, it can appear to be a more effective treatment."

Random Sampling

The background of the slide is white with abstract green geometric shapes on the right side. These shapes include overlapping triangles and polygons in various shades of green, from light lime to dark forest green. A thin, light gray line also extends from the bottom right towards the center of the slide.

Sampling

Sampling with replacement:

1. Draw one element uniformly at random from list.
2. Return the element to the list.
3. Repeat

Sampling without replacement:

Same as above without step 2

What does *uniformly at random* mean?

Sampling

Sampling with or without replacement:

- All samples are equally likely.
- Uniform distribution!

$P(\text{sample having a certain property}) =$

Sampling

Sampling with or without replacement:

- All samples are equally likely.
- Uniform distribution!

$$P(\text{sample having a certain property}) = \frac{\# \text{ samples having property}}{\# \text{ possible samples}}$$

Practice Problems

Example 5. There are 20 students in a class. A computer program selects a random sample of students by drawing 5 students at random **with replacement**. What is the chance that a particular student is among the 5 selected students?

Sample space: sequences of length S
with entries $\{1, 2, \dots, 20\}$

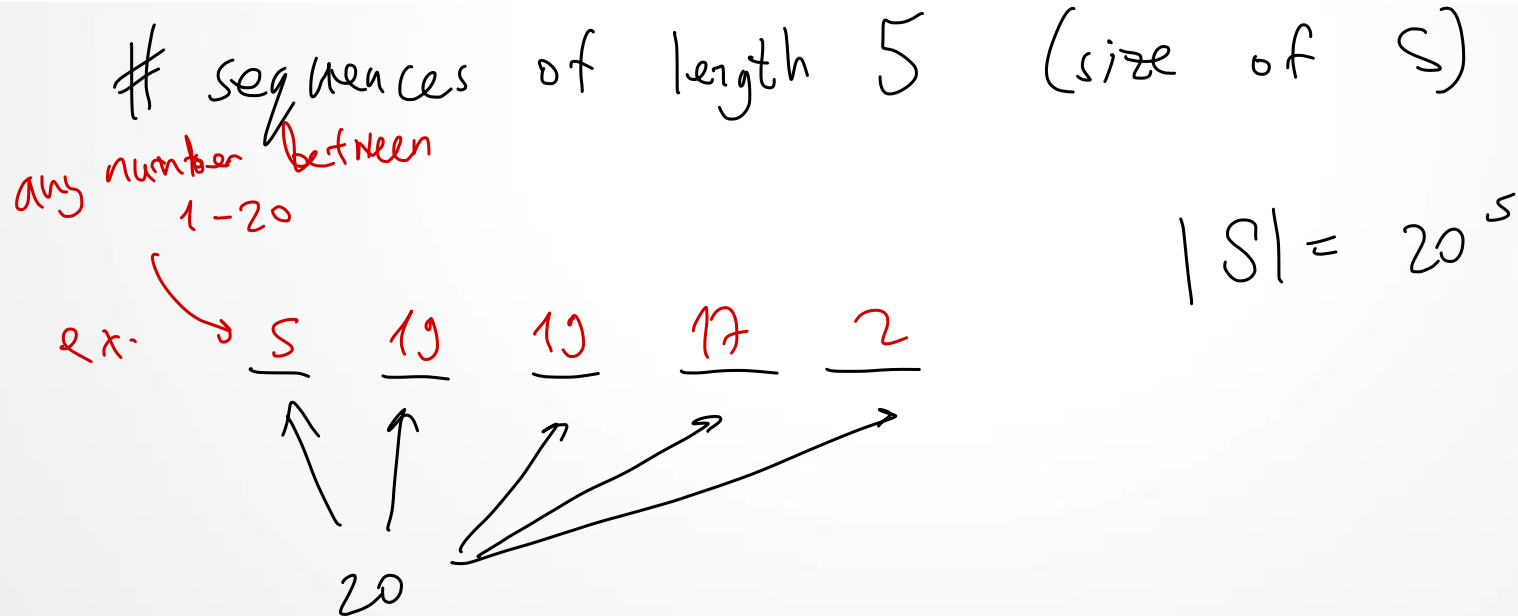
Examples: $3, 12, 7, 20, 14$
 $3, 3, 3, 20, 1$

particular student: 17

$$\frac{\# \text{ sequences of length } S \text{ that include } 17}{\# \text{ sequences of length } S}$$

Practice Problems

Part 1. Denominator. If you draw a sample of size 5 at random with replacement from a population of size 20, how many different sequences of individuals could you draw?



Practice Problems

Part 2. Numerator. If you draw a sample of size 5 at random with replacement from a population of size 20, how many different sequences of individuals include a particular person?

sequences of length 5 that include 17

ex: 17 3 2 20 8

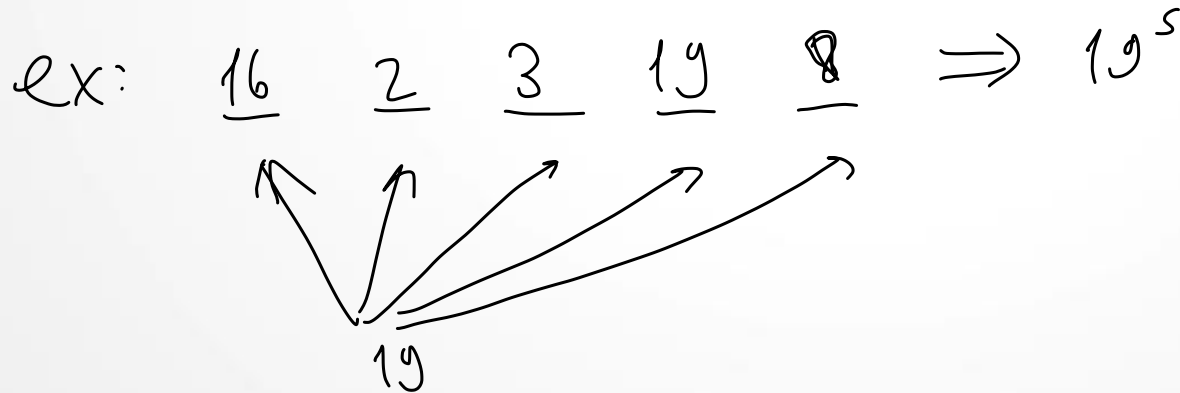
16 17 18 19 20

17 17 17 17 17

Practice Problems

Using the complement. If you draw a sample of size 5 at random with replacement from a population of size 20, how many different sequences of individuals **do not** include a particular person?

sequences that don't include 17



Practice Problems

Example 5. There are 20 students in a class. A computer program selects a random sample of students by drawing 5 students at random **with replacement**. What is the chance that a particular student is among the 5 selected students?

$$P(\text{sequence of length } 5 \text{ including } 17) = \frac{\# \text{ seq inc. } 17}{\# \text{ seq in } S} = \frac{\# \text{ seqs in } S - \# \text{ seqs in } S \text{ w/o } 17}{\# \text{ seqs in } S}$$

$$\Rightarrow \frac{20^5 - 19^5}{20^5} = 1 - \left(\frac{19}{20}\right)^5 \approx 0.226$$

$$1 - P(\text{sequence of length } 5 \text{ not including } 17)$$

Practice Problems

Example 6. There are 20 students in a class. A computer program selects a random sample of students by drawing 5 students at random **without replacement**. What is the chance that a particular student is among the 5 selected students?

ex: 16, 17, 18, 19, 20 ✓

ex: 17, 17, 17, 36 ✗

Which probability will be higher?

- A. Probability of including a particular student when sampling with replacement.
- B. Probability of including a particular student when sampling without replacement.
- C. Both probabilities are the same.

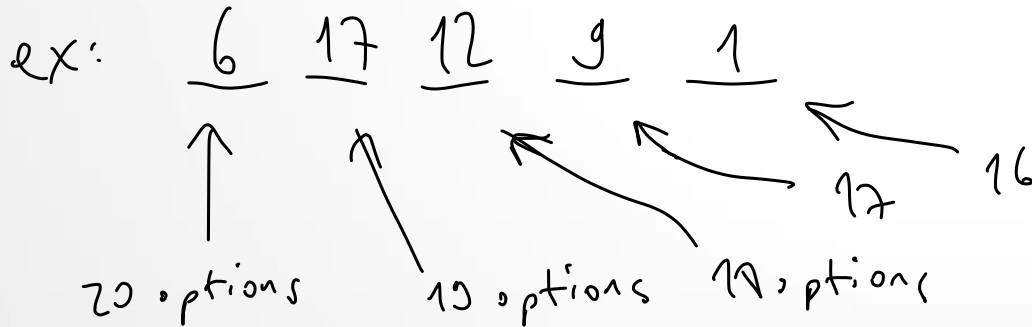
S : sequences of length 5 without repeats
 $\{1, 2, \dots, 20\}$

Practice Problems

Part 1. Denominator. If you draw a sample of size 5 at random without replacement from a population of size 20, how many different sequences of individuals could you draw?

sequences in S

$$20 \cdot 19 \cdot 18 \cdot 17 \cdot 16 = \frac{20!}{15!}$$



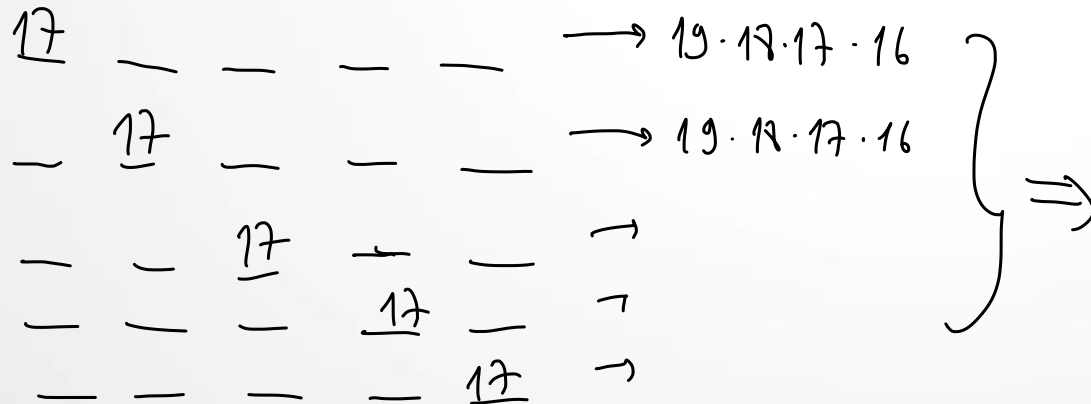
$$\frac{1 \cdot 2 \cdot 3 \cdot 4 \cdot \dots \cdot 15 \cdot 16 \cdot 17 \cdot 18 \cdot 19 \cdot 20}{1 \cdot 2 \cdot 3 \cdot 4 \cdot \dots \cdot 15}$$

Practice Problems

Part 2. Numerator. If you draw a sample of size 5 at random without replacement from a population of size 20, how many different sequences of individuals include a particular person?

ex: 2, 3, 4, 17, 19

17, 16, 3, 4, 8



$$19! - 15! = (1 \cdot 2 \cdot 3 \cdot \dots \cdot 15) \cdot 16 \cdot 17 \cdot 18 \cdot 19$$

$$- (1 \cdot 2 \cdot 3 \cdot \dots \cdot 15) =$$

$$(16 \cdot 17 \cdot 18 \cdot 19 - 1) \cdot 15!$$

$$5 \cdot 19 \cdot 18 \cdot 17 \cdot 16 = 5 \cdot \frac{19!}{15!}$$

Practice Problems

Using the complement. If you draw a sample of size 5 at random without replacement from a population of size 20, how many different sequences of individuals **do not** include a particular person?

sequences of length 5 not including 17

$$\begin{array}{ccccccccc} \text{ex} & \underline{2} & \underline{20} & \underline{4} & \underline{3} & \underline{8} & \Rightarrow & 19 \cdot 18 \cdot 17 \cdot 16 \cdot 15 = & \frac{19!}{14!} \\ & \uparrow & \uparrow & \uparrow & \uparrow & \uparrow & & & \\ & 19 & 18 & 17 & 16 & 15 & & & \end{array}$$

Practice Problems

Example 6. There are 20 students in a class. A computer program selects a random sample of students by drawing 5 students at random **without replacement**. What is the chance that a particular student is among the 5 selected students?

$$\begin{aligned}
 P(\text{selecting } 17 \text{ in seq. of length } 5 \text{ w.o. replacement}) &= \frac{\# \text{ seq in } S \text{ w } 17}{\# \text{ seq in } S} && \left| \frac{\# \text{ seqs in } S - \# \text{ seqs in } S \text{ w.o. } 17}{\# \text{ seqs in } S} \right. \\
 &= \frac{5 \cdot \frac{19!}{\cancel{15!}}}{\frac{20!}{\cancel{15!}}} = 5 \cdot \frac{19!}{20!} = \frac{5}{20} && \left| \frac{20 \cdot 19 \cdot 18 \cdot 17 \cdot 16 - 19 \cdot 18 \cdot 17 \cdot 16 \cdot 15}{20 \cdot 19 \cdot 18 \cdot 17 \cdot 16} \right. \\
 &= \frac{1}{4} = 0.25 > 0.226 && \left| \frac{20 - 15}{20} = \frac{5}{20} = \frac{1}{4} \right.
 \end{aligned}$$

Summary

- When we sample uniformly, whether with or without replacement, each possible sample is equally likely.
- Probability questions become counting questions:

$$P(\text{sample having a certain property}) = \frac{\# \text{ samples having property}}{\# \text{ possible samples}}$$

- **Next time:** combinatorics, or counting principles