

LECTURE 18

Sampling

How to sample effectively, and how to quantify the samples we collect.

CSCI 3022

Maribeth Oscamou

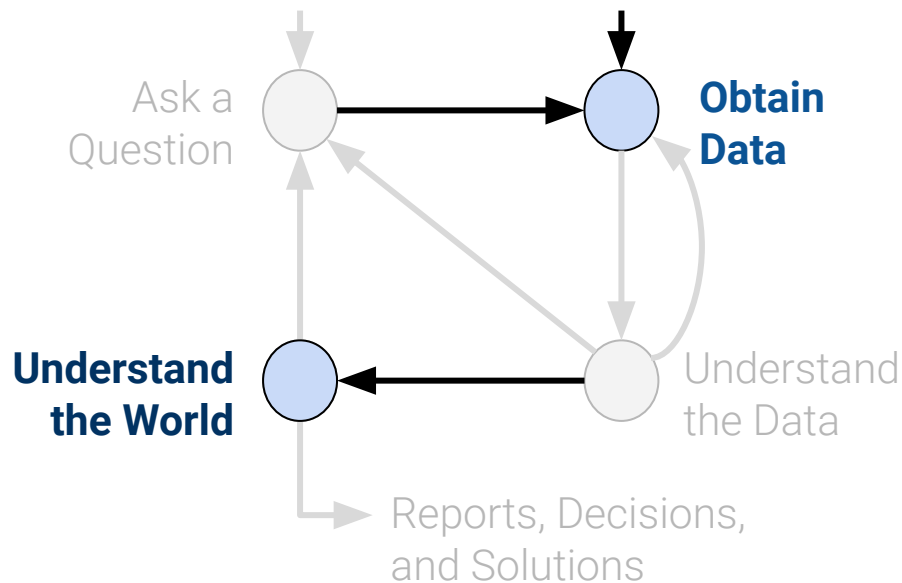
Content credit: [Acknowledgments](#)

Today's Topic: Sampling

Understanding the sampling process is what lets us go from **describing the data** to **understanding the world**

Without knowing / assuming something about how the data were collected:

- There is no connection between the **sample** and the **population**
- The **data set** doesn't tell us about the **world behind the data**



Today's Roadmap

CSCI 3022

- Censuses and Surveys
- Sampling: Definitions
- Sampling Bias: A Case Study
- Probability Samples
- Multinomial Probabilities

Sampling: Definitions

CSCI 3022

- Sampling: Definitions
- Sampling Bias: A Case Study
- Probability Samples
- Multinomial Probabilities

Censuses and Surveys

A **census** is “an official count or survey of a **population**, typically recording various details of individuals.”

A **survey** is a set of questions.

- For instance: census workers survey individuals and households.

What is asked, and how it is asked, can affect:

- How the respondent answers.
- **Whether** the respondent answers.

There are entire courses on surveying!

FiveThirtyEight

Politics Sports Science Podcasts Video Interactives

JUN. 27, 2019, AT 12:42 PM

The Supreme Court Stopped The Census Citizenship Question — For Now

NATIONAL

Citizenship Question To Be Removed From 2020 Census In U.S. Territories

August 9, 2019 · 3:23 PM ET

[FiveThirtyEight](#), [NPR](#)

Sampling from a finite population

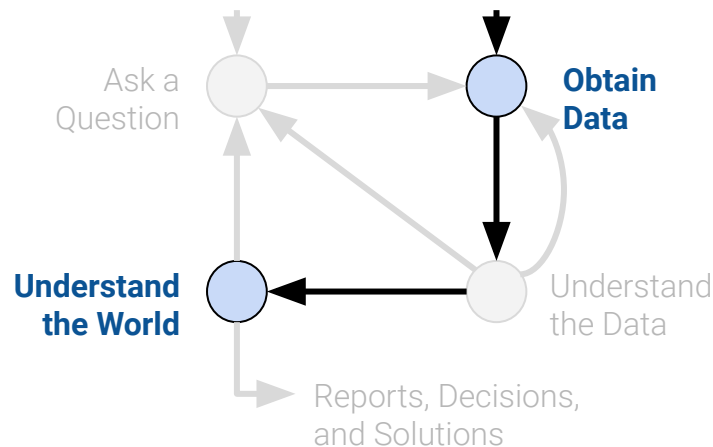
A census is great, but expensive and difficult to execute.

- Would **all** voters be willing to participate in a voting census prior to an actual election?

A **sample** is (usually) a subset of the population.

- Samples are often used to make **inferences about the population**.
- How you draw the sample will affect your accuracy.
- Two common sources of error:
 - **chance error**: random samples can vary from what is expected, in any direction.
 - **bias**: a systematic error in one direction.
 - Could come from our sampling scheme, and survey methods.

Inference: drawing conclusions (and quantifying their reliability) about a population based on a sample. [Data 8 book](#)



Other kinds of populations

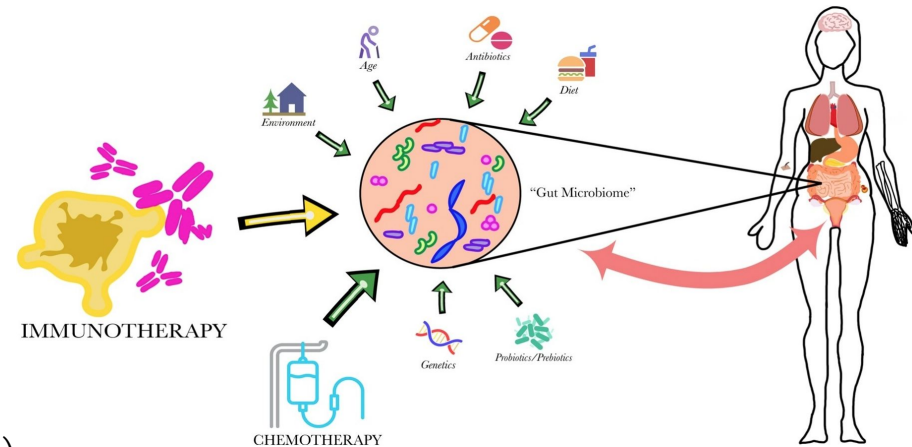
The individuals in a population are not always people!

Could be

- **Bacteria** in your gut (sampled using DNA sequencing)
- **Trees** of a certain species
- **Small businesses** receiving a microloan
- **Published results** in a journal / field ([example](#))

In any of these cases we might examine a sample and try to draw an inference about the population it came from.

- Simplest example: what % have some binary property (like voting intention)?



Population, sample, and sampling frame

Population: The group that you want to learn something about.

Sampling Frame: The list from which the sample is drawn.

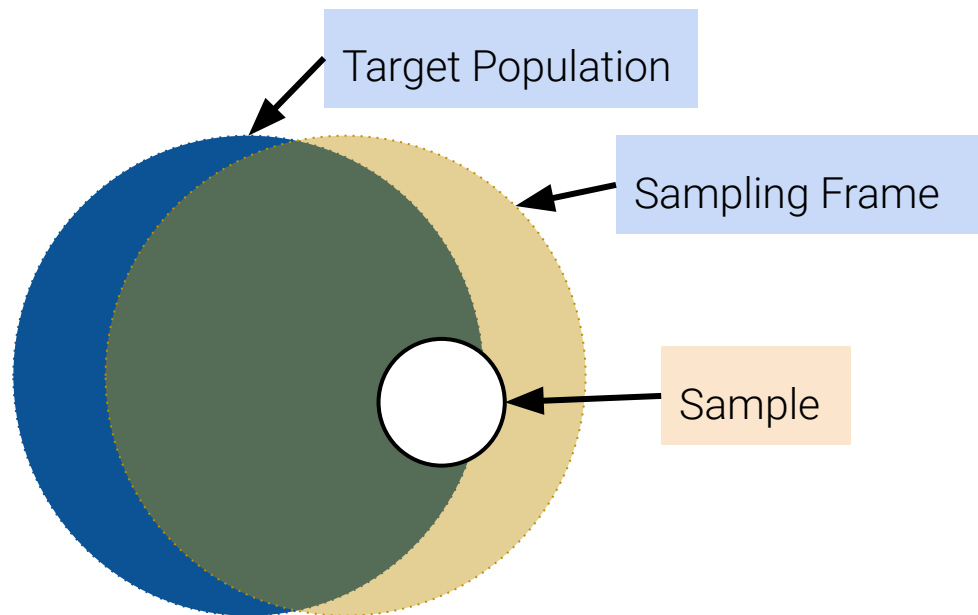
- If you're sampling people, the sampling frame is the set of all people that could possibly end up in your sample.

Sample: Who you actually end up sampling.

- A subset of your sampling frame.

There may be individuals in your **sampling frame** (and hence, your sample) that are **not** in your population!

Similarly, there might be individuals in your target population that are not in your sampling frame.



Bias: A Case Study

CSCI 3022

- Censuses and Surveys
- Sampling: Definitions
- **Sampling Bias: A Case Study**
- Probability Samples
- Multinomial Probabilities

Case study: 1936 Presidential Election



Roosevelt (D)



Landon (R)

In 1936, President Franklin D. Roosevelt (left) went up for re-election against Alf Landon (right). As is usual, **polls** were conducted in the months leading up to the election to try and predict the outcome.

(Election result spoiler: Landon was not a [U.S. President](#))

The Literary Digest: Election Prediction

The *Literary Digest* was a magazine. They had successfully predicted the outcome of 5 general elections coming into 1936.

They sent out their survey to **10,000,000** individuals, who they found from:

- Phone books.
- Lists of magazine subscribers.
- Lists of country club members.

	% Roosevelt	# surveyed
Actual election	61%	All voters (~45,000,000)
The Literary Digest poll	43%	10,000,000



How could this have happened?
They surveyed 10 million people!

The Literary Digest: What happened?

(1) The Literary Digest sample was **not representative** of the population.

- The Digest's **sampling frame**: people in the phonebook, subscribed to magazines, and went to country clubs.
- These people were more affluent and tended to vote Republican (Landon).

	% Roosevelt	# surveyed
Actual election	61%	All voters (~45,000,000)
The Literary Digest poll	43%	10,000,000

(2) Only 2.4 million people **actually filled out the survey!**

- 24% response rate (low).
- Who knows how the 76% **non-respondents** would have polled?

The Literary Digest

NEW YORK OCTOBER 31, 1936

Topics of the day

LANDON, 1,293,669; ROOSEVELT, 972,897

Final Returns in The Digest's Poll of Ten Million Voters

Well, the great battle of the ballots in the Poll of ten million voters, scattered throughout the forty-eight States of the

Republican National Committee purchased THE LITERARY DIGEST?" And all types and varieties, including: "Have the Jews purchased

returned and let the people of draw their conclusions as to o So far, we have been right in Will we be right in the current as Mrs. Roosevelt said concerni dent's reelection, is in the 'lap. "We never make any claims tion but we respectfully refer minion of one of the most an

Gallup's Poll: Election Prediction

George Gallup, a rising statistician, also made predictions about the 1936 elections.
His estimate was **much** closer despite having a smaller **sample size** of “only” 50,000

(Also more than necessary!)

George Gallup also predicted what The Literary Digest was going to predict, within 1%, with a **sample size of only 3000 people**.

- He predicted the Literary Digest’s **sampling frame** (phonebook, magazine subscribers, country clubs).
- So he sampled those same individuals!

	% Roosevelt	# surveyed
Actual election	61%	All voters (~45,000,000)
The Literary Digest poll	43%	10,000,000
George Gallup’s poll	56%	50,000
George Gallup’s prediction of Digest’s prediction	44%	3,000

Samples, while convenient, are subject to chance error and **bias**.

No Quiz this week!

HW 7 due FRIDAY (not Thursday) at 11:59pm MT

COMMON NOTATION for expected value and variance:

$$\mu = \mu_X = E[X]$$

$$\sigma^2 = \sigma_X^2 = Var[X]$$

Selection Bias

- Systematically excluding (or favoring) particular groups.
- **Example:** The Literary Digest poll excludes people not in phone books.
- **How to avoid:** Examine the sampling frame and the method of sampling.

Response Bias

- People don't always respond truthfully, or questions lead to certain responses.
- **Example:** Asking citizenship questions on the census survey→illegal immigrants might not answer truthfully
- **How to avoid:** Response bias exists in ANY survey. However, we can try to minimize it by examining the nature of questions and the method of surveying.

Non-response Bias

- People don't always respond → People who don't respond aren't like the people who do!
- **Example:** Only 2.4m out of 10m people responded to The Literary Digest poll.
- **How to avoid:** Keep your surveys short, and be persistent.

Probability Samples

CSCI 3022

- Censuses and Surveys
- Sampling: Definitions
- Bias: A Case Study
- **Probability Samples**
- Multinomial Probabilities

A **huge sample size** does not fix a **bad sampling method**!

We want the sample to be **representative** of the population.

Think about **tasting soup**: if it's **well-stirred**, a spoonful is all you need!

- Don't just try to get a BIG sample. If your method of sampling is BAD, and your sample is BIG, what you'll have is a BIG BAD sample

Easiest way to to get a representative sample is by using **randomness**.



- Random (aka Probability) sample:
 - Before the sample is drawn, you have to know the selection probability of every group of people in the population
 - Not all individuals / groups have to have equal chance of being selected

Common Non-Random Sample

A **convenience sample** is whoever you can get ahold of.

Example: Suppose we have a cage of mice, and each week, we want to measure the weights of these mice. To do so, we take a convenience sample of the first 5 of the mice we can grab and weigh them.



- Not a good idea for inference!
- Haphazard \neq random.
- Sources of bias can introduce themselves in ways you may not think of!

Sample of Convenience (NOT random)

A **convenience sample**:

- **Example:** *sample consists of whoever visits your website*
- Just because you think you're **sampling "randomly"**, doesn't mean you have a random sample.
- If you can't figure out **ahead of time**
 - what's the population
 - what's the **chance of selection**, for each group in the populationthen you **don't have a random sample**

Warning:

- Haphazard \neq **random**.
- Many potential sources of bias!

Probability Sample (aka Random Sample)

Why sample at random?

1. (As mentioned before) To get more representative samples → **reduce bias**
 - Random samples **can** produce biased estimates of population quantities. (For example, if we're estimating the maximum of a population)
2. More importantly, with random samples we can **estimate** the **bias** and **chance error** → **quantify uncertainty**

Probability Sample (aka Random Sample)

Why sample at random?

1. (As mentioned before) To get more representative samples → **reduce bias**
 - Random samples **can** produce biased estimates of population quantities.
2. More importantly, with random samples we can **estimate** the **bias** and **chance error** → **quantify uncertainty**

For a **probability sample**,

- We have to be able to provide the **chance** that any specified **set** of individuals will be in the sample.
- All individuals in the population **need not** have the same chance of being selected.
- Because we know all the probabilities, we will be able to **measure the errors**.

The real world is usually more complicated!

- Election polling: When Gallup calls, most people don't answer.
- Bacteria: We don't know the probability a given bacterium will get into a microbiome sample.

If the sampling / measurement process isn't fully under our control, we try to **model it**.

Example Scheme 1: Probability Sample

Suppose I have 3 TA's (**A**lan, **B**ennett, **C**eline):

I decide to sample 2 of them as follows:

- I choose **A** with probability 1.0
- I choose either **B** or **C**, each with probability 0.5.

All subsets of 2: {**A, B**} {**A, C**} {**B, C**}

Probabilities: 0.5 0.5 0

This is a **probability sample** (though not a great one).

- Of the 3 people in the population, I know the chance of getting each subset.
- Suppose I'm measuring the average distance TA's live from campus.
 - This scheme does not see the entire population!
 - My estimate using the single sample I take has some **chance error** depending on if I see AB or AC.
 - This scheme **biases** towards A's response

Common random sampling schemes

A **random sample with replacement** is a sample drawn **uniformly** at random **with** replacement.

- Random doesn't always mean "uniformly at random," but in this specific context, it does.
- Some individuals in the population might get picked more than once



A **simple random sample (SRS)** is a sample drawn **uniformly** at random **without** replacement.

- **Every individual has the same chance of being selected.**
- Every pair has the same chance as every other pair.
- Every triple has the same chance as every other triple.
- And so on.

A raffle could use either sampling scheme, depending on if winners are eligible for multiple prizes.

Example Scheme 2: Simple Random Sample?

We have the following sampling scheme:

- A class roster has 1100 students listed alphabetically.
- Pick one of the first 10 students on the list at random (e.g. [Student 8](#)).
- To create your sample, take that student and every 10th student listed after that (e.g. [Students 8, 18, 28, 38](#), etc).

1. Is this a probability sample?

2. Does each student have the same probability of being selected?

3. Is this a simple random sample?

Example Scheme 2: Simple Random Sample?

Consider the following sampling scheme:

- A class roster has 1100 students listed alphabetically.
- Pick one of the first 10 students on the list at random (e.g. [Student 8](#)).
- To create your sample, take that student and every 10th student listed after that (e.g. [Students 8, 18, 28, 38](#), etc).

1. Is this a probability sample?

Yes.

For a sample $[n, n + 10, n + 20, \dots, n + 1090]$, where $1 \leq n \leq 10$, the probability of that sample is $1/10$.

Otherwise, the probability is 0.

Only 10 possible samples!

2. Does each student have the same probability of being selected?

Yes.

Each student is chosen with probability $1/10$.

3. Is this a simple random sample?

No.

The chance of selecting (8, 18) is $1/10$; the chance of selecting (8, 9) is 0.

This method is called a **systematic sample**

Common random sampling schemes

A **random sample with replacement** is a sample drawn **uniformly** at random **with** replacement.

- Random doesn't always mean "uniformly at random," but in this specific context, it does.
- Some individuals in the population might get picked more than once

A **systematic sample**: Order the sample frame. Choose an integer k . Sample every k th unit in the sample frame.

A **simple random sample (SRS)** is a sample drawn **uniformly** at random **without** replacement.

- **Every individual has the same chance of being selected.**
- Every pair has the same chance as every other pair.
- Every triple has the same chance as every other triple.
- And so on.

A **stratified random sample**: if you have a heterogeneous population that can be broken up into homogeneous groups, randomly sample from each group proportionate to their prevalence in the population.

Barbie vs. Oppenheimer

On July 21st, two highly anticipated movies arrive in theaters: **Barbie** and **Oppenheimer**. 3678970

We want to know which movie will prevail on opening day, in Berkeley.

Demo



[NY Times](#), [GQ](#)

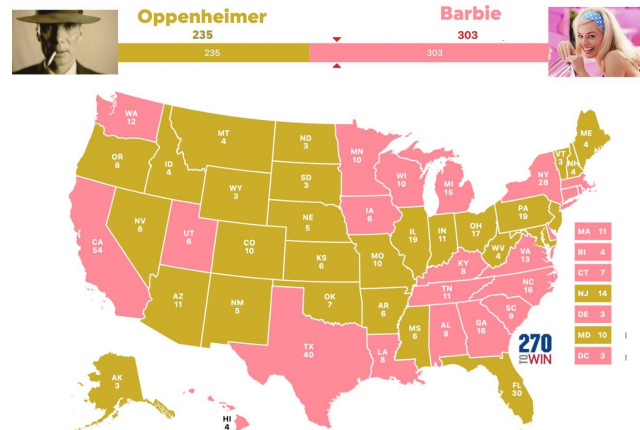
Imaginary Barbie Land Pollster

Suppose we took a sample of Berkeley residents to predict the ³⁶⁷⁸⁹⁷⁰ box office outcomes.

- We poll all **retirees** for their preference.
- Even if they answer truthfully, this is a **convenience sample**.

Then, suppose July 21st has passed (it has!).

- How “off” is our sample estimate from the actual outcome?
- How would a random sample with replacement have performed?



Answer: Barbie won!

- Opening weekend tally was:
 - **Barbie**: \$155M
 - **Oppenheimer**: \$82.4M

3678970

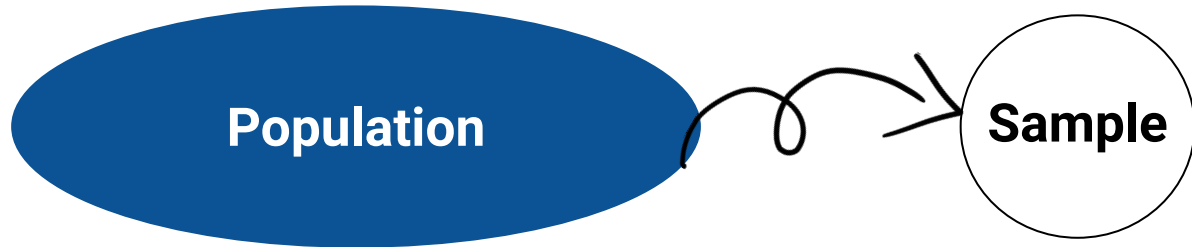


Demo

Multinomial Probabilities

CSCI 3022

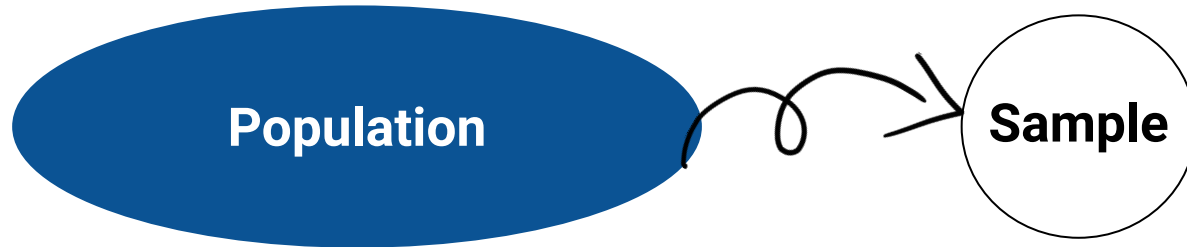
- Sampling: Definitions
- Sampling Bias: A Case Study
- Probability Samples
- **Multinomial Probabilities**



If we have a probability sample (aka a random sample):

- We can quantify error and bias.
- **Given the population distribution**, we can compute the probability of us getting a **particular sample**.

Note: We almost **never** know the population distribution! But this is a good start.



If we have a probability sample:

- We can quantify error and bias (to be covered later).
- **Given the population distribution**, we can compute the probability of us getting a **particular sample**.

Note: We almost **never** know the population distribution! But this is a good start.

Special case: Random sampling with replacement of a **categorical population** produces **Multinomial Probabilities**.

A very common approximation for sampling

A common situation in data science:

- We have an enormous population.
- We can only afford to sample a relatively small number of individuals.

As the **population gets very large** compared to the sample, then random sampling **with** replacement becomes a **good approximation** to RS **without**.

Example: Suppose there are 10,000 people in a population.
Exactly 7,500 of them like Reese's; the other 2,500 like Snickers.

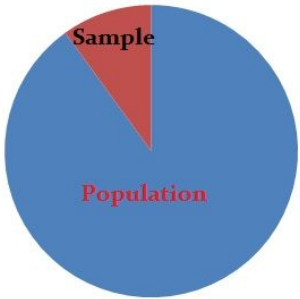
What is the probability that in a random sample of 20, **all people like Reese's**?

$$\text{SRS (Random Sample Without Replacement)} \quad \overbrace{\left(\frac{7500}{10000}\right)}^{0.75} \overbrace{\left(\frac{7499}{9999}\right)}^{0.74997} \cdots \overbrace{\left(\frac{7482}{9982}\right)}^{0.7495} \overbrace{\left(\frac{7481}{9981}\right)}^{0.7495} \approx .003151$$

$$\text{Random Sample With Replacement} \quad \left(\frac{7500}{10000}\right)^{20} \approx .003171$$

Probabilities of sampling with replacement are much easier to compute!

10% Rule



As the **population gets very large** compared to the sample, then random sampling **with** replacement becomes a **good approximation** to random sampling **without replacement**.

10% rule: If sample size <= 10% of population size:
Then we assume (approximate) independence (i.e. treat samples as if they were drawn randomly with replacement).

SRS (Random Sample Without Replacement) $\left(\frac{\overbrace{7500}^{0.75}}{10000}\right) \left(\frac{\overbrace{7499}^{0.74997}}{9999}\right) \cdots \left(\frac{\overbrace{7482}^{0.7495}}{9982}\right) \left(\frac{\overbrace{7481}^{0.7495}}{9981}\right) \approx .003151$ $20 < 0.10 \cdot (10000)$

Random Sample With Replacement $\left(\frac{7500}{10000}\right)^{20} \approx .003171$

Probabilities of sampling with replacement are much easier to compute!

Drawing samples from categorical distribution

Binomial and multinomial probabilities arise when we:

- Sample at random, **with replacement**.
- Sample a fixed number (n) times.
- Sample from a **categorical distribution**.
 - If 2 categories, **Binomial**:
Bag of marbles: 60% blue 40% not blue
 - If >2 categories, **Multinomial**:
Bag of marbles: 60% blue 30% green 10% red

Goal: **Count the number of each category** that end up in our sample.

- [`np.random.multinomial`](#) returns these counts.
- We'll derive the multinomial probabilities in this section as a review of probability.

Multinomial Probabilities

Suppose we get a sample of 4 **blue** marbles, 2 **green** marbles, and 1 **red** marble. What is the probability of getting this sample?

Q1. What is $P(\text{bgbbbgrr})$?

Use product rule to determine probability for a particular **order**:

$$P(\text{bgbbbgrr}) = 0.6 \times 0.3 \times 0.6 \times 0.6 \times 0.6 \times 0.3 \times 0.1 = (0.6)^4(0.3)^2(0.1)^1$$

Q2. What is $P(4 \text{ blue}, 2 \text{ green}, 1 \text{ red})$?

$$\frac{7!}{4! 2! 1!} (0.6)^4 (0.3)^2 (0.1)^1$$

multinomial probability

We use the **addition rule** and **multiplication rule**:

of ways to choose 4 of 7 places to write **b**, then choose 2 places to write **g**, (other 1 get filled with **r**)

For a particular outcome (say, Q1), probability of this **ordered series** of **b**'s, **g**'s, and **r**'s



Multinomial Probabilities: generalized

If we are drawing at random with replacement **n** times, from a population broken into three separate categories (where $p_1 + p_2 + p_3 = 1$):

- Category 1, with proportion **p₁** of the individuals.
- Category 2, with proportion **p₂** of the individuals.
- Category 3, with proportion **p₃** of the individuals.

Then, the **multinomial probability** of drawing **k₁** individuals from Category 1, **k₂** individuals from Category 2, and **k₃** individuals from Category 3 (where $k_1 + k_2 + k_3 = n$) is

$$\frac{n!}{k_1!k_2!k_3!} p_1^{k_1} p_2^{k_2} p_3^{k_3}$$

i.i.d. random variables

Consider n variables X_1, X_2, \dots, X_n .

X_1, X_2, \dots, X_n are **independent and identically distributed** if

- X_1, X_2, \dots, X_n are independent, and
- All have the same PMF (if discrete) or PDF (if continuous).
 - $\Rightarrow E[X_i] = \mu$ for $i = 1, \dots, n$
 - $\Rightarrow \text{Var}(X_i) = \sigma^2$ for $i = 1, \dots, n$

Same thing:

i.i.d.

iid

IID

Quick check

Are X_1, X_2, \dots, X_n i.i.d. with the following distributions?

1. $X_i \sim \text{Exp}(\lambda)$, X_i independent
2. $X_i \sim \text{Exp}(\lambda_i)$, X_i independent
3. $X_i \sim \text{Exp}(\lambda)$, $X_1 = X_2 = \dots = X_n$
4. $X_i \sim \text{Bin}(n_i, p)$, X_i independent

