# CAAP Statistics - Lec02

Jul 07, 2022

# Review

- From Chronic Fatigue Syndrome studies last time..

|  | | Good outcome | | |
|---|---|---|---|---|
|  | | Yes | No | Total |
| Group | Treatment | 19 | 8 | 27 |
|  | Control | 5 | 21 | 26 |
|  | Total | 24 | 29 | 53 |

- Proportion with good outcomes in treatment group

$$19/27 \approx 0.70 \rightarrow 70\%$$

- Proportion with good outcomes in control group

$$5/26 \approx 0.19 \rightarrow 19\%$$

# Learning Objectives

- Variable Types
- Observational Data
- What is sampling? (Sample vs Population)
- Experimental Data
  - Treatment group vs Control group
  - Randomization

# Data Basics

# Classroom survey

A survey was conducted on students in an introductory statistics course. Below are a few of the questions on the survey, and the corresponding variables the data from the responses were stored in:

- **gender**: What is your gender?
- **intro_extra**: Are you an introvert or an extrovert?
- **sleep**: How many hours do you sleep at night, on average?
- **bedtime**: What time do you usually go to bed?
- **countries**: How many countries have you visited?
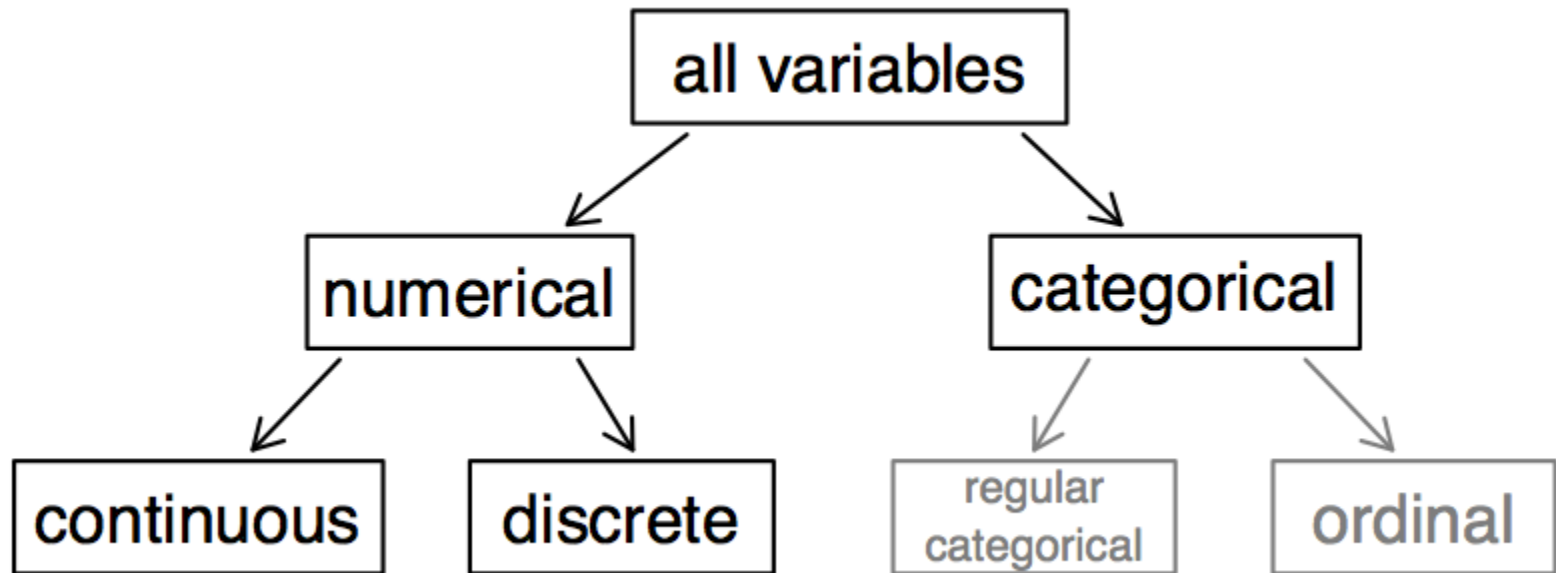- **dread**: On a scale of 1-5, how much do you dread being here?

# Data matrix

Data collected on students in a statistics class on a variety of variables:

| | *variable*<br>↓ | | | |
|------|--------|-------------|-----|-------|
| Stu. | gender | intro_extra | ⋯ | dread |
| 1 | male | extravert | ⋯ | 3 |
| 2 | female | extravert | ⋯ | 2 |
| 3 | female | introvert | ⋯ | 4 |
| 4 | female | extravert | ⋯ | 2 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| 86 | male | extravert | ⋯ | 3 |

← *observation*

# Types of variables

# Types of variables (cont.)

|   | gender | sleep | bedtime | countries | dread |
|---|--------|-------|---------|-----------|-------|
| 1 | male   | 5     | 12-2    | 13        | 3     |
| 2 | female | 7     | 10-12   | 7         | 2     |
| 3 | female | 5.5   | 12-2    | 1         | 4     |
| 4 | female | 7     | 12-2    |           | 2     |
| 5 | female | 3     | 12-2    | 1         | 3     |
| 6 | female | 3     | 12-2    | 9         | 4     |

- gender:

# Types of variables (cont.)

| | gender | sleep | bedtime | countries | dread |
|---|--------|-------|---------|-----------|-------|
| 1 | male | 5 | 12-2 | 13 | 3 |
| 2 | female | 7 | 10-12 | 7 | 2 |
| 3 | female | 5.5 | 12-2 | 1 | 4 |
| 4 | female | 7 | 12-2 | | 2 |
| 5 | female | 3 | 12-2 | 1 | 3 |
| 6 | female | 3 | 12-2 | 9 | 4 |

- gender: *categorical*

# Types of variables (cont.)

| | gender | sleep | bedtime | countries | dread |
|---|--------|-------|---------|-----------|-------|
| 1 | male | 5 | 12-2 | 13 | 3 |
| 2 | female | 7 | 10-12 | 7 | 2 |
| 3 | female | 5.5 | 12-2 | 1 | 4 |
| 4 | female | 7 | 12-2 | | 2 |
| 5 | female | 3 | 12-2 | 1 | 3 |
| 6 | female | 3 | 12-2 | 9 | 4 |

- gender: *categorical*
- sleep:

# Types of variables (cont.)

| | gender | sleep | bedtime | countries | dread |
|---|---|---|---|---|---|
| 1 | male | 5 | 12-2 | 13 | 3 |
| 2 | female | 7 | 10-12 | 7 | 2 |
| 3 | female | 5.5 | 12-2 | 1 | 4 |
| 4 | female | 7 | 12-2 | | 2 |
| 5 | female | 3 | 12-2 | 1 | 3 |
| 6 | female | 3 | 12-2 | 9 | 4 |

- gender: *categorical*
- sleep: *numerical, continuous*

# Types of variables (cont.)

| | gender | sleep | bedtime | countries | dread |
|---|---|---|---|---|---|
| 1 | male | 5 | 12-2 | 13 | 3 |
| 2 | female | 7 | 10-12 | 7 | 2 |
| 3 | female | 5.5 | 12-2 | 1 | 4 |
| 4 | female | 7 | 12-2 | | 2 |
| 5 | female | 3 | 12-2 | 1 | 3 |
| 6 | female | 3 | 12-2 | 9 | 4 |

- gender: *categorical*
- sleep: *numerical, continuous*
- bedtime:

# Types of variables (cont.)

| | gender | sleep | bedtime | countries | dread |
|---|---|---|---|---|---|
| 1 | male | 5 | 12-2 | 13 | 3 |
| 2 | female | 7 | 10-12 | 7 | 2 |
| 3 | female | 5.5 | 12-2 | 1 | 4 |
| 4 | female | 7 | 12-2 | | 2 |
| 5 | female | 3 | 12-2 | 1 | 3 |
| 6 | female | 3 | 12-2 | 9 | 4 |

- gender: *categorical*
- sleep: *numerical, continuous*
- bedtime: *categorical, ordinal*

# Types of variables (cont.)

| | gender | sleep | bedtime | countries | dread |
|---|--------|-------|---------|-----------|-------|
| 1 | male | 5 | 12-2 | 13 | 3 |
| 2 | female | 7 | 10-12 | 7 | 2 |
| 3 | female | 5.5 | 12-2 | 1 | 4 |
| 4 | female | 7 | 12-2 | | 2 |
| 5 | female | 3 | 12-2 | 1 | 3 |
| 6 | female | 3 | 12-2 | 9 | 4 |

- gender: *categorical*
- sleep: *numerical, continuous*
- bedtime: *categorical, ordinal*
- countries:

# Types of variables (cont.)

| | gender | sleep | bedtime | countries | dread |
|---|--------|-------|---------|-----------|-------|
| 1 | male   | 5     | 12-2    | 13        | 3     |
| 2 | female | 7     | 10-12   | 7         | 2     |
| 3 | female | 5.5   | 12-2    | 1         | 4     |
| 4 | female | 7     | 12-2    |           | 2     |
| 5 | female | 3     | 12-2    | 1         | 3     |
| 6 | female | 3     | 12-2    | 9         | 4     |

- **gender**: *categorical*
- **sleep**: *numerical, continuous*
- **bedtime**: *categorical, ordinal*
- **countries**:  numerical, discrete

# Types of variables (cont.)

| | gender | sleep | bedtime | countries | dread |
|---|--------|-------|---------|-----------|-------|
| 1 | male | 5 | 12-2 | 13 | 3 |
| 2 | female | 7 | 10-12 | 7 | 2 |
| 3 | female | 5.5 | 12-2 | 1 | 4 |
| 4 | female | 7 | 12-2 | | 2 |
| 5 | female | 3 | 12-2 | 1 | 3 |
| 6 | female | 3 | 12-2 | 9 | 4 |

- gender: *categorical*
- sleep: *numerical, continuous*
- bedtime: *categorical, ordinal*
- countries: *numerical, discrete*
- dread:

# Types of variables (cont.)

| | gender | sleep | bedtime | countries | dread |
|---|---|---|---|---|---|
| 1 | male | 5 | 12-2 | 13 | 3 |
| 2 | female | 7 | 10-12 | 7 | 2 |
| 3 | female | 5.5 | 12-2 | 1 | 4 |
| 4 | female | 7 | 12-2 | | 2 |
| 5 | female | 3 | 12-2 | 1 | 3 |
| 6 | female | 3 | 12-2 | 9 | 4 |

- gender: *categorical*
- sleep: *numerical, continuous*
- bedtime: *categorical, ordinal*
- countries: *numerical, discrete*
- dread: *categorical, ordinal*

# Practice

What type of variable is a telephone area code?

(a) numerical, continuous
(b) numerical, discrete
(c) categorical
(d) categorical, ordinal

# Practice

What type of variable is a telephone area code?

(a) numerical, continuous

(b) numerical, discrete

(c) *categorical*

(d) categorical, ordinal

# Relationships among variables

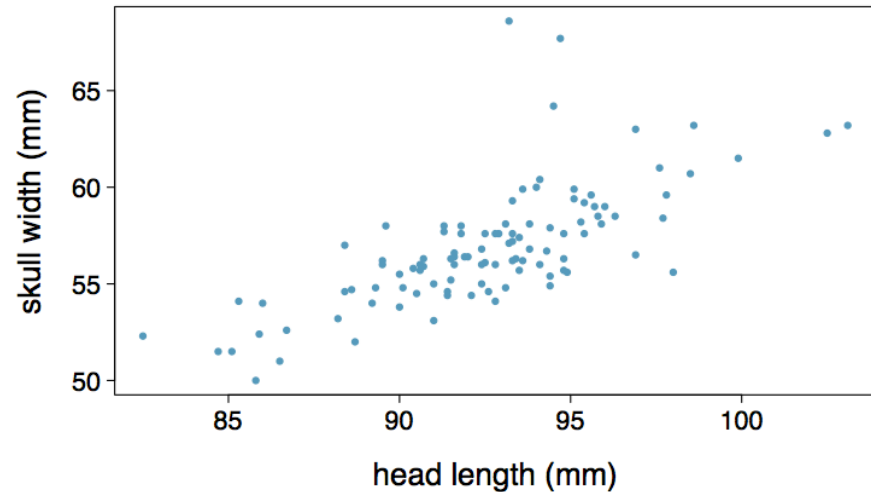Does there appear to be a relationship between the hours of study per week and the GPA of a student?

# Relationships among variables

Does there appear to be a relationship between the hours of study per week and the GPA of a student?



Can you spot anything unusual about any of the data points?

# Relationships among variables

Does there appear to be a relationship between the hours of study per week and the GPA of a student?



Can you spot anything unusual about any of the data points?

There is one student with *GPA* > 4.0, this is likely a data error.

# Associated vs. independent

- When two variables show some connection with one another, they are called *associated* variables.
  - Associated variables can also be called *dependent* variables and vice-versa.
- If two variables are not associated, i.e. there is no evident connection between the two, then they are said to be *independent*.

# Practice

Based on the scatterplot on the right, which of the following statements is correct about the head and skull lengths of possums?



(a) There is no relationship between head length and skull width, i.e. the variables are independent.
(b) Head length and skull width are positively associated.
(c) Skull width and head length are negatively associated.
(d) A longer head causes the skull to be wider.
(e) A wider skull causes the head to be longer.

# Practice

Bas[...]
the [...]
follo[...]
corr[...]
sku[...]

(a) [...]

(b) Head length and skull width are positively associated.
(c) Skull width and head length are negatively associated.
(d) A longer head causes the skull to be wider.
(e) A wider skull causes the head to be longer.

# Practice

Based on the scatterplot on the right, which of the following statements is correct about the head and skull lengths of possums?



(a) There is no relationship between head length and skull width, i.e. the variables are independent.
(b) *Head length and skull width are positively associated.*
(c) Skull width and head length are negatively associated.
(d) A longer head causes the skull to be wider.
(e) A wider skull causes the head to be longer.

# Observational studies and sampling strategies

# Populations and Samples



PHYS ED | AUGUST 29, 2012, 12:01 AM | 21 Comments

**Finding Your Ideal Running Form**

By GRETCHEN REYNOLDS

David De Lossy/Getty Images

http://well.blogs.nytimes.com/2012/08/29/finding-your-ideal-running-form

*Research Question*: Can people become better, more efficient runners on their own, merely by running?

# Populations and Samples



PHYS ED | AUGUST 29, 2012, 12:01 AM | 21 Comments

**Finding Your Ideal Running Form**
By GRETCHEN REYNOLDS

David De Lossy/Getty Images

http://well.blogs.nytimes.com/2012/08/29/finding-your-ideal-running-form

*Research Question*: Can people become better, more efficient runners on their own, merely by running?

*Population of Interest*: All people

# Populations and Samples

## Finding Your Ideal Running Form
By GRETCHEN REYNOLDS



David De Lossy/Getty Images

http://well.blogs.nytimes.com/2012/08/29/finding-your-ideal-running-form

*Research Question*: Can people become better, more efficient runners on their own, merely by running?

*Population of Interest*: All people

*Sample*:  Group of adult women who recently joined a running group

# Populations and Samples



**PHYS ED** | AUGUST 29, 2012, 12:01 AM | 21 Comments

## Finding Your Ideal Running Form

By GRETCHEN REYNOLDS

David De Lossy/Getty Images

http://well.blogs.nytimes.com/2012/08/29/finding-your-ideal-running-form

*Research Question*: Can people become better, more efficient runners on their own, merely by running?

*Population of Interest*: All people

*Sample*:  Group of adult women who recently joined a running group

*Population to which results can be generalized*:  Adult women, if the data are randomly sampled

# Census

- Wouldn't it be better to just include everyone and "sample" the entire population?
    - This is called a *census*.

# Census

- Wouldn't it be better to just include everyone and "sample" the entire population?
    - This is called a *census*.

- There are problems with taking a census:
    - It can be difficult to complete a census: there always seem to be some individuals who are hard to locate or hard to measure. *And these difficult-to-find people may have certain characteristics that distinguish them from the rest of the population.*
    - Populations rarely stand still. Even if you could take a census, the population changes constantly, so it's never possible to get a perfect measure.
    - Taking a census may be more complex than sampling.

# Anecdotal evidence and early smoking research

- Anti-smoking research started in the 1930s and 1940s when cigarette smoking became increasingly popular. While some smokers seemed to be sensitive to cigarette smoke, others were completely unaffected.

- Anti-smoking research was faced with resistance based on anecdotal evidence such as "My uncle smokes three packs a day and he's in perfectly good health", evidence based on a limited sample size that might not be representative of the population.

- It was concluded that "smoking is a complex human behavior, by its nature difficult to study, confounded by human variability."

- In time researchers were able to examine larger samples of cases (smokers), and trends showing that smoking has negative health impacts became much clearer.

Brandt, **The Cigarette Century** (2009), Basic Books.

# Exploratory analysis to inference

- Sampling is natural.

# Exploratory analysis to inference

- Sampling is natural.

- Think about sampling something you are cooking - you taste (examine) a small part of what you're cooking to get an idea about the dish as a whole.

# Exploratory analysis to inference

- Sampling is natural.

- Think about sampling something you are cooking - you taste (examine) a small part of what you're cooking to get an idea about the dish as a whole.

- When you taste a spoonful of soup and decide the spoonful you tasted isn't salty enough, that's *exploratory analysis*.

# Exploratory analysis to inference

- Sampling is natural.

- Think about sampling something you are cooking - you taste (examine) a small part of what you're cooking to get an idea about the dish as a whole.

- When you taste a spoonful of soup and decide the spoonful you tasted isn't salty enough, that's *exploratory analysis*.

- If you generalize and conclude that your entire soup needs salt, that's an *inference*.

# Exploratory analysis to inference

- Sampling is natural.

- Think about sampling something you are cooking - you taste (examine) a small part of what you're cooking to get an idea about the dish as a whole.

- When you taste a spoonful of soup and decide the spoonful you tasted isn't salty enough, that's *exploratory analysis*.

- If you generalize and conclude that your entire soup needs salt, that's an *inference*.

- For your inference to be valid, the spoonful you tasted (the sample) needs to be *representative* of the entire pot (the population).
  - If your spoonful comes only from the surface and the salt is collected at the bottom of the pot, what you tasted is probably not representative of the whole pot.
  - If you first stir the soup thoroughly before you taste, your spoonful will more likely be representative of the whole pot.

# Sampling bias

- *Non-response*: If only a small fraction of the randomly sampled people choose to respond to a survey, the sample may no longer be representative of the population.

# Sampling bias

- *Non-response*: If only a small fraction of the randomly sampled people choose to respond to a survey, the sample may no longer be representative of the population.
- *Voluntary response*: Occurs when the sample consists of people who volunteer to respond because they have strong opinions on the issue. Such a sample will also not be representative of the population.

# Sampling bias

- **Non-response:** If only a small fraction of the randomly sampled people choose to respond to a survey, the sample may no longer be representative of the population.
- **Voluntary response:** Occurs when the sample consists of people who volunteer to respond because they have strong opinions on the issue. Such a sample will also not be representative of the population.
- **Convenience sample:** Individuals who are easily accessible are more likely to be included in the sample.

# Sampling bias example:
## Landon vs. Franklin D. Roosevelt

A historical example of a biased sample yielding misleading results



In 1936, Landon sought the Republican presidential nomination opposing the re-election of FDR.

# The Literary Digest Poll

- The Literary Digest polled about 10 million Americans, and got responses from about 2.4 million.
- The poll showed that Landon would likely be the overwhelming winner and FDR would get only 43% of the votes.
- **Election result:  FDR won, with 62% of the votes.**
- The magazine was completely discredited because of the poll, and was soon discontinued.

# The Literary Digest Poll - what went wrong?

- The magazine had surveyed
  - its own readers,
  - registered automobile owners, and
  - registered telephone users.

- These groups had incomes well above the national average of the day (remember, this is Great Depression era) which resulted in lists of voters far more likely to support Republicans than a truly *typical* voter of the time, i.e. the sample was not representative of the American population at the time.

# Large samples are preferable, but...

- The Literary Digest election poll was based on a sample size of 2.4 million, which is huge, but since the sample was *biased*, the sample did not yield an accurate prediction.

- Back to the soup analogy: If the soup is not well stirred, it doesn't matter how large a spoon you have, it will still not taste right. If the soup is well stirred, a small spoon will suffice to test the soup.

# Practice

A school district is considering whether it will no longer allow high school students to park at school after two recent accidents where students were severely injured. As a first step, they survey parents by mail, asking them whether or not the parents would object to this policy change. Of 6,000 surveys that go out, 1,200 are returned. Of these 1,200 surveys that were completed, 960 agreed with the policy change and 240 disagreed. Which of the following statements are true?

I.   Some of the mailings may have never reached the parents.

II.  The school district has strong support from parents to move forward with the policy approval.

III. It is possible that majority of the parents of high school students disagree with the policy change.

IV.  The survey results are unlikely to be biased because all parents were mailed a survey.

(a) Only I               (b) I and II               (c) I and III               (d) III and IV               (e) Only IV

# Practice

A school district is considering whether it will no longer allow high school students to park at school after two recent accidents where students were severely injured. As a first step, they survey parents by mail, asking them whether or not the parents would object to this policy change. Of 6,000 surveys that go out, 1,200 are returned. Of these 1,200 surveys that were completed, 960 agreed with the policy change and 240 disagreed. Which of the following statements are true?

I. Some of the mailings may have never reached the parents.

II. The school district has strong support from parents to move forward with the policy approval.

III. It is possible that majority of the parents of high school students disagree with the policy change.

IV. The survey results are unlikely to be biased because all parents were mailed a survey.

(a) Only I       (b) I and II       (c) I and III       (d) III and IV       (e) Only IV
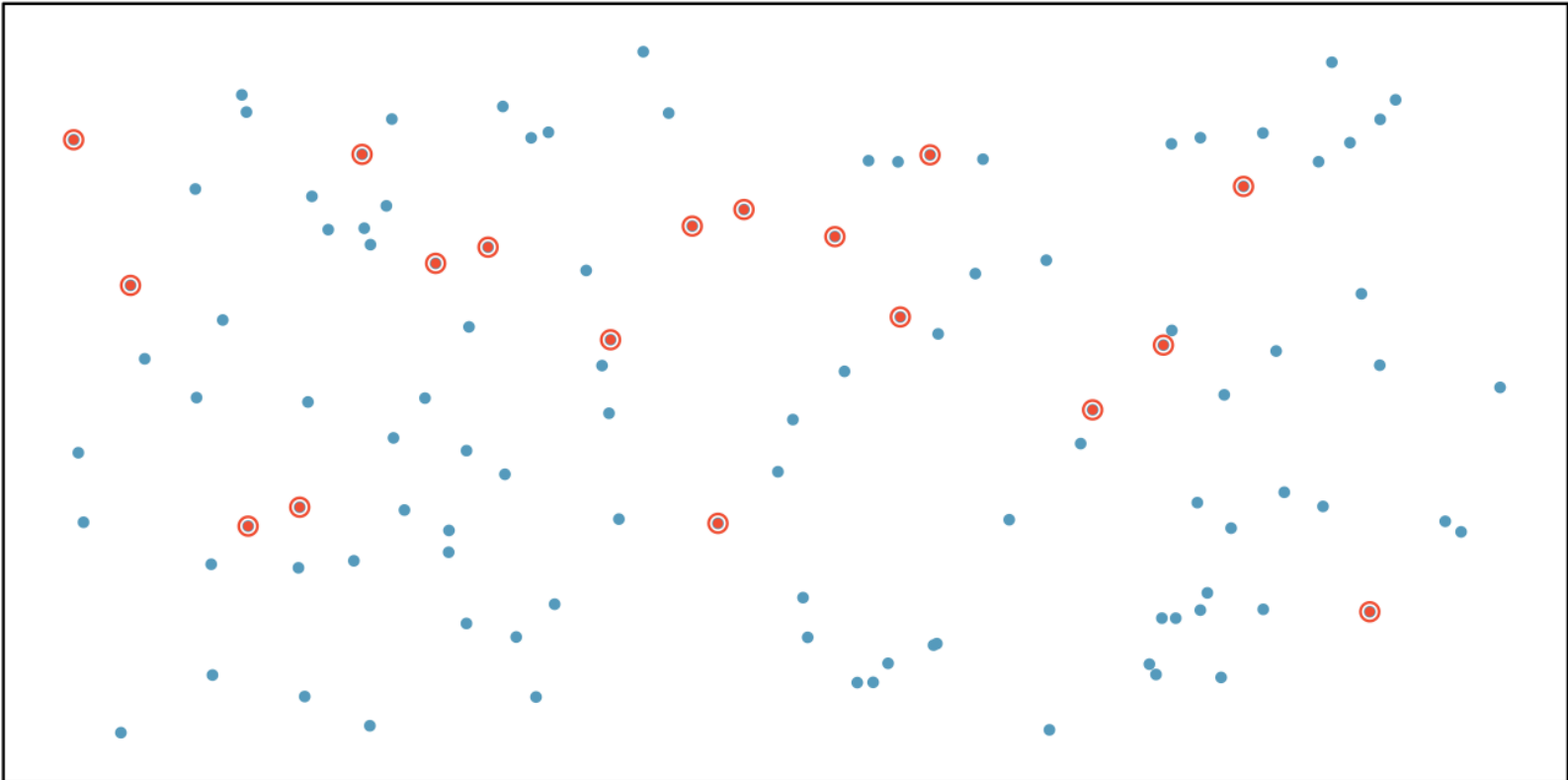
# Observational studies

- Researchers collect data in a way that does not directly interfere with how the data arise.

- Results of an observational study can generally be used to establish an association between the explanatory and response variables.

# Obtaining Good Samples

- Almost all statistical methods are based on the notion of implied randomness.

- If observational data are not collected in a random framework from a population, these statistical methods -- the estimates and errors associated with the estimates -- are not reliable.

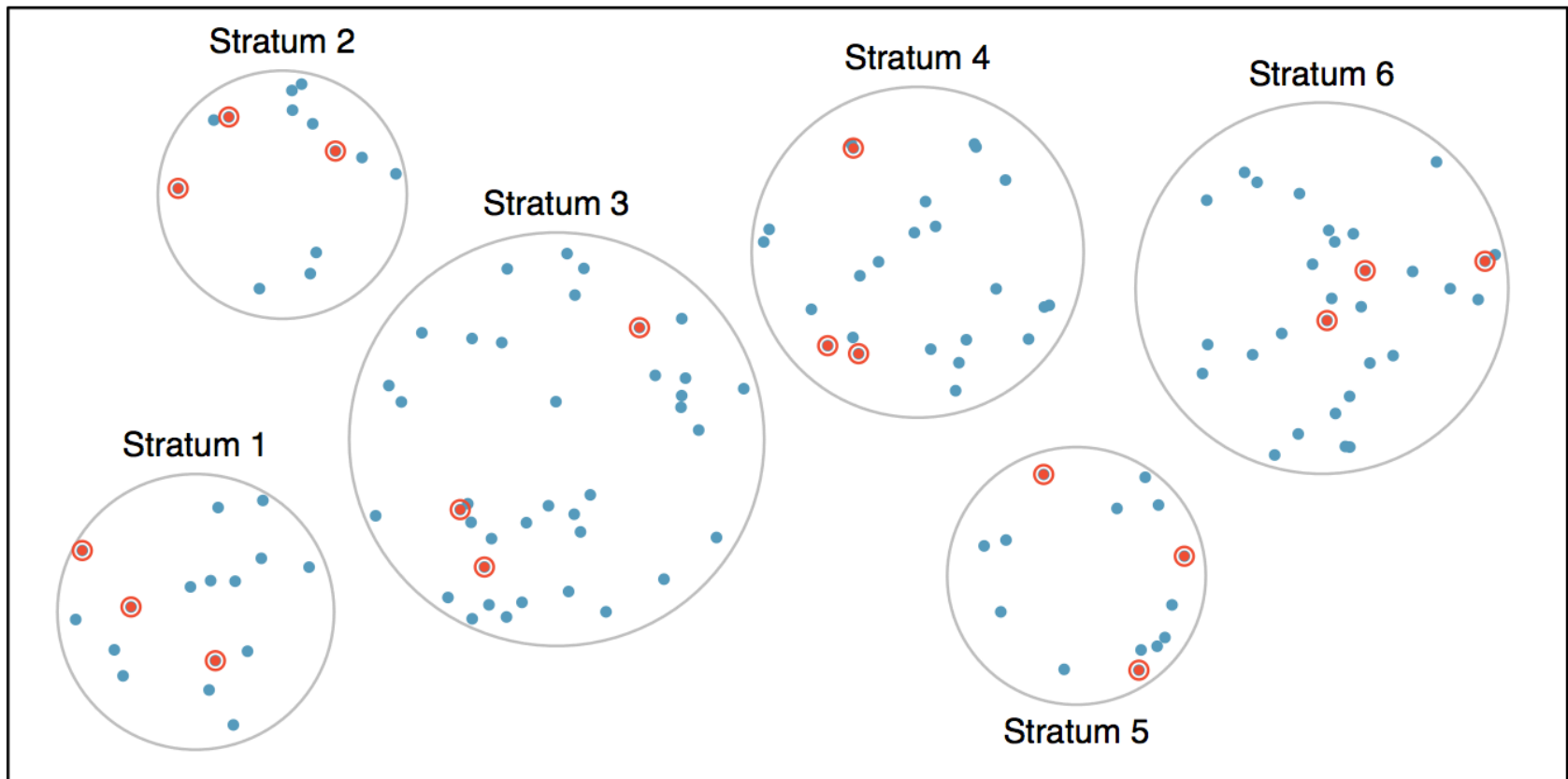- Most commonly used random sampling techniques are *simple*, *stratified*, and *cluster* sampling.

# Simple Random Sample

Randomly select cases from the population, where there is no implied connection between the points that are selected.
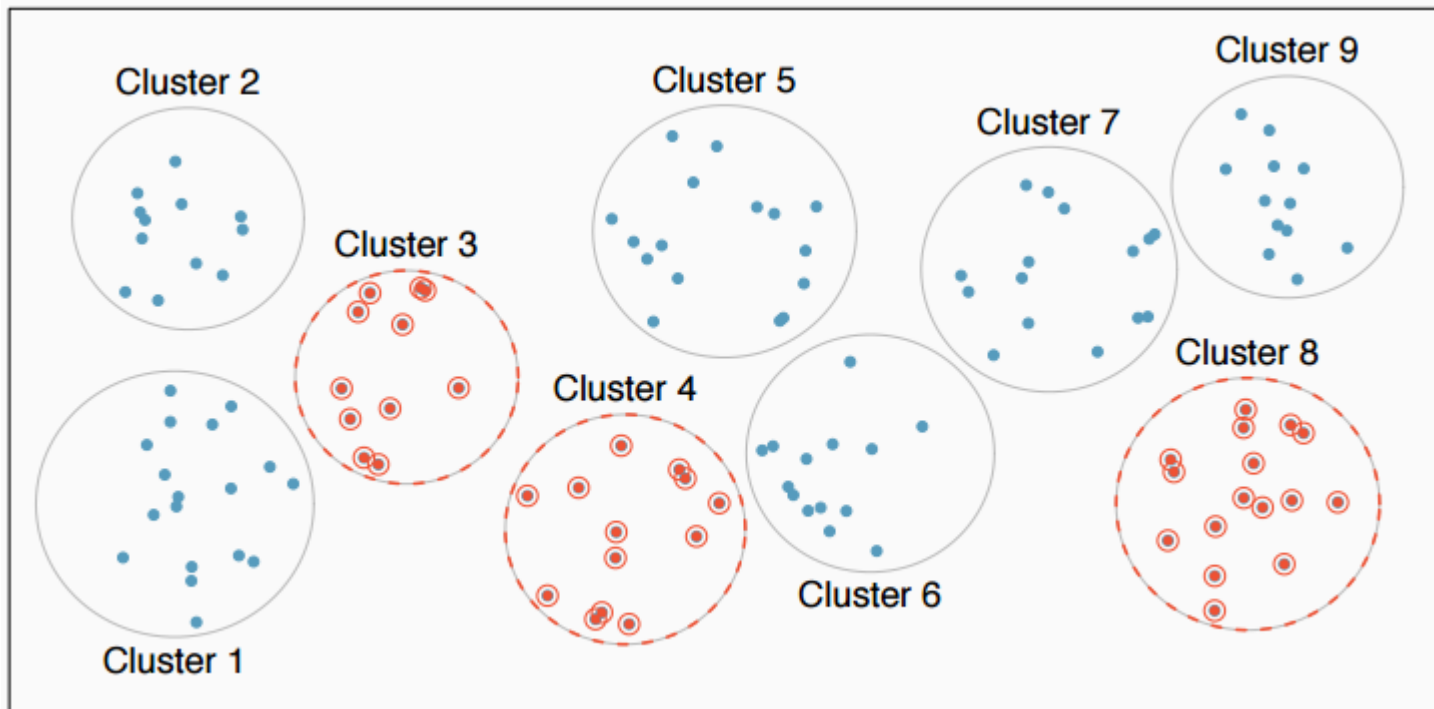
# Stratified Sample

*Strata* are made up of similar observations. We take a simple random sample from <u>each</u> stratum.
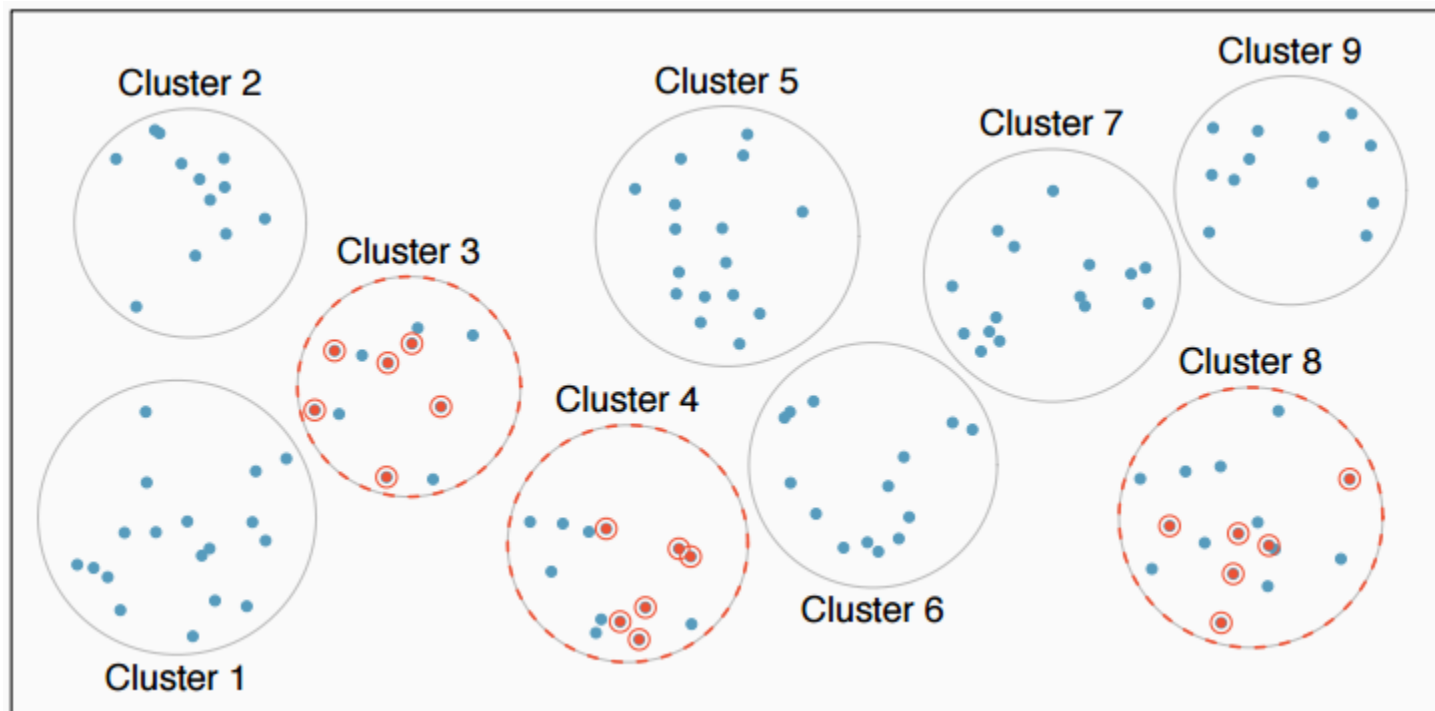
# Cluster Sample

*Clusters* are usually not made up of homogeneous observations. We take a simple random sample of clusters, and then sample all observations in that cluster. Usually preferred for economical reasons.

# Multistage Sample

*Clusters* are usually not made up of homogeneous observations. We take a simple random sample of clusters, and then take a simple random sample of observations from the sampled clusters

# Practice

A city council has requested a household survey be conducted in a suburban area of their city. The area is broken into many distinct and unique neighborhoods, some including large homes, some with only apartments. Which approach would likely be the *least* effective?

(a) Simple random sampling
(b) Cluster sampling
(c) Stratified sampling
(d) Blocked sampling

# Practice

A city council has requested a household survey be conducted in a suburban area of their city. The area is broken into many distinct and unique neighborhoods, some including large homes, some with only apartments. Which approach would likely be the *least* effective?

(a) Simple random sampling

(b) *Cluster sampling*

(c) Stratified sampling

(d) Blocked sampling

# Experiments

# Principles of experimental design

1.  **Control**: Compare treatment of interest to a control group.
2.  **Randomize**: Randomly assign subjects to treatments, and randomly sample from the population whenever possible.
3.  **Replicate**: Within a study, replicate by collecting a sufficiently large sample. Or replicate the entire study.
4.  **Block**: If there are variables that are known or suspected to affect the response variable, first group subjects into blocks based on these variables, and then randomize cases within each block to treatment groups.

# More on Blocking



- We would like to design an experiment to investigate if energy gels makes you run faster:

# More on Blocking



- We would like to design an experiment to investigate if energy gels makes you run faster:
  - Treatment: energy gel
  - Control: no energy gel

# More on Blocking

- We would like to design an experiment to investigate if energy gels makes you run faster:
  - Treatment: energy gel
  - Control: no energy gel

- It is suspected that energy gels might affect pro and amateur athletes differently, therefore we block for pro status:

# More on Blocking



- We would like to design an experiment to investigate if energy gels makes you run faster:
  - Treatment: energy gel
  - Control: no energy gel

- It is suspected that energy gels might affect pro and amateur athletes differently, therefore we block for pro status:
  - Divide the sample to pro and amateur
  - Randomly assign pro athletes to treatment and control groups
  - Randomly assign amateur athletes to treatment and control groups
  - Pro/amateur status is equally represented in the resulting treatment and control groups

# More on Blocking



- We would like to design an experiment to investigate if energy gels makes you run faster:

  ○ Treatment: energy gel
  ○ Control: no energy gel

- It is suspected that energy gels might affect pro and amateur athletes differently, therefore we block for pro status:

  ○ Divide the sample to pro and amateur
  ○ Randomly assign pro athletes to treatment and control groups
  ○ Randomly assign amateur athletes to treatment and control groups
  ○ Pro/amateur status is equally represented in the resulting treatment and control groups

Why is this important? Can you think of other variables to block for?

# Practice

A study is designed to test the effect of light level and noise level on exam performance of students. The researcher also believes that light and noise levels might have different effects on males and females, so wants to make sure both genders are equally represented in each group. Which of the below is correct?

A. There are 3 explanatory variables (light, noise, gender) and 1 response variable (exam performance)

B. There are 2 explanatory variables (light and noise), 1 blocking variable (gender), and 1 response variable (exam performance)

C. There is 1 explanatory variable (gender) and 3 response variables (light, noise, exam performance)

D. There are 2 blocking variables (light and noise), 1 explanatory variable (gender), and 1 response variable (exam performance)

# Practice

A study is designed to test the effect of light level and noise level on exam performance of students. The researcher also believes that light and noise levels might have different effects on males and females, so wants to make sure both genders are equally represented in each group. Which of the below is correct?

A. There are 3 explanatory variables (light, noise, gender) and 1 response variable (exam performance)

B. *There are 2 explanatory variables (light and noise), 1 blocking variable (gender), and 1 response variable (exam performance)*

C. There is 1 explanatory variable (gender) and 3 response variables (light, noise, exam performance)

D. There are 2 blocking variables (light and noise), 1 explanatory variable (gender), and 1 response variable (exam performance)

# Difference Between Blocking and Explanatory Variables

- Factors are conditions we can impose on the experimental units.

- Blocking variables are characteristics that the experimental units come with, that we would like to control for.

- Blocking is like stratifying, except used in experimental settings when randomly assigning, as opposed to when sampling.

# More Experimental Design Terminology...

- Placebo: fake treatment, often used as the control group for medical studies

- Placebo effect: experimental units showing improvement simply because they believe they are receiving a special treatment

- Blinding: when experimental units do not know whether they are in the control or treatment group

- Double-blind: when both the experimental units and the researchers who interact with the patients do not know who is in the control and who is in the treatment group

# Practice

What is the main difference between observational studies and experiments?

A. Experiments take place in a lab while observational studies do not need to.

B. In an observational study we only look at what happened in the past.

C. Most experiments use random assignment while observational studies do not.

D. Observational studies are completely useless since no causal inference can be made based on their findings.

# Practice

What is the main difference between observational studies and experiments?

A. Experiments take place in a lab while observational studies do not need to.

B. In an observational study we only look at what happened in the past.

C. *Most experiments use random assignment while observational studies do not.*

D. Observational studies are completely useless since no causal inference can be made based on their findings.

# Random Assignment vs. Random Sampling



|  | Random assignment | No random assignment |  |
|---|---|---|---|
| *ideal experiment* |  |  | *most observational studies* |
| **Random sampling** | Causal conclusion, generalized to the whole population. | No causal conclusion, correlation statement generalized to the whole population. | **Generalizability** |
| **No random sampling** | Causal conclusion, only for the sample. | No causal conclusion, correlation statement only for the sample. | **No generalizability** |
| *most experiments* | Causation | Correlation | *bad observational studies* |

# Let's discuss!

# Chia Seeds and Weight Loss

Chia Pets – those terra-cotta figurines that sprout fuzzy green hair – made the chia plant a household name. But chia has gained an entirely new reputation as a diet supplement. In one 2009 study, a team of researchers recruited 38 men and divided them randomly into two groups: treatment or control. They also recruited 38 women, and they randomly placed half of these participants into the treatment group and the other half into the control group. One group was given 25 grams of chia seeds twice a day, and the other was given a placebo. The subjects volunteered to be a part of the study. After 12 weeks, the scientists found no significant difference between the groups in appetite or weight loss.

(a) What type of study is this?

(b) What are the experimental and control treatments in this study?

(c) Has blocking been used in this study? If so, what is the blocking variable?

(d) Has blinding been used in this study?

# Eat better, feel better?

In a public health study on the effects of consumption of fruits and vegetables on psychological well-being in young adults, participants were randomly assigned to three groups: (1) diet- as-usual, (2) an ecological momentary intervention involving text message reminders to increase their fruits and vegetable consumption plus a voucher to purchase them, or (3) a fruit and vegetable intervention in which participants were given two additional daily servings of fresh fruits and vegetables to consume on top of their normal diet. Participants were asked to take a nightly survey on their smartphones. Participants were student volunteers at the University of Otago, New Zealand. At the end of the 14-day study, only participants in the third group showed improvements to their psychological well-being across the 14-days relative to the other groups.

(a)  What type of study is this?
(b)  Identify the explanatory and response variables.
(c)  Comment on whether the results of the study can be generalized to the population.
(d)  Comment on whether the results of the study can be used to establish causal relationships.
(e)  A newspaper article reporting on the study states, "The results of this study provide proof that giving young adults fresh fruits and vegetables to eat can have psychological benefits, even over a brief period of time." How would you suggest revising this statement so that it can be supported by the study?