

# Chapter 1: Introduction to data

Wen-Han Hwang

( Slides primarily developed by Mine Çetinkaya-Rundel from OpenIntro.)



Institute of Statistics  
National Tsing Hua University



# Chapter 1: Introduction to data

Wen-Han Hwang  
( Slides primarily developed by Mine Çetinkaya-Rundel from OpenIntro.)



# Outline

- 1 Case study
- 2 Data basics
- 3 Sampling principles and strategies
- 4 Experiments

# Treating Chronic Fatigue Syndrome

- Objective: Evaluate the effectiveness of cognitive-behavior therapy for chronic fatigue syndrome.
- Participant pool: 142 patients who were recruited from referrals by primary care physicians and consultants to a hospital clinic specializing in chronic fatigue syndrome.
- Actual participants: Only 60 of the 142 referred patients entered the study. Some were excluded because they didn't meet the diagnostic criteria, some had other health issues, and some refused to be a part of the study.

Deale et. al. *Cognitive behavior therapy for chronic fatigue syndrome: A randomized controlled trial*. The American Journal of Psychiatry 154.3 (1997).

- Patients randomly assigned to treatment and control groups, 30 patients in each group:
  - **Treatment:** Cognitive behavior therapy – collaborative, educative, and with a behavioral emphasis. Patients were shown on how activity could be increased steadily and safely without exacerbating symptoms.
  - **Control:** Relaxation – No advice was given about how activity could be increased. Instead progressive muscle relaxation, visualization, and rapid relaxation skills were taught.

# Results

The table below shows the distribution of patients with good outcomes at 6-month follow-up. Note that 7 patients dropped out of the study: 3 from the treatment and 4 from the control group.

	<i>Good outcome</i>		Total
	Yes	No	
<i>Group</i>			
Treatment	19	8	27
Control	5 <sub>15</sub>	21 <sub>11</sub>	26
Total	24	29	53

- Proportion with good outcomes in treatment group:

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

- Proportion with good outcomes in the control group:

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

$$\frac{15}{26} = 58\%$$

# Understanding the results

Do the data show a “real” difference between the groups?

- Suppose you flip a coin 100 times. While the chance a coin lands heads in any given coin flip is 50%, we probably won't observe exactly 50 heads. This type of fluctuation is part of almost any type of data generating process.
- The observed difference between the two groups ( $70 - 19 = 51\%$ ) may be real, or may be due to natural variation.
- Since the difference is quite large, it is more believable that the difference is real.
- We need statistical tools to determine if the difference is so large that we should reject the notion that it was due to chance.

# Generalizing the results

Are the results of this study generalizable to all patients with chronic fatigue syndrome?

These patients had specific characteristics and volunteered to be a part of this study, therefore they may not be representative of all patients with chronic fatigue syndrome. While we cannot immediately generalize the results to all patients, this first study is encouraging. The method works for patients with some narrow set of characteristics, and that gives hope that it will work, at least to some degree, with other patients.



# 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: Do you consider yourself introverted or extraverted?
- 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

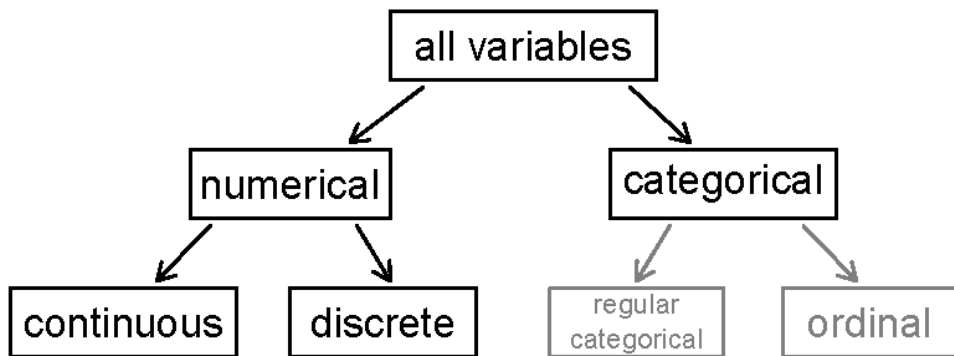
Data collected on students in a statistics class on a variety of variables:

variable  
↓

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.)

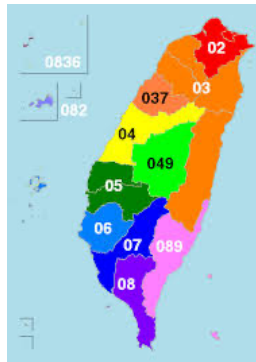
	gender	sleep	bedtime	countries	dread
1	male	5	12-2	13	3
2	female	7	10-12	7	2
3	female	5.5	12-2	1	4
4	female	7	12-2	3	2
5	female	3	12-2	1	3
6	female	3	12-2	9	4

- gender: categorical
- sleep: numerical, continuous
- bedtime: categorical, ordinal
- countries: numerical, discrete
- dread: categorical, ordinal - could also be used as numerical

# Practice

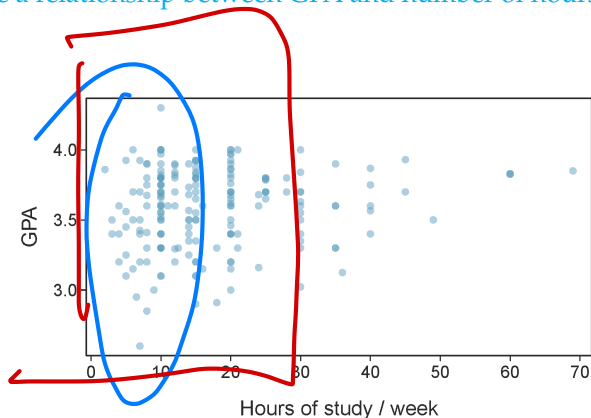
What type of variable is a telephone area code?

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



# Relationships among variables

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



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. Hours = 70

# Explanatory and response variables

$$f(x) = y$$

- To identify the explanatory variable in a pair of variables, identify which of the two is suspected of affecting the other:

explanatory variable  $\xrightarrow{\text{might affect}}$  response variable

- Labeling variables as explanatory and response does not guarantee the relationship between the two is actually causal, even if there is an association identified between the two variables. We use these labels only to keep track of which variable we suspect affects the other.
- Sometimes, the explanatory variable is called the independent variable, and the response variable is called the dependent variable.

$$y = f(x)$$

# Two primary types of data collection

- **Observational studies:** Collect data in a way that does not directly interfere with how the data arise (e.g. surveys).
  - Can provide evidence of a naturally occurring association between variables, but they cannot by themselves show a causal connection.
- **Experiment:** Researchers randomly assign subjects to various treatments in order to establish causal connections between the explanatory and response variables.



### **### Observational Studies: '**

#### **Longitudinal Cohort Study on Smoking and Lung Cancer\*\***

**- Researchers follow two groups of people for several years: one group of smokers and one group of non-smokers. They collect data on various health outcomes, including the incidence of lung cancer, without intervening. This study can provide evidence of an association between smoking and lung cancer but cannot definitively prove causation due to potential confounding factors.**

### **### Experiments:**

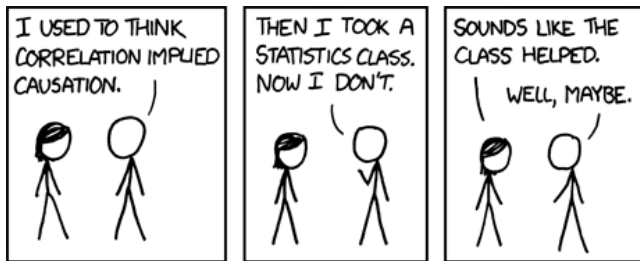
#### **Randomized Controlled Trial (RCT) on a New Drug for High Blood Pressure\*\***

**- Participants are randomly assigned to either receive the new drug or a placebo. Neither the participants nor the researchers know who is receiving the actual drug versus the placebo (double-blind). This setup allows researchers to directly measure the drug's effect on blood pressure while controlling for other variables.**

**Observational studies are excellent for identifying associations and generating hypotheses, while experiments are the gold standard for testing causality.**

# Association vs. causation

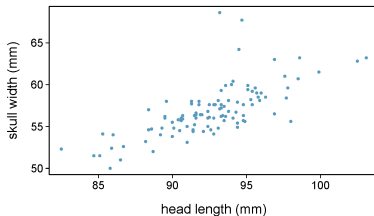
- When two variables show some connection with one another, they are called **associated** variables.
  - Associated variables can also be called **dependent** variables and vice-versa.
- If two variables are not associated, i.e. there is no evident connection between the two, then they are said to be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



<http://xkcd.com/552/>

# 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.

# 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

# 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.

- Wouldn't it 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.

# Illegal Immigrants Reluctant To Fill Out Census Form

by PETER O'DOWD

March 31, 2010 4:00 AM

from KJZZ



Listen to the Story



Morning Edition

3 min 48 sec

+ Playlist

↓ Download



There is an effort underway to make sure Hispanics are accurately counted in the 2010 Census. Phoenix has some of the country's "hardest-to-count" districts. Some Latinos, especially illegal residents, fear that participating in the count will expose them to immigration raids or government harassment.

<http://www.npr.org/templates/story/story.php?storyId=125380052>

# 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.
- **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.



cnn.com, Jan 14, 2012

- **Convenience sample:** Individuals who are easily accessible are more likely to be included in the sample.

# Sampling bias example: Alf Landon vs. Franklin D. Roosevelt (FDR)

A historical example of a biased sample yielding misleading results:

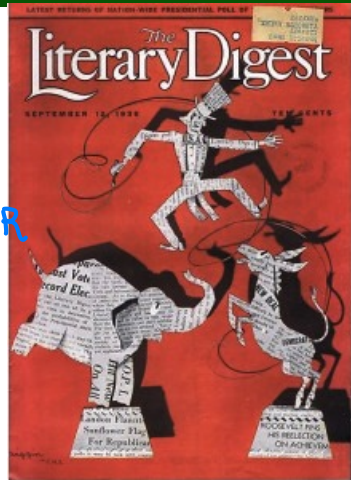


In 1936, Landon sought the Republican presidential nomination opposing the re-election of FDR.



# 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 >> FDR
- 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?

$$\frac{960}{1200}$$

- 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.

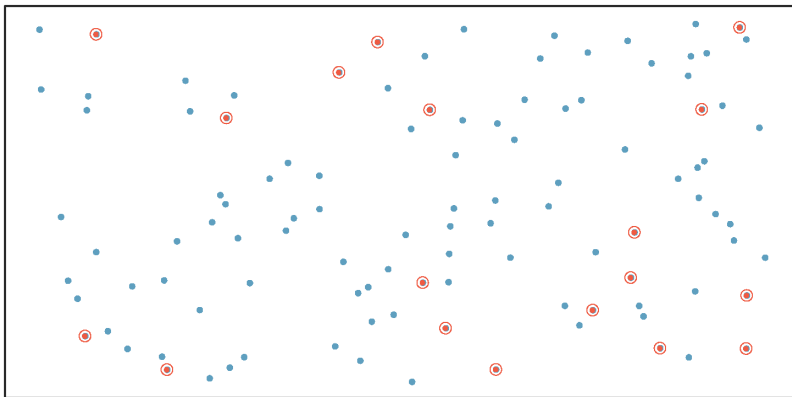
# Obtaining good samples

- 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.



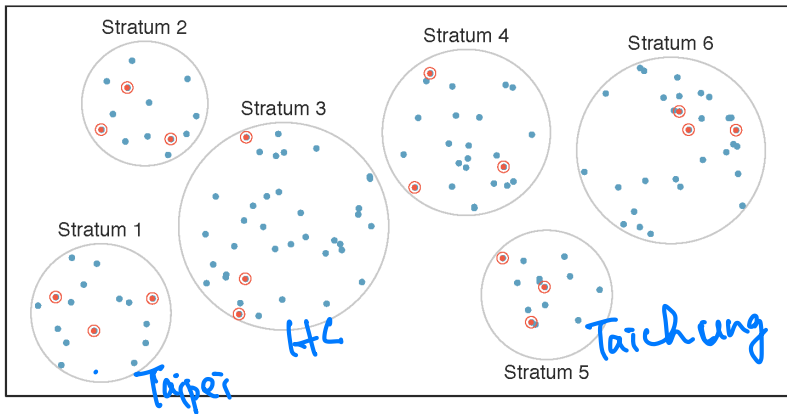
# Simple random sample

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



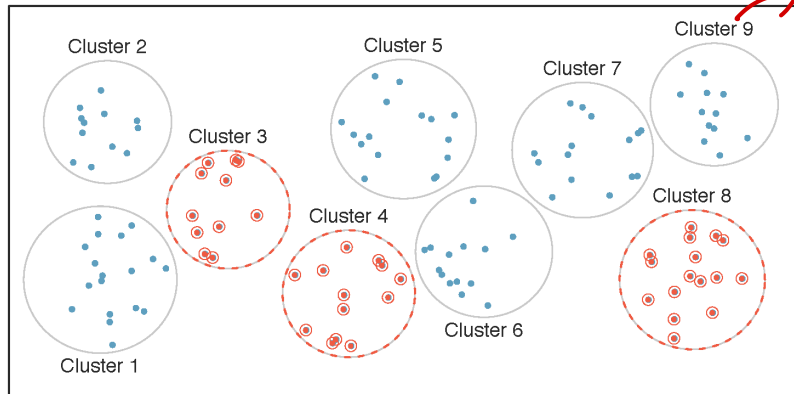
# Stratified sample

**Strata** are made up of similar observations. We take a simple random sample from each 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.

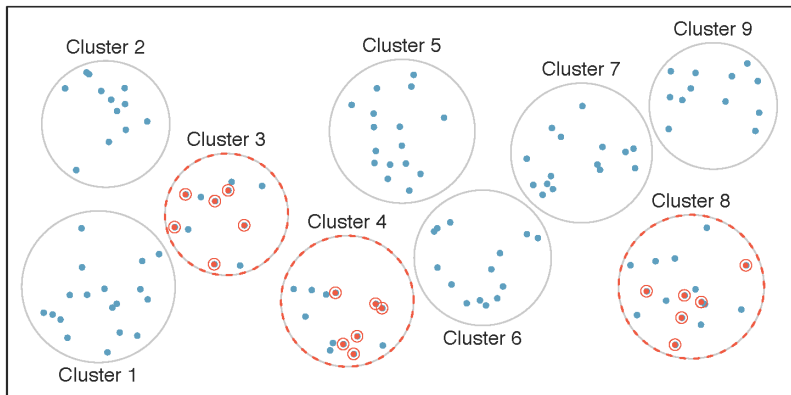


family

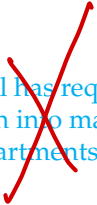
class

# 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

Skip.

Not well defined.

# Principles of experimental design

- ① **Control:** Control for the (potential) effect of variables other than the ones directly being studied.
- 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:
  - 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)



# Difference between blocking and explanatory variables

- Explanatory variables, or factors, are the conditions or treatments we intentionally vary to study their effects.
- Blocking variables are pre-existing traits of units we group together to minimize their impact on the experiment's outcome.
- Unlike stratifying in surveys, blocking in experiments helps ensure any observed differences are due to the treatment, not these traits.

# 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

## What is the main difference between observational studies and experiments?

- Ⓐ Experiments take place in a lab while observational studies do not need to.
- Ⓑ In an observational study we only look at what happened in the past. **Most experiments use random assignment while observational studies do not.**
- Ⓒ Observational studies are completely useless since no causal inference can be made based on their findings.

**Additional Note:** Random sampling is selecting participants from a larger population so that every individual has an equal chance of being included. This helps in generalizing the results to the larger population. Random assignment involves randomly assigning these selected participants to different groups or conditions in an experiment. This ensures that each group is comparable at the start, allowing us to attribute any differences in outcomes to the treatment effect.

# Random assignment vs. random sampling

<i>ideal experiment</i>	Random assignment	No random assignment	<i>most observational studies</i>
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
<i>most experiments</i>	Causation	Correlation	<i>bad observational studies</i>