

## Lecture 1: Laying the Foundations + Terminology

Chapters 1.1-1.2

1 / 22

### Goals for Today

- ▶ Go over the syllabus
- ▶ Show some fun examples
- ▶ Discuss how to evaluate the efficacy of a treatment
- ▶ Describe the different kinds of variables we'll consider

2 / 22

## What is statistics?

(Direct from text) The general scientific process of investigation can be summed up as follows:

1. Identify the scientific question or problem
2. Collect relevant data on the topic
3. Analyze the data
4. Form a conclusion and **communicate it**

Statistics concerns itself with points 2 through 4.

3 / 22

## Your Majors

Biology	11	Economics	5
History	4	Environmental Studies	3
Mathematics	3	Psychology	3
Biochem and Molecular Biology	2	Chemistry	2
International Policy Studies	2	Linguistics	2
Undecided	2	Anthropology	1
Economics/Mathematics	1	Environmental Studies-Hist	1
Environmental Studies-Pol Sci	1	Physics	1
Sociology	1		

4 / 22

## Example: 2012 Election - Nate Silver's Predictions vs Actual Results



Nate Silver's Map



The Actual Map

5 / 22

## Example: Brain & Breast Cancer in Western Washington

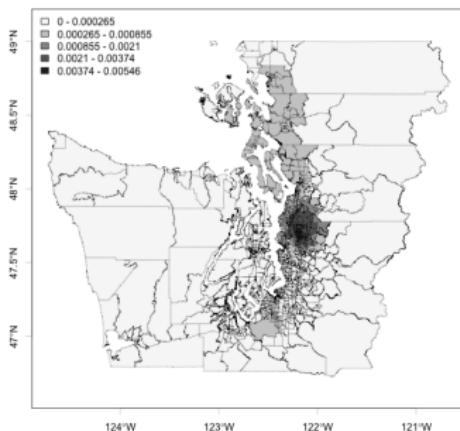
My PhD dissertation involved detecting cancer “clusters”: areas of residual spatial variation of disease risk.

We modeled the (Bayesian) probability of cluster membership for each of the  $n = 887$  census tracts in Western Washington in 2000, using cancer data from 1995–2005, controlling for age, race, and gender.

6 / 22

## Brain Cancer Controlling for Age, Race, & Gender

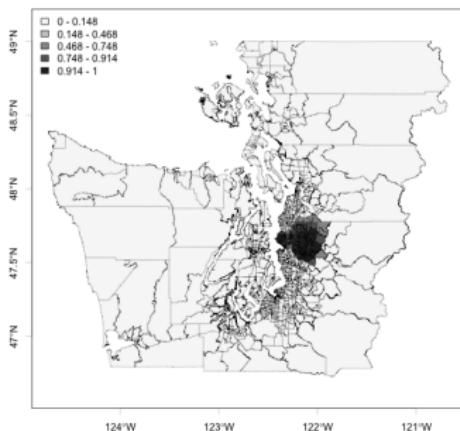
Brain Cancer



7 / 22

## Breast Cancer Controlling for Age, Race, & Gender

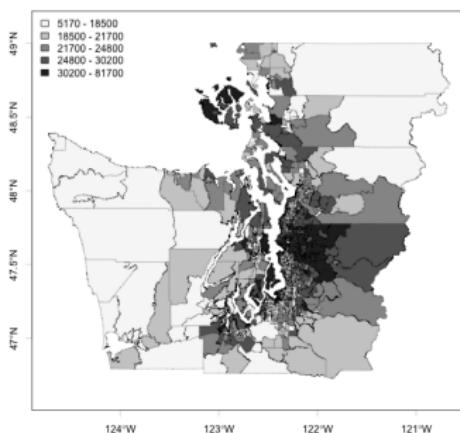
Breast Cancer



8 / 22

## Income per Capita Quintiles

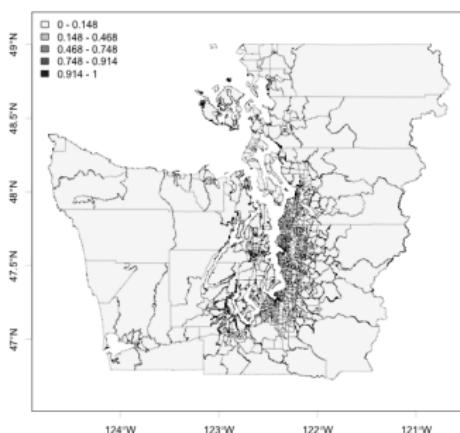
Income Per Capita



9 / 22

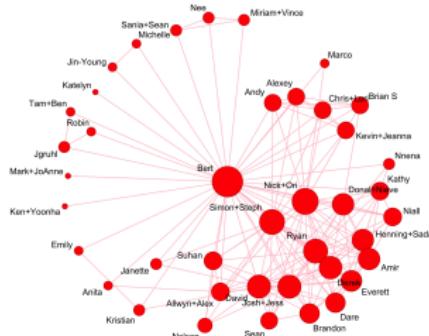
## Breast Cancer Adjusted for Income as Well

Breast Cancer Adjusted for Income



10 / 22

## Example: Social Network Display of a Recent Party I Had



11 / 22

Say we want answer the following questions:

- ▶ Does a new kind of cognitive therapy alter levels of depression in patients?
- ▶ Or you question the effectiveness of antioxidants in preventing cancer.
- ▶ Will reassuring potential new users to a gambling website that we won't spam them increase the sign-up rate?

12 / 22

## Evaluating the efficacy of a 'treatment'

In all the above cases, you are questioning the efficacy of a treatment/intervention. One way to evaluate the efficacy is via an experiment where you define

- ▶ A control group: the "business as usual" baseline group
- ▶ A treatment group: the group that receives/is subject to the treatment/intervention

and make comparisons.

13 / 22

## Website Experiments

**Control:**

Join BettingExpert

Username:

Email:

Password:

I accept the [Terms and Conditions](#)

**Sign up +**



**Treatment:**

Join BettingExpert

Username:

Email:

Password:

I accept the [Terms and Conditions](#)

**(2019 policy) - we will never spam you!**

**Sign up +**

14 / 22

## Example of a treatment vs control

Two other examples in the media of late

- ▶ Facebook's tinkering with user's emotions ([link](#))
- ▶ OkCupid's admission that they experiment on human beings ([link](#))

15 / 22

## Variables

A **variable** is a description of any characteristic whose value may change from one unit in the population to the next:

16 / 22

## Data

At its simplest, data are presented in a data table or matrix where (almost always) each

- ▶ row corresponds to [cases](#) or [units of observation/analysis](#)
- ▶ column represents the variables corresponding to a particular observation

It is almost always the case that

- ▶  $n$  is the number of observations
- ▶  $p$  is the number of variables

17 / 22

## Data Summaries

Consider the variable "federal spending per capita" in each of the 3,143 counties in the US. One can hardly digest this:

```
[1] 6.068095 6.139862 8.752158 7.122016 5.130910 9.973062 9.311835 15.439218
[9] 8.613707 7.104621 6.324061 10.640378 9.781442 8.982702 6.840035 20.330684
[17] 9.687698 11.080738 7.839761 9.461856 9.650295 7.760627 25.774791 13.948106
...
[3121] 7.520731 10.246400 3.106800 17.679572 4.824044 7.247212 8.484211 8.794626
[3129] 9.829593 8.100945 17.090715 4.855849 6.621378 22.587359 10.813260 11.422522
[3137] 9.580265 4.368986 5.062138 6.236968 4.549105 8.713817 6.694784
```

18 / 22

## Data Summaries

We can't interpret all the data at once; we need to boil it down via [summary statistics](#), single numbers summarizing a large amount of data.

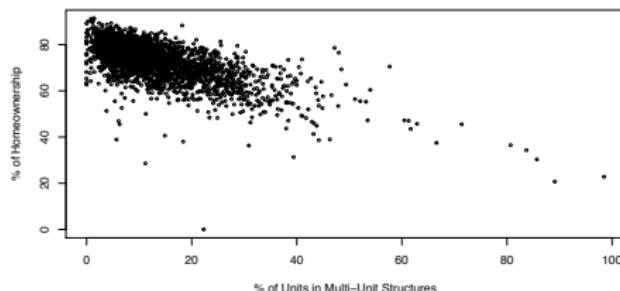
Using the `summary()` command in R:

```
Min. 1st Qu. Median    Mean 3rd Qu.   Max.   NA's
0.000  6.964  8.669  9.991 10.860 204.600        4
```

19 / 22

## Relationships between variables

We can best display the relationship between two variables using a scatterplot AKA [bivariate plot](#):



20 / 22

## Relationships between variables

Almost always we are interested in the relationship between two or more variables.

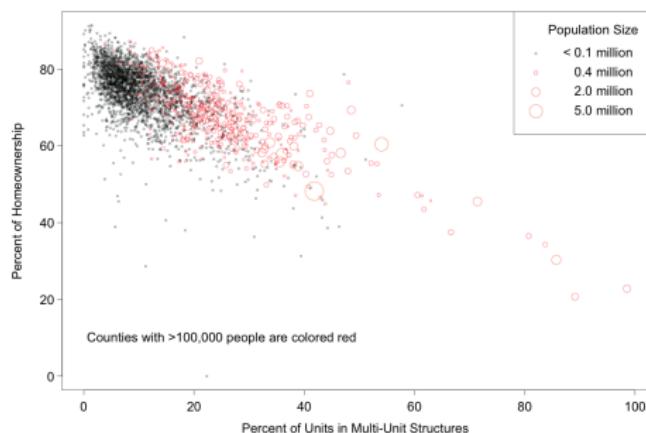
A pair of variables are either related in some way ([associated](#)) or not ([independent](#)). No pair of variables are both associated and independent.

We can have either a [negative association](#) (as the value of one variable increases, the other decreases) or a [positive association](#).

21 / 22

## Relationships between variables

We can consider a third variable in the previous plot.



22 / 22

## Lecture 2: Sampling and Bias

### Chapter 1.3

1 / 17

## Goals for Today

- ▶ Understand important considerations about data collection, in particular **sampling**.
- ▶ Food for thought about the next lecture:  
explanatory/response variables and causality.

2 / 17

## Populations and Samples

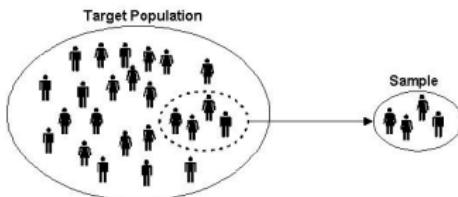
We want to make statements about some aspect of a [study](#)/target population.

1. What proportion of Oregonians smoke?
2. What are the sexual behaviors of males and female Americans in 1948?
3. What proportion of the Reed community believes they have personally experienced offensive, hostile, or intimidating conduct on campus?

3 / 17

## Populations and Samples

It is often not feasible to collect data for every case in the population. If so, we take a [sample](#) of cases.



If the sample is [representative](#) of the desired population then our results will be [generalizable](#).

4 / 17

## Populations and Samples

So say we take a representative sample of 1000 Oregonians and poll their smoking habits. We can then generalize the results to the [entire](#) population of Oregon.

One example of a non-representative sample is a [biased sample](#).

[How do we take a representative sample?](#) In its simplest form, you need to [randomly](#) sample from the entire population. But this is easier said than done.

5 / 17

## Comment on the Representativeness of These Samples:

1. The Royal Air Force wants to study how resistant their airplanes are to bullets. They study the bullet holes on all the airplanes on the tarmac after an air battle against the Luftwaffe (German Air Force).
2. I want to know the average income of Reed graduates in the last 10 years. So I get the records of 10 randomly chosen Reedies. They all answer and I take the average.
3. Imagine it's 1993 i.e. almost all households have landlines. You want to know the average number of people in each household in Portland. You randomly pick out 500 phone numbers from the phone book and conduct a phone survey.
4. You want to know the prevalence of illegal downloading of TV shows among Reed students. You get the emails of 100 randomly chosen Reedies and ask them "How many times did you download a pirated TV show last week?"

6 / 17

## Statistics in Society: Alfred Kinsey

In the mid 20th century, biologist/sexologist Alfred Kinsey wanted to study human sexuality.



At the time sexuality was an extremely taboo subject, very little research had been conducted at that point and Kinsey was astonished at the public's general ignorance.

7 / 17

## Statistics in Society: Kinsey's Questions/Research Problem

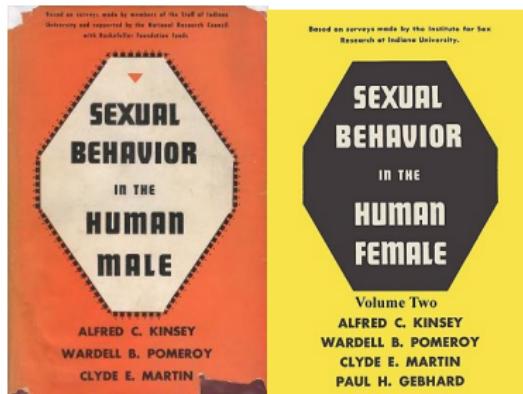
What type of questions was Kinsey interested in? Using his 300 question survey, he hoped to address...

1. What percentage of Americans engaged in premarital and extramarital sex?
2. What were the homosexual tendencies of American males?
3. How common were oral sex and masturbation?
4. ...

8 / 17

## Statistics in Society: Kinsey Reports

The results were published two books on human sexual behavior known as the “Kinsey Reports”: Sexual Behavior in the Human Male (1948) and Female (1953).



9 / 17

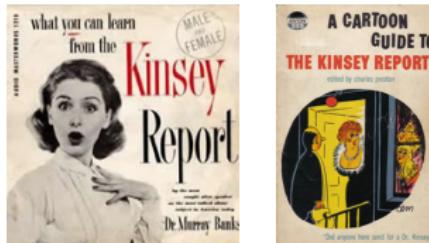
## Statistics in Society: Conclusions of Kinsey Reports

Kinsey claimed, among other things

1. 85% of white men had had premarital sex, 50% had had extra-marital sex
2. Kinsey wrote in 1948 that **one in ten** white men were more or less, exclusively homosexual for at least three years between the ages of 16 and 55.
3. Kinsey reported that oral sex was very common (70% of couples did it), masturbation was very common (almost 63%/92% of women/men did it)

## Statistics in Society: Reaction to Kinsey Reports

Needless to say, people were taken quite aback.



There was also a huge conservative backlash against the reports.

11 / 17

## Statistics in Society: Kinsey's Methods

What were his data collection methods? How did he sample his data? Focusing on the male report, my understanding is that

1. He did in fact base his conclusions on a very large sample size of 5300 males.
2. He sought out volunteers to answer his 300 question survey.
3. He recruited new people by asking previous respondents if they knew other people. This led to a large proportion of his sample to include prison populations and male prostitutes.

What could be some issues?

12 / 17

## Response of the American Statistical Association

The American Statistical Association criticized the sampling procedure. In particular, John Tukey, one of the most eminent statisticians of the time, said

*"A random selection of three people would have been better than a group of 300 chosen by Mr. Kinsey."*

Even though the Kinsey Report was groundbreaking and contributed much to the field of sexology by bringing many topics to the forefront, Kinsey's statements were not generalizable to the general public.

13 / 17

## Reed Campus Climate Survey

During the 2012-2013 academic year Reed contracted Rankin & Associates Consulting to conduct the Campus Climate Survey to "examine the learning, living, and working environment at Reed College."

On page v and iii of the Executive Summary:

[http://www.reed.edu/institutional\\_diversity/campus\\_climate.html](http://www.reed.edu/institutional_diversity/campus_climate.html)

14 / 17

## Examples of Different Types of Bias:

1. **Volunteer bias:** individuals who are more willing to participate have a higher chance of being sampled.
2. **Survival bias:** large segments of the population who “died” are not sampled.
3. **Selection bias:** some individuals are more likely to be selected for study than others.
4. **Convenience sample bias:** individuals who are easily accessible are more likely to be included.

15 / 17

## Moral of the Story

For you:

1. **the consumer of statistics:** Ask yourself what was the study design?
  - ▶ Who is the study population?
  - ▶ Who are the respondents and how were they selected?
2. **the producer of statistics:** think about **how** you will collect your data beforehand. If you want your results to generalize **beyond** just your sample to your study population, your sampling scheme has to be as representative as feasible.

16 / 17

## Explanatory and Response Variables

Example: A medical doctor pours over some his patients' medical records and observes:



He then posits the following **causal** relationship:

- ▶ **Explanatory variable:** sleeping with shoes on
- ▶ **Response variable:** waking up with headaches

What's wrong with hypotheses?

# Lecture 3: Observational Studies + Randomized Experiments + Confounding + Simpson's Paradox

## Chapter 1.4

1 / 26

### Goals for Today

- ▶ We illustrate the difference between
  - ▶ an observational study
  - ▶ a randomized experiment, where the treatment is assigned at random.
- ▶ Introduce the notion of confounding AKA lurking variables
- ▶ Discuss Simpson's Paradox (not in textbook).

2 / 26

## Going Back to Previous Example

Going back to the study on



- ▶ The explanatory variable was: sleeping with your shoes on
- ▶ The response variable was: waking up with a headache
- ▶ The doctor hypothesized a causal relationship

3 / 26

## Confounding Variable AKA Lurking Variable

This is an example of confounding. A confounding variable affects both the explanatory and response variable. So if:

4 / 26

## Controlling for Potential Confounding

One way to control for (i.e. take into account) confounding is to do an exhaustive search for all such variables. This is not always practical.

Another way is via an experiment where we randomly assign individuals to a treatment or a control group in a randomized experiment.

5 / 26

## Back to Shoes and Headaches

So imagine we recruit 10,000 people for our study and randomly assign 5000 people to each of:

- ▶ Treatment: sleep with shoes on
- ▶ Control: sleep with shoes off

In this table

Group	n	# with headache
Treatment	5000	$n_1$
Control	5000	$n_2$
Total	10,000	$n_1 + n_2$

$n_1$  and  $n_2$  won't be very different.

6 / 26

## Observational Studies vs Randomized Experiments

The key word from the study design above was **randomly assign**.

- ▶ **Observational studies:** a study where researchers have **no control** over who receives the treatment
- ▶ **Randomized experiments:** a study where researchers not only have control over who receives the treatment, but also make the assignments **at random**.

7 / 26

## Observational Studies vs Randomized Experiments

**Conclusion:** The study introduced at the end of the last lecture is an **observational study**, so we cannot conclude that wearing shoes when you sleep **causes** you wake up with a headache.

**Mantra:** **Correlation is not causation** Just because two variables appear to be associated/correlated, does not mean that one is **causing the other**.

- ▶ Spurious correlations: <http://www.tylervigen.com/>
- ▶ Saturday Morning Breakfast Cereal:  
<http://www.smbc-comics.com/?id=3129>

8 / 26

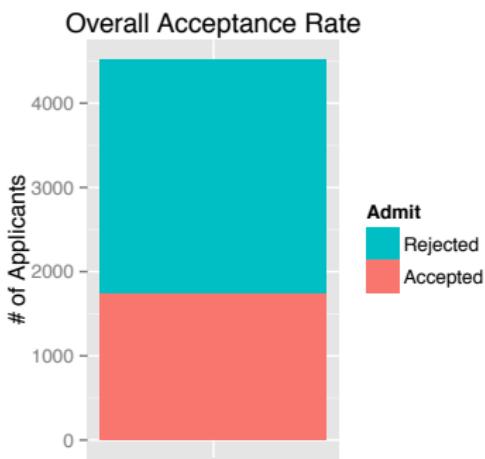
## Well-Known Example of Confounding

A famous example of an unaccounted for confounding variable having serious repercussions was when the UC Berkeley was sued in 1973 for bias against women who had applied for admission to graduate schools.

Let's consider the  $n = 4526$  people who applied to the 6 largest departments.

9 / 26

Of the  $n = 4526$  applicants:



10 / 26

Split the counts by gender:



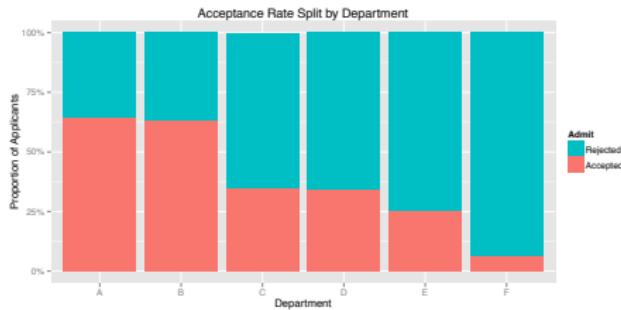
11 / 26

Look at proportions instead of counts:



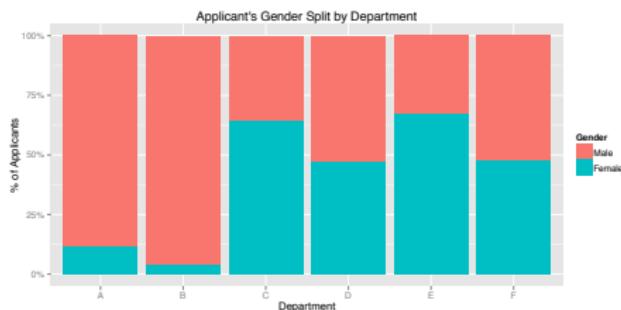
12 / 26

## What was the “competitiveness” of departments?



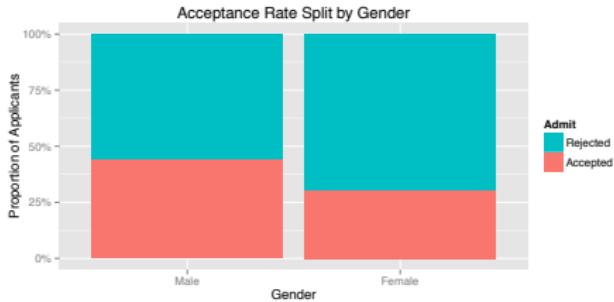
13 / 26

## Where were the women applying?



14 / 26

So while in aggregate things looked like this:



15 / 26

You need to account for department!



16 / 26

## Bickel et al.'s (1975) Explanation

There was the presence of a confounding variable: competitiveness of applying to the department, which is a function

- ▶ number of applicants
- ▶ number of available slots

So it wasn't that departments were discriminating against women, rather:

- ▶ women tended to apply to departments with high competition and hence lower admission rates, primarily the humanities.
- ▶ men tended to apply to departments with low competition and hence higher admission rates, primarily the sciences.

17 / 26

## Bickel et al.'s (1975) Explanation

In fact, Bickel et al. found that "If the data are properly pooled...there is a small but statistically significant bias in favor of women."

This was the exact opposite claim of the lawsuit. This is known as Simpson's Paradox.

18 / 26

## Simpson's Paradox

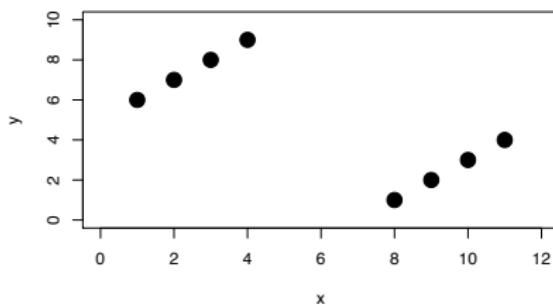
(From Wikipedia) Simpson's paradox occurs when a trend that appears in different groups of data disappears when these groups are combined, and the [reverse trend](#) appears for the aggregate data.

This is due to a confounding variable.

19 / 26

## A Graphical Illustration of Simpson's Paradox

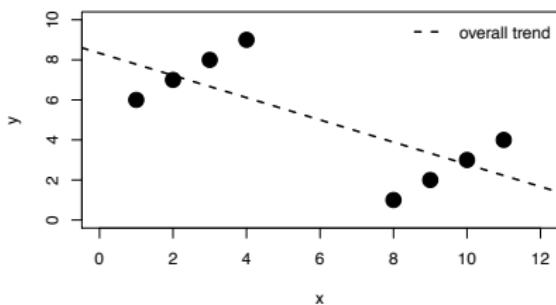
Say we have the following points:



20 / 26

## A Graphical Illustration of Simpson's Paradox

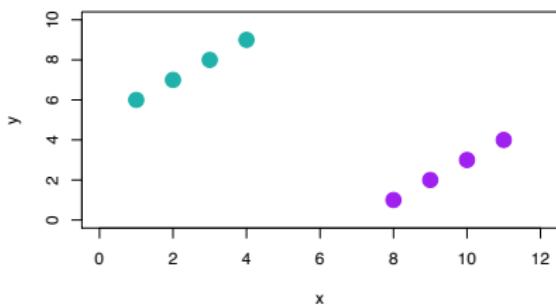
Overall, if we fit a single line, the explanatory variable  $x$  is negatively related with the outcome variable  $y$ :



21 / 26

## A Graphical Illustration of Simpson's Paradox

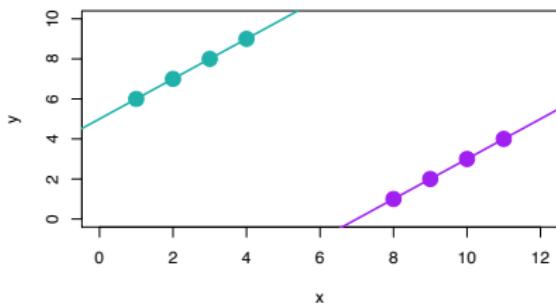
But say we consider a confounding variable, in this case color, and fit two separate lines for each group:



22 / 26

## A Graphical Illustration of Simpson's Paradox

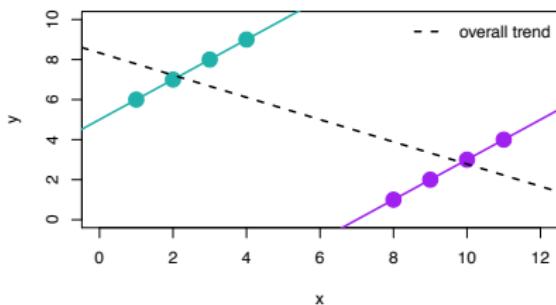
The subgroups now exhibit a [positive relationship](#)!



23 / 26

## A Graphical Illustration of Simpson's Paradox

i.e. the trend in aggregate is the [reverse](#) of the trend in the subgroups (teal & purple).



24 / 26

## Bickel et al.'s (1975) Conclusion

"The bias in the aggregated data stems not from any pattern of discrimination on the part of admissions committees, which seem quite fair on the whole, but apparently from prior screening at earlier levels of the educational system."

"Women are shunted by their socialization and education toward fields of graduate study that are generally more crowded, less productive of completed degrees, and less well funded, and that frequently offer poorer professional employment prospects."

The original paper can be found [here](#).

25 / 26

## Next time

We will discuss

- ▶ Specific types of sampling beyond just [simple random sampling](#), as this is not always feasible
- ▶ Experimental design: some key principles to keep in mind when evaluating the efficacy of treatments.

26 / 26

## Lecture 4: Sampling Methods + Design of Experiments

Chapter 1.4.2 + 1.5

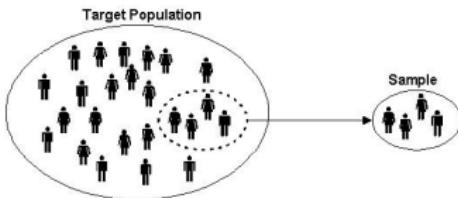
1 / 22

### Goals for Today

- ▶ Discuss different types of sampling
- ▶ Designing experiments
- ▶ Very important example: clinical trials
- ▶ Example of my own designed experiment: Fried Chicken Face Off

2 / 22

## Recall from Lecture 1.3: Population and Samples

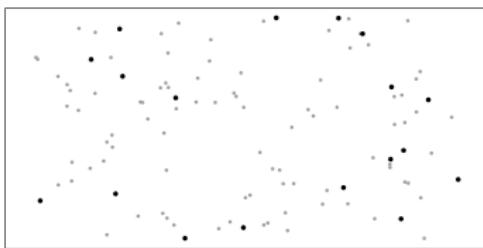


If the sample is representative of the desired population then our results are **generalizable**.

How do we take a representative (i.e. unbiased) sample? You **randomly** sample from the population.

3 / 22

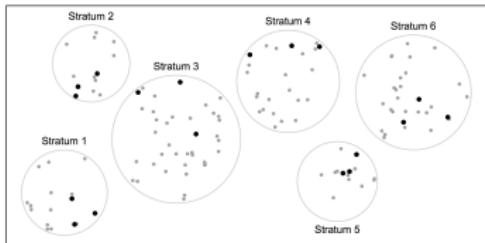
### 1. Simple Random Sampling



**Most granular sampling:** Where every individual in the population has the same probability of being sampled. Here, all dots are members of the population, and the bolder dots are sampled.

4 / 22

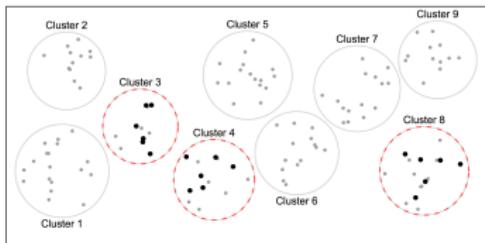
## 2. Stratified Sampling



**Divide and conquer:** The population is divided into strata, and we sample from each strata. For example, each strata could be a census tract in Oregon, and we sample 3 individuals from each strata.

5 / 22

## 3. Cluster Sampling



**Two stage sampling:** Very similar to stratified sampling in its process, except that there is no requirement to sample from every cluster. First the clusters in red were chosen at random, and then we sample from them.

6 / 22

## Three Different Types of Sampling

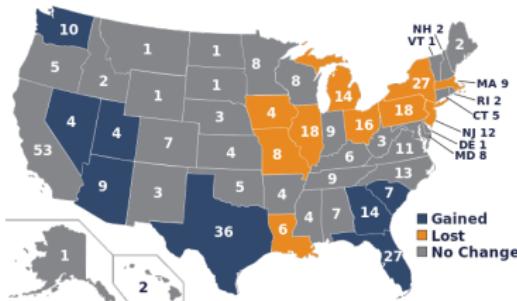
1. Simple random sampling: most granular sampling
2. Stratified sampling: divide and conquer
3. Cluster sampling: two-stage sampling

The mathematics behind the stratified and cluster sampling are more complicated to account for the hierarchies involved. Ex: for stratified sampling use the Horvitz-Thompson estimator.

7 / 22

## Statistics in Society: The US Census

The purpose of the decennial US census is [congressional apportionment](#): the 435 seats in the US House of Representatives get distributed to the 50 states in proportion to their population.  
After the 2010 census:



8 / 22

## Statistics in Society: The US Census

President Bill Clinton's administration planned on using sampling in the 2000 census. In an article dated in 1996:

The screenshot shows a news article from The New York Times. At the top, there is a navigation bar with links for HOME PAGE, TODAY'S PAPER, VIDEO, MOST POPULAR, TIMES TOPICS, and MOST RECENT. Below the navigation bar, the site's logo 'The New York Times' is displayed next to the word 'U.S.'. A search bar with the placeholder 'Search All NYTimes.com' and a 'Go' button are also present. On the left side of the main content area, there is a sidebar with a 'More Like This' section containing links to related articles: 'U.S. CENSUS BUREAU REJECTS REVISION TO COUNTING'S TALLY', 'The Nation: Sample Case; You Fill Up My Census, Even If I...', 'Lessons From the Election That Shock America', 'Find More Stories', and 'Statistical Methods'. The main article title is 'In a First, 2000 Census Is to Use Sampling'. Below the title, it says 'By STEVEN A. HOLMBERG Published: February 23, 1996'. The article discusses how the Census Bureau plans to use sampling to count nearly 90 percent of the United States population in 2000, relying on statistical sampling methods to determine the number remaining. It notes that this is the first time the official tally of the American population, done every 10 years and used to apportion seats in the House of Representatives, will be based in part on a scientifically determined estimate rather than the actual head count conducted through a mass direct-mail campaign. The article quotes Census Bureau officials and the Census Bureau Director, Martha Farnsworth Riche, discussing the goals of reducing costs and increasing accuracy.

In a First, 2000 Census Is to Use Sampling

By STEVEN A. HOLMBERG  
Published: February 23, 1996

To cut costs and improve accuracy, the Census Bureau said today that it would actually count only 90 percent of the United States population in 2000 and rely on statistical sampling methods to determine the number remaining.

The plans, announced at the Commerce Department, mean that for the first time the official tally of the American population, done every 10 years and used to apportion seats in the House of Representatives, will be based in part on a scientifically determined estimate rather than the actual head count conducted through a mass direct-mail campaign.

Census Bureau officials say the revised method is needed to keep costs down and to avoid a repeat of the 1990 census, which missed record numbers of people that had been traditionally hard to count, mainly members of ethnic and racial minorities.

"What we intend to do is to meet our twin goals of reducing costs and increasing accuracy is to make a much greater use of widely accepted scientific statistical methods, and sampling is first and foremost among them," said Martha Farnsworth Riche, the Census Bureau Director.

9 / 22

## Statistics in Society: The US Census

However, Article I, Section 2 of the US Constitution states: *The actual Enumeration shall be made within three Years after the first Meeting of the Congress of the US, and within every subsequent Term of ten Years...*

As such, the Supreme Court ruled 5-4 in 1999 that

- ▶ sampling could not "under any circumstances" be used to reapportion U.S. House seats
- ▶ could be used for other purposes such as redrawing state legislative districts or allocating federal funds to cities and states

10 / 22

## Statistics in Society: The Census

THE WALL STREET JOURNAL

U.K. EDITION Friday, May 15, 2009 As of 4:42 PM EDT

Home World U.S. Business Tech Markets Market Data Your Money Opinion Life & Culture N.Y. Real Estate Management

Sets & Wines Politics & Policy Washington Wire Economy Health Law Rollout WSJ/NBC News Poll Journal Report Columns & Blogs

TOP STORIES IN POLITICS 1 of 12 2 of 12 3 of 12

How Serious Are Risks If U.S. Doesn't Act? White House Wants 'Hard Look' at Syria Weapons Offer Senators Factor Voters Into Their Syria Equation Discord Over Military Strike Imperils President's Agenda

POLITICS | May 15, 2009, 4:42 p.m. EDT

### Census Nominee Rules Out Statistical Sampling in 2010

Article Comments Tweet 1

E-mail Print Save Share Facebook Twitter LinkedIn

A A Available to WSJ.com Subscribers

By TIMOTHY J. ALBERTA

WASHINGTON—President Barack Obama's nominee to head the Census Bureau on Friday ruled out using statistical sampling to adjust the results of the 2010 census, quelling Republican concerns and making his confirmation likely next week.

Robert Groves, director of the University of Michigan's Survey Research Center and a former Census Bureau official, is an expert on statistical sampling, the practice of extrapolating a larger population from a smaller slice of it. Proponents of sampling say it helps produce a more accurate tally of the population, especially when it comes to traditionally undercounted groups, such as minorities living in urban areas.

But many Republican lawmakers insist that sampling violates the Constitution, which calls for an "actual Enumeration" of the population every 10 years. Critics also say the use of sampling would politicize the traditionally nonpolitical Census Bureau.

11 / 22

## Principles Of Designing Experiments

Switching gears...

(Wikipedia) In general usage, **design of experiments (DOE)** or **experimental design** is the design of any information-gathering exercises where variation is present, whether under the full control of the experimenter or not.

However, in statistics, these terms are usually used for **controlled experiments**: experiments where there is a control and treatment group.

## Principles Of Designing Experiments

1. **Controlling:** We want to control for differences between the two groups.
2. **Randomization:** We randomize individuals into treatment vs control so that any differences in uninteresting variables even out in the long run.
3. **Replication:** The more cases we observe, the more “precise” the results.
4. **Blocking:** Researchers sometimes know or suspect that variables, other than the treatment, influence the response. In this case, they may first group individuals based on this variable into blocks and then randomize cases within each block.

13 / 22

## Clinical Trials

To evaluate the efficacy of a drug, they must be subject to a **clinical trial**. The gold standard for a clinical trial is **randomized controlled trial**. i.e. randomized control and treatment groups.

- ▶ **Blinded study:** When researchers do not inform patients which group, or arm, they are in
- ▶ **Double blinded study:** When the person administering the treatment/control themselves do not know which group the patient is in.
- ▶ **Placebo:** Fake treatment. Sometimes the **thought** alone of having a treatment can influence behavior / health

14 / 22

## Example of Mine: Ezell's Famous Chicken

In Seattle's Central District lies



From Wikipedia: Oprah Winfrey called it her favorite fried chicken. There are a number of photos of her on the wall of the original restaurant proclaiming her love of the chicken. It is also said she has the chicken flown to her in Chicago when she has a craving.

15 / 22

## Example of Mine: Ezell's Famous Chicken

One day I was raving about Ezell's Chicken. My friend Nick accused me of being another person "buying into the hype"; that if people were subjected to a blinded taste test, Ezell's would fare no better than KFC. So...



vs



We set up a "Fried Chicken Face Off" where we would have individuals try both kinds of chicken and rate which one they liked more.

16 / 22

## Design of Experiment Principles in Place

**Goal:** Evaluate which kind of chicken, Ezell's or KFC, that people prefer in a blinded taste test. (Not if participant can determine which chicken came from which restaurant.)

**Question:** What principles of the design of experiments should be put in place to this end?

17 / 22

## Design of Experiment Principles in Place

The design principles we put in place:

- ▶ **Single blinded:** The taster doesn't know which (Ezell's or KFC) chicken they are eating, but the server does.
- ▶ **Randomizing** which kind of meat (wing, breast, leg) between tasters. Each taster would try two kinds of meat.
- ▶ **Controlling for which kind of meat within a taster:** i.e. if you eat a KFC wing, you will necessarily eat an Ezell's wing
- ▶ **Randomizing** which order of chicken you eat: KFC first or not

18 / 22

## Design of Experiment Principles in Place

The design principles we put in place:

- ▶ **Controlling for temperature:** hence we're picking a place that is central to both Ezell's and KFC given the traveling required.
- ▶ **Controlling for visual look:** We thought blind-folds were a bit excessive
- ▶ **Controlling for kind of batter:** we can't do KFC crispy chicken b/c Ezell's doesn't have that type of batter. This is a limitation of the study b/c some feel the crispy chicken is better, but we have no choice.
- ▶ Just one **replicate** of each kind of meat.

19 / 22

## Results

Final score: KFC 8, Ezell's 4.

Some notes:

- ▶ Even though people were "blinded", most knew which the two pieces were from KFC.
- ▶ People generally felt the chicken meat from Ezell's was better, and this was magnified as the chicken went cold.
- ▶ However, they felt the skin was better at KFC. Given that fried chicken is what it is b/c of the skin, people voted for KFC.
- ▶ Future metrics need to consider the chicken and the skin separately, as well as the "overall experience" scores. i.e. this face off should be viewed as a **pilot study**

20 / 22

## Caution: Grad Students NOT at Work



21 / 22

## Next time

Examining and visualizing numerical data

22 / 22

## Lecture 5: Visualizing Numerical Data

Chapter 1.6 + 1.7

1 / 1

### Goals for Today

- ▶ Visualizing numerical data
  - ▶ Two famous historical examples of data visualization
  - ▶ Reed's 2013 entering class
- ▶ Histograms
- ▶ Measures of Central Tendency: Mean, Median, and Mode
- ▶ Measure of Spread: Sample variance and sample standard deviation

2 / 1

## Famous Example 1: Napoleon's March on Russia in 1812

In 1812, Napoleon led a French invasion of Russia, at one point marching on Moscow.



3 / 1

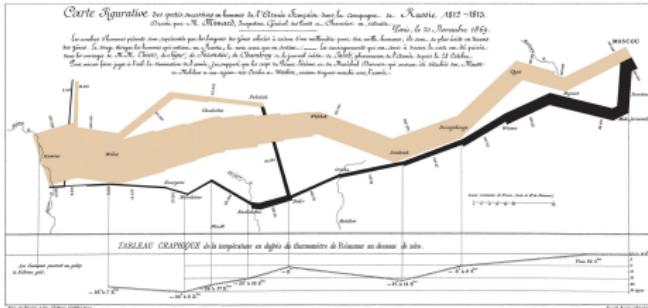
## Famous Example 1: Napoleon's March on Russia in 1812

The advance and retreat on Moscow was an unmitigated disaster:



4 / 1

## Famous Example 1: Napoleon's March on Russia in 1812



5 / 1

## Famous Example 1: Napolean's March on Russia in 1812

Why is this visualization big deal?

On a two-dimensional page, it displays 6 variables (in others words, 6 dimensions of information) at once:

6 / 1

## Famous Example 2: 1854 Broad Street Cholera Outbreak

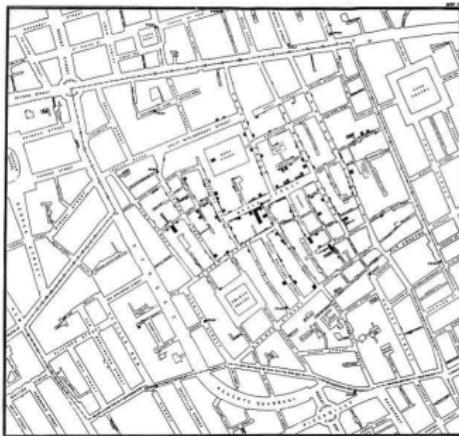
On August 31 1854, an epidemic of cholera began in the Soho neighborhood of London. Over the next three days 127 people near Broad Street had died.

Dr. John Snow, a physician, was a student of the disease. (From Wikipedia) Snow was a skeptic of the then-dominant [miasma theory](#) that stated that diseases such as cholera or the Black Death were caused by pollution or a noxious form of “bad air.”

Snow created the following map to investigate:

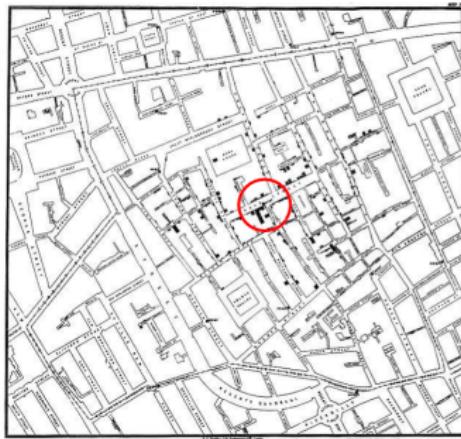
7 / 1

## Famous Example 2: 1854 Broad Street Cholera Outbreak



8 / 1

## Famous Example 2: 1854 Broad Street Cholera Outbreak



9 / 1

## Famous Example 2: 1854 Broad Street Cholera Outbreak

He identified the source of the outbreak as water from the [Broad Street Pump](#), which was near a cesspit that began to leak.



This led to discovering that cholera was transmitted by food and water being contaminated by fecal matter and not via the air. This was a watershed moment in the emerging field of epidemiology.

10 / 1

## Histograms

In the `openintro` package, the `email150` dataset contains a random sample of 50 emails, in which researchers try to identify emails as spam. One variable is the # of characters:

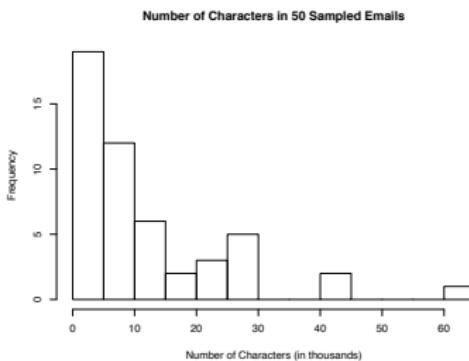
Characters	0-4.999	5-9.999	10-14.999	...	60-64.999
(in 1000's)					
Count	1	19	12	6	...
					1

So each of the intervals 0-5, 5-10, 10-15, etc. are buckets/bins and we count the number of emails in each bucket/bin.

11 / 1

## Histograms

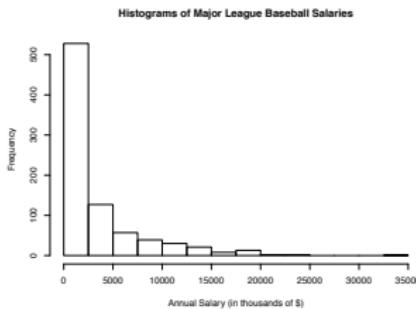
Histograms provide a description of the shape of the [distribution](#) of data.



12 / 1

## Skew and Long Tail

Also in the `openintro` package is MLB salary data in 2010. If we plot a histogram:

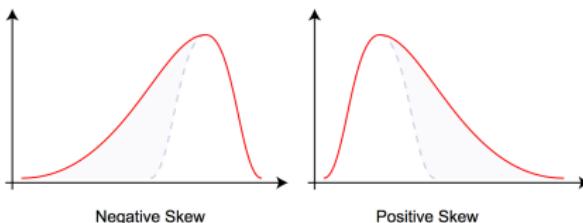


The data has a **long tail** to the right: data is **right-skewed**. i.e. a small number of players who make a **VERY** large amount of money.

13 / 1

## Trick to Remembering Which Skew is Which

- ▶ Long tail to the right: data is **right-skewed** AKA **positively-skewed**
- ▶ Long tail to the left: data is **left-skewed** AKA **negatively-skewed**



14 / 1

## Reed's 2013 US-Originating Entering Class

What can we do about skewed data?

<http://rpubs.com/rudeboybert/reed2013>

15 / 1

## Mean

The mean, AKA average, is a common way to measure the center of the data. So for example, the mean of 1, 2, 5, 3, and 7 is

$$\frac{1 + 2 + 5 + 3 + 7}{5} = 3.6$$

We label the sample mean  $\bar{x}$  (pronounced "x bar"):

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

where  $x_1, x_2, \dots, x_n$  are the  $n$  observed/sampled values.

16 / 1

## Median

The median, however, is the middle number.

Two cases:

- ▶ Odd number of values: the median of (1, 3, 5, 8, 10) is 5.
- ▶ Even number of values: the median of (1, 3, 5, 8) is the average of the middle two values:  $\frac{3+5}{2} = 4$

But why use the median at all?

17 / 1

## Mean vs Median: Imaginary Scenario

- ▶ Say at company X, there 5 employees: the CEO and everyone else.
- ▶ The CEO earns \$1000 an hour, while the others earn \$20, \$21, \$30, and \$40 an hour.
- ▶ The employees complain that they are paid too little.
- ▶ The CEO counters that the mean hourly salary is  $\bar{x} = \frac{20+21+30+40+1000}{5} = 222.20$  an hour, which is really high.

18 / 1

## Mean vs Median: Imaginary Scenario

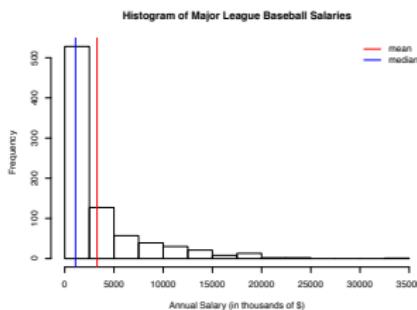
The CEO's extreme salary is inflating the mean. A more appropriate measure is the median hourly salary of 30.

Medians are less sensitive to (i.e. more robust to) outliers than the mean.

Ex: the “median home price” is typically used, because it isn’t as sensitive as the mean to the few very expensive houses.

19 / 1

## Mean vs Median: Back to MLB Salary Data



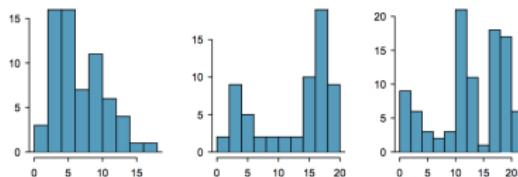
20 / 1

## Mode

A **mode** is the value that appears the most often in a data set. So out of (1, 3, 3, 5, 6), the modal value is 3.

Modes also describe **peaks**, but this can get subjective.

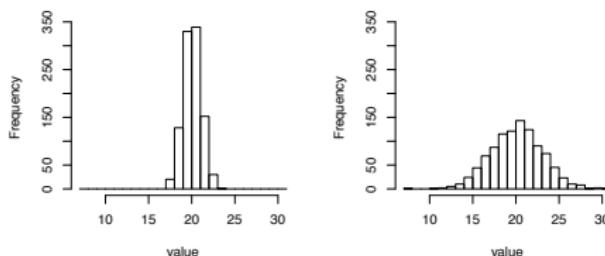
A distribution can be **unimodal**, **bimodal**, or **multimodal**:



21 / 1

## Measure of Spread

Next, consider the following two histograms: Both have mean of about 20. What is the difference between them?



22 / 1

## Measure of Spread

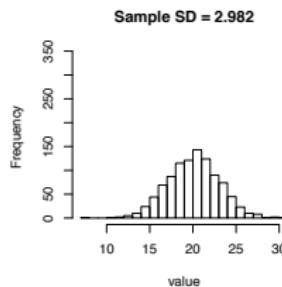
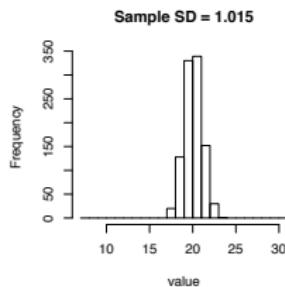
We need a measure of **spread/variability**. The **sample variance  $s^2$**  is roughly the average squared distance from the mean.

The **sample standard deviation  $s$**  is the square root of the sample variance. The sample standard deviation is useful when considering how close the data are to the mean.

23 / 1

## Measure of Spread

Back to example:



24 / 1

## How to Compute the Sample Standard Deviation

Read section 1.6.4. The formula really doesn't make much intuitive sense, but is the way it is due to mathematical convenience. Fortunately there is an R command that computes it for you: `sd()`

25 / 1

## Next Time

- ▶ Another simple data visualization tool: boxplots
- ▶ Examining/Visualizing Categorical Data

26 / 1

# Lecture 6: Visualizing Numerical and Categorical Data

Chapter 1.6+1.7

1 / 1

## Goals for Today

- ▶ Rule of thumb for standard deviations
- ▶ Population vs sample mean/variance/standard deviations
- ▶ Percentiles and Quartiles
- ▶ Boxplots
- ▶ Piecharts, barplots, mosaicplots

2 / 1

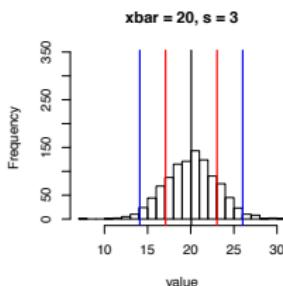
## Rule of Thumb for Standard Deviations

If the data distribution is bell-shaped, then

- ▶ about  $\frac{2}{3}$  of the data will be within one SD of the mean (book says 70%).
- ▶ about 95% of the data will be within two SD.

3 / 1

## Example

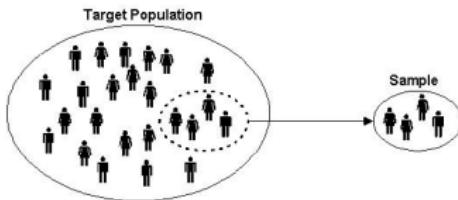


- ▶ black line is mean  $\bar{x}$
- ▶ red lines mark about  $\frac{2}{3}$ :  
 $[\bar{x} - s, \bar{x} + s] = [20 - 3, 20 + 3] = [17, 23]$ .
- ▶ blue lines mark about 95%:  
 $[\bar{x} - 2s, \bar{x} + 2s] = [20 - 6, 20 + 6] = [14, 26]$ .

4 / 1

## Population vs Sample Mean/Variance/Standard Deviation

Recall the notion of taking a **representative sample** from a **study/target population**. Say we are interested in the income of the individuals.



5 / 1

## Population vs Sample Mean/Variance/Standard Deviation

- ▶ The **sample mean  $\bar{x}$**  is the mean income of the 4 sampled people.
- ▶ The **population mean  $\mu$**  is the mean income of all 24 people in the target population.
- ▶ We say  $\bar{x}$  **estimates  $\mu$** . If the sample is representative, then  $\bar{x}$  estimates  $\mu$  with high **accuracy** i.e. it is unbiased.

6 / 1

## Population vs Sample Mean/Variance/Standard Deviation

	True Population Value	Sample Value
Mean	$\mu$	$\bar{x}$
Variance	$\sigma^2$	$s^2$
Standard Deviation	$\sigma$	$s$

The sample value is used to estimate the (true) population value.

7 / 1

## Percentiles

A percentile (%'ile) indicates the value below which a given %'age of observations fall.

SAT Scores from 2012

<http://media.collegeboard.com/digitalServices/pdf/research/SAT-Percentile-Ranks-2012.pdf>

So for example, if you scored 700 in critical reading, 95% of college-bound seniors who took the test did worse.

8 / 1

## Quartiles

Quartiles split up the data into 4 intervals, each with about one quarter of the data:

- ▶ The lower quartile is the 25th %'ile
- ▶ The median is the 50th %'ile
- ▶ The upper quartile is the 75th %'ile

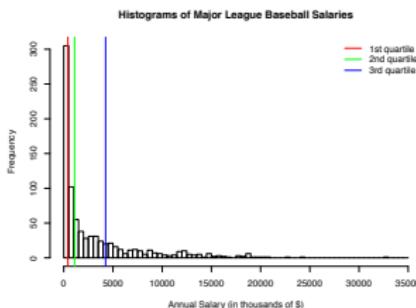
The interquartile range (IQR) is another measure of the spread of a sample:

$$\text{IQR} = \text{upper quartile} - \text{lower quartile}$$

9 / 1

## MLB Data Quartiles

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
400.0	418.3	1094.0	3282.0	4250.0	33000.0



The IQR is  $(\text{3rd Quartile} - \text{1st Quartile}) = 4250.0 - 418.3 = 3831.7$   
i.e the distance between the red and blue line.

10 / 1

## Robust Statistics (Chapter 1.6.6)

Robust estimates are statistics where extreme observations (outliers) have less effect on their values, i.e. are more resistant to their effect. The median and IQR are two examples.

Example: Old scoring system in figure skating: drop the highest & lowest scores and then take the average.

Say we have a figure skater who gets judged by countries V-Z:

Country	V	W	X	Y	Z
Score	4.0	5.2	5.2	5.3	6.0

Drop the 4.0 and 6.0, then the final score is:  $\frac{5.2+5.2+5.3}{3} = 5.23$

11 / 1

## Boxplots

Boxplots are visual summaries of a sample  $x_1, \dots, x_n$  that bring to light unusual values (potential outliers):

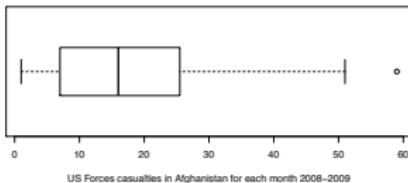
Example: # US Forces casualties in the war in Afghanistan for each month from 2008-2009:

7, 1, 7, 5, 16, 28, 20, 22, 27, 16, 1, 3, 14, 15, 13, 6, 12, 24, 44, 51, 37, 59, 17, 17

12 / 1

## Boxplots

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.00	7.00	16.00	19.25	24.75	59.00



Page 29 of text describes the length of the **whiskers**: they capture data that is no more than  $1.5 \times IQR$  of both ends of the box.

13 / 1

## Outliers Are Relatively Extreme

An **outlier** is an observation that appears extreme relative to the rest of the data.

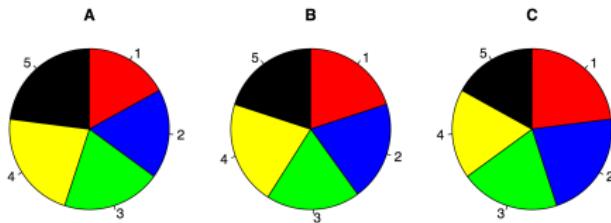
Why it is important to look for outliers? Examination of data for possible outliers serves many useful purposes, including

- ▶ Identifying strong skew in the distribution.
- ▶ Identifying data collection or entry errors.
- ▶ Providing insight into interesting properties of the data.

14 / 1

## Piecharts

Say we have the following piecharts represent the polling from a local election with five candidates (1-5) at three different time points A, B, and C:

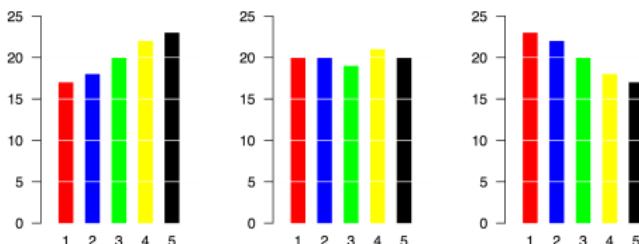


Answer the following questions:

- ▶ In the first race, is candidate 5 doing better than candidate 4?
- ▶ Who did better between time A and time B, candidate 2 or candidate 4?

15 / 1

## Barplots Instead

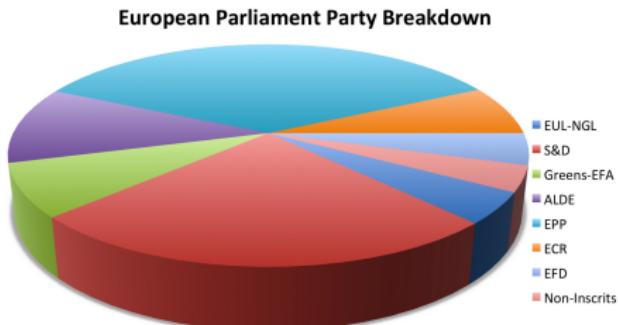


Answers:

- ▶ Candidate 5 is doing better than 4
- ▶ Between A and B, candidate 2 went from about 17% to 20% while candidate 4 went from about 22% to 21%. So candidate 2 did better

16 / 1

## 3D Piecharts Can Be Deceiving



EEP (teal) has 266 seats, whereas S&D (red) has 190 seats.

17 / 1

## Titanic Survival Data

Typing data(Titanic) in R loads the survival and death counts, split by each of the following categories:

- ▶ Class: 1st, 2nd, 3rd, or crew (4 levels)
- ▶ Gender (2 levels)
- ▶ Age: Child or adult (2 levels)

i.e.  $4 \times 2 \times 2 = 16$  possible groups to consider.

### Questions

- ▶ What was the effect of class (1st, 2nd, 3rd, crew) on your chances of survival?
- ▶ Did the “women and children” first lifeboat policy hold?

18 / 1

## Frequency Table

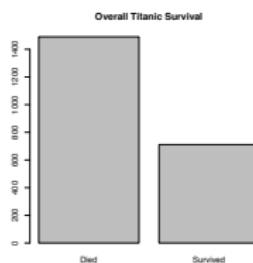
A table summarizing a single categorial variable is called a **frequency table**. Overall:

Died	1490
Survived	711
Total	2201

19 / 1

## Barplot

**Barplots** are ways to display categorial variables:



20 / 1

## Contingency Table

A table that [cross-classifies](#) two categorical variables is a [contingency table](#). Now let's split survival by class: 1st, 2nd, 3rd, and crew.

Before:

Died	1490
Survived	711
Total	2201

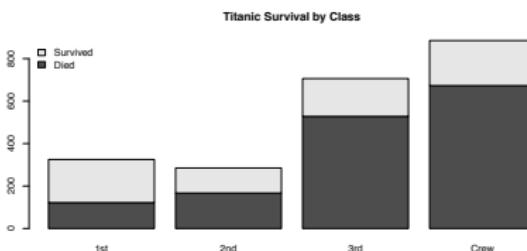
After:

	1st	2nd	3rd	Crew	Total
Died	122	167	528	673	1490
Survived	203	118	178	212	711
Total	325	285	706	885	2201

21 / 1

## Stacked Barplot

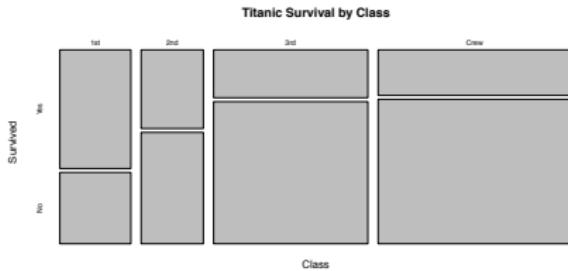
[Stacked barplots](#) are one way to display values from a contingency table:



22 / 1

## Mosaic Plots

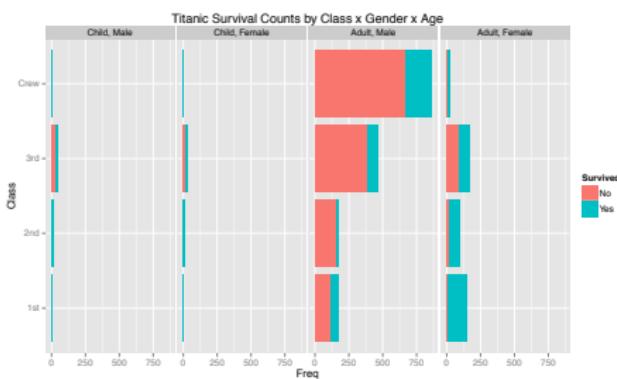
Mosaic plots are similar, but the widths of the bars now reflect proportions:



23 / 1

## Stacked Barplots

Using the `ggplot2` package, we can plot survivals by class, age, and gender all at once.



24 / 1

## Standardized/Normalized Stacked Barplots

Instead of raw counts, we can expand each bar to reflect proportions (i.e. standardize/normalize them).



## Lecture 7: Probability

### Chapter 2.x

1 / 19

## Outcomes

Probability forms the theoretical backbone of statistics. We use probability to characterize randomness.

We often frame probability in terms of a [random process](#) giving rise to an [outcome](#).

Typical examples

- ▶ Die roll: 6 outcomes
- ▶ Coin Flip: 2 outcomes

2 / 19

## Disjoint AKA Mutually Exclusive Outcomes

Two outcomes are **disjoint** (AKA mutually exclusive) if they cannot both occur at the same time.

Die example:

- ▶ Rolling a 1 and a 2 are disjoint.
- ▶ Rolling a 1 and rolling “an odd number” are not disjoint.

3 / 19

## Addition Rule of Probability

If  $A_1$  and  $A_2$  are disjoint outcomes, then

$$P(A_1 \text{ or } A_2) = P(A_1) + P(A_2)$$

Ex: Rolling 1 and 2 are disjoint, so:

$$P(\text{rolling 1 or 2}) = P(\text{rolling 1}) + P(\text{rolling 2}) = \frac{1}{6} + \frac{1}{6}$$

4 / 19

## General Addition Rule of Probability

If  $A_1$  and  $A_2$  are two outcomes (not necessarily disjoint), then

$$P(A_1 \text{ or } A_2) = P(A_1) + P(A_2) - P(A_1 \text{ and } A_2)$$

Venn diagram:

## General Addition Rule of Probability

Events are just combinations of outcomes. Ex: Deck of cards

- ▶  $A_1$  = event we draw a diamond
- ▶  $A_2$  = event we draw a face card

These two events are not disjoint, as there are 3 diamond face cards. Venn diagram:

## General Addition Rule of Probability

$$\begin{aligned} P(A_1 \text{ or } A_2) &= P(\text{diamond or a face card}) \\ &= P(\text{diamond}) + P(\text{face card}) - \\ &\quad P(\text{diamond AND face card}) \\ &= \frac{13}{52} + \frac{3 \times 4}{52} - \frac{3}{52} = \frac{22}{52} = 42.3\% \end{aligned}$$

7 / 19

## Sample Space and the Complement of Events

A die has 6 possible outcomes. The sample space is the set of all possible outcomes  $S = \{1, 2, \dots, 6\}$ .

Say event  $A$  is the event of rolling an even number i.e.  $A = \{2, 4, 6\}$ . The complement of event  $A$  is  $A^c = \{1, 3, 5\}$  i.e. getting an odd number.

Thm

$$P(A) + P(A^c) = 1$$

8 / 19

## Independence

Two processes are **independent** if knowing the outcome of one provides no useful information about the outcome of the other. Otherwise they are dependent.

Consider:

1. Die rolls
2. You get a movie recommendation from your friend Robin, but then their significant other Sam also recommends it.
3. You compare test scores from two Grade 9 students in the same class. Then same school. Then same school district. Then same city. Then same state.

9 / 19

## Independence

We say that events  $A$  and  $B$  are **independent** if

$$P(A \text{ and } B) = P(A) \times P(B)$$

Ex: Dice rolls are independent:

$$\begin{aligned} P(\text{rolling 1 and then 6}) &= P(\text{rolling 1}) \times P(\text{rolling 6}) \\ &= \frac{1}{6} \times \frac{1}{6} = \frac{1}{36} \end{aligned}$$

10 / 19

## Conditional Probability

The conditional probability of an event  $A$  given the event  $B$ , is defined by

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

11 / 19

## Example

Let's suppose I take a random sample of 100 Reed students to study their smoking habits.

	Smoker	Not Smoker	Total
Male	19	41	60
Female	12	28	40
Total	31	69	100

- ▶ What is the probability of a randomly selected male smoking?

$$P(S|M) = \frac{P(S \text{ and } M)}{P(M)} = \frac{19/100}{60/100} = \frac{19}{60}$$

- ▶ What is the probability that a randomly selected smoker is female?

$$P(F|S) = \frac{P(F \text{ and } S)}{P(S)} = \frac{12/100}{31/100} = \frac{12}{31}$$

12 / 19

## Put It Together! Independence and Conditional Prob.

If  $A$  and  $B$  are independent events, then

$$P(A \text{ and } B) = P(A) \times P(B)$$

then

$$P(A|B) = \frac{P(A \text{ and } B)}{P(B)} = \frac{P(A) \times P(B)}{P(B)} = P(A)$$

i.e.  $P(A|B) = P(A)$ : the event  $B$  occurring has no bearing on the probability of  $A$

13 / 19

## Gambler's Fallacy: Roulette



You can bet on individual numbers, sets of numbers, or [red vs black](#). Let's assume no 0 or 00, so that  $P(\text{red}) = P(\text{black}) = \frac{1}{2}$ .

14 / 19

## Gambler's Fallacy: Roulette

One of the biggest cons in casinos: spin history boards.



Let's ignore the numbers and just focus on what color occurred.

Note: the white values on the left are **black** spins.

15 / 19

## Gambler's Fallacy: Roulette

Let's say you look at the board and see that the last 4 spins were **red**.

You will always hear people say "Black is due!"

Ex. on the 5th spin people think:

$$\begin{aligned} P(\text{black}_5 \mid \text{red}_1 \text{ and } \text{red}_2 \text{ and } \text{red}_3 \text{ and } \text{red}_4) &> \\ P(\text{red}_5 \mid \text{red}_1 \text{ and } \text{red}_2 \text{ and } \text{red}_3 \text{ and } \text{red}_4) \end{aligned}$$

16 / 19

## Gambler's Fallacy: Roulette

But assuming the wheel is not rigged, spins are independent i.e.  
 $P(A|B) = P(A)$ . So:

$$P(\text{black}_5 | \text{red}_1 \text{ and } \text{red}_2 \text{ and } \text{red}_3 \text{ and } \text{red}_4) = P(\text{black}_5) = \frac{1}{2}$$

$$P(\text{red}_5 | \text{red}_1 \text{ and } \text{red}_2 \text{ and } \text{red}_3 \text{ and } \text{red}_4) = P(\text{red}_5) = \frac{1}{2}$$

17 / 19

## Next Week's Lab

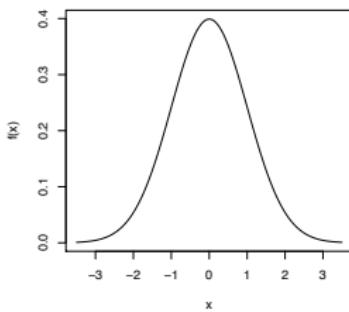
Basketball players who make several baskets in succession are described as having a "hot hand." This refutes the assumption that each shot is **independent** of the next.

We are going to investigate this claim with data from a particular basketball player: Kobe Bryant of the Los Angeles Lakers in the 2009 NBA finals.

18 / 19

## Next Time

Discuss the Normal Distribution



## Lecture 8: Normal Distribution

### Chapter 3.1

1 / 23

## Goals for Today

- ▶ Define the normal distribution in terms of its **parameters**
- ▶ Review:  $\frac{2}{3}$  / 95% / 99.7% rule
- ▶ Standardizing normal observations to **z-scores**

2 / 23

## Normal Distribution

From text page 118:

Many variables are nearly normal, but none are exactly normal. Thus the normal distribution, while not perfect for any single problem, is very useful for a variety of problems.

We will use it in data exploration and to solve important problems in statistics.

3 / 23

## Normal Distribution

Normal distributions:

1. are symmetric
2. are unimodal and bell-shaped
3. have area under the curve 1

4 / 23

## Normal Distribution

A normal curve can be described by two parameters:

- ▶ the mean  $\mu$ . i.e. the center
- ▶ the standard deviation (SD)  $\sigma$ . i.e. the measure of spread

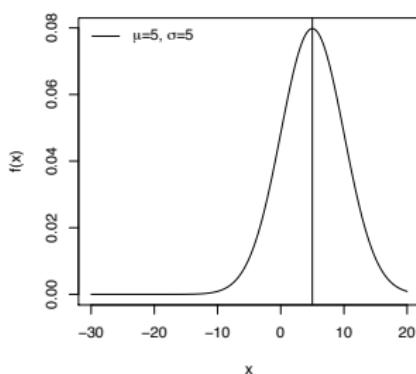
$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

Recall these were the population mean and the population SD.

5 / 23

## Normal Distribution

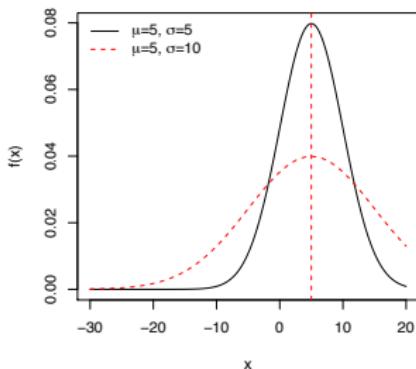
$\mu$  (mean) specifies the center,  $\sigma$  (standard deviation) the spread.



6 / 23

## Normal Example

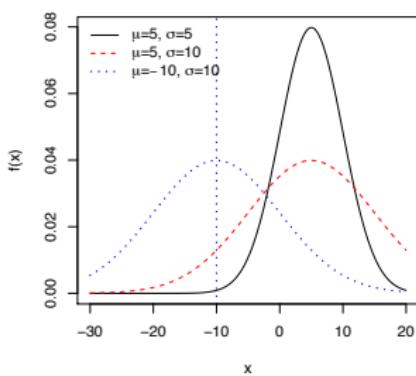
$\mu$  (mean) specifies the center,  $\sigma$  (standard deviation) the spread.



7 / 23

## Normal Example

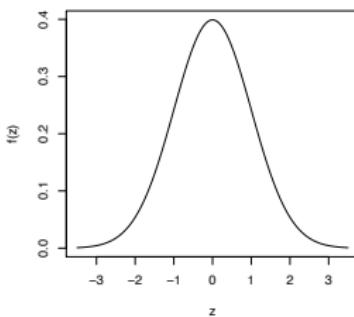
$\mu$  (mean) specifies the center,  $\sigma$  (standard deviation) the spread.



8 / 23

## Standardized Normal Distribution

If  $\mu = 0$  and  $\sigma = 1$ , this is the standard normal distribution:



9 / 23

## Rules of Thumb

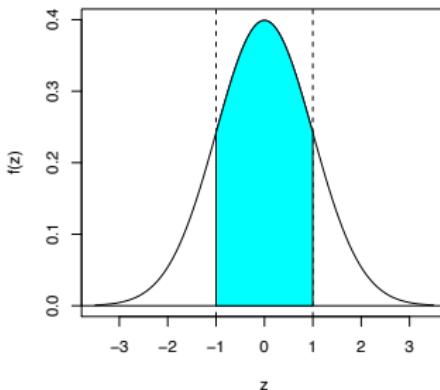
Recall if a distribution is normal, then:

1. Approx.  $\frac{2}{3}$ 's of the data are within  $\pm 1$  SD of the mean
2. Approx. 95% of the data are within  $\pm 2$  SD of the mean
3. Also approx. 99.7% of the data are within  $\pm 3$  SD of the mean

10 / 23

Ex: Standard Normal  $\mu = 0, \sigma = 1$

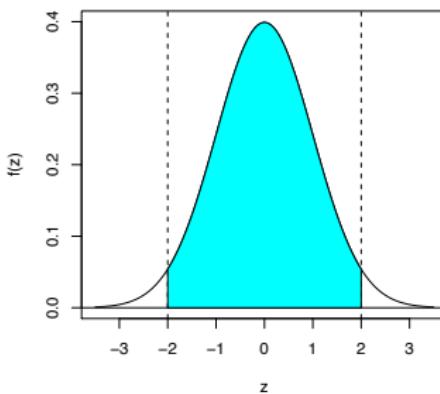
Cyan Area is Two-Thirds



11 / 23

Ex: Standard Normal  $\mu = 0, \sigma = 1$

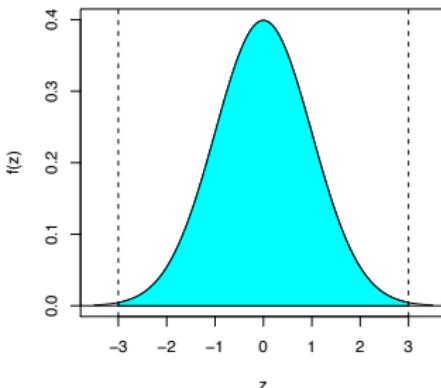
Cyan Area is 95%



12 / 23

Ex: Standard Normal  $\mu = 0, \sigma = 1$

Cyan Area is 99.7%



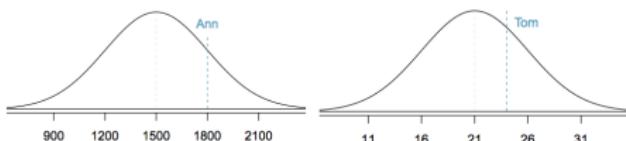
13 / 23

## Motivating Example

From text: Say Ann scores 1800 on the SAT and Tom scores 24 on the ACT. Say both tests scores were normally distributed with:

	SAT	ACT
Mean $\mu$	1500	21
SD $\sigma$	300	5

Question: Who did relatively better?



14 / 23

## z-scores

The **z-score AKA standardized observation** of an observation  $x$  is the number of SD it falls above or below the mean.

The z-score for an observation  $x$  that follows a distribution with mean  $\mu$  and SD  $\sigma$ :

$$z = \frac{x - \mu}{\sigma}$$

15 / 23

## z-scores

Why is the z-score  $z = \frac{x - \mu}{\sigma}$  called the **standardized observation**?

1. The observations are **centered** at  $\mu$ .  
re-center the  $x$  observations to 0 by subtracting  $\mu$ .
2. The observations have **spread**  $\sigma$ .  
re-scale the **spread** of the  $x - \mu$  values to be 1 by dividing by  $\sigma$ .

So we can compare observations from **any** normally distributed data with  $(\mu, \sigma)$

i.e. we've **standardized the observations** to make them comparable.

16 / 23

## Back to Example

- ▶ Ann scored 1800.  $z = \frac{1800 - 1500}{300} = +1$  standard deviation from the mean
- ▶ Tom scored 24.  $z = \frac{24 - 21}{5} = +0.6$  standard deviation from the mean

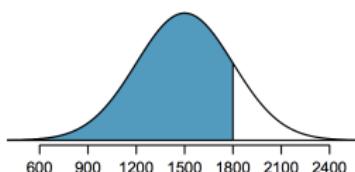
So Ann did relatively better.

17 / 23

## Percentiles

Recall a **percentile** (%'ile) indicates the value below which a given %'age of observations fall below.

**Question:** What %'ile is Ann's SAT score of 1800?  
i.e. what is the blue shaded area?



18 / 23

## Percentiles

Because the total area under the curve is 1, the area to the left of  $z$  represents the %'ile of the observation:



- ▶ The blue shaded area on the left plot will be less than 0.5. We have %'iles less than the 50th %'ile.
- ▶ The blue shaded area on the right plot will be greater than 0.5. We have %'iles greater than the 50th %'ile.

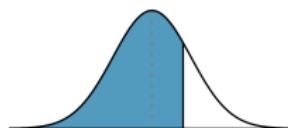
19 / 23

## Normal Probability Table

A **normal probability table** allows you to:

- ▶ identify the %'ile corresponding to a z-score
- ▶ or vice versa: the z-score corresponding to a %'ile

The normal probability tables on page 409 represent z-scores and %'iles corresponding to area to the left:



20 / 23

## Normal Probability Table

Z	Second decimal place of Z									
	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

- ▶ **Red case:** Given a z-score of 0.43. A lookup tells us the area to the left of  $z=0.43$  is 0.6664, i.e. the 66th %'ile
- ▶ **Blue case:** We want the z-score that is the 80th %'ile.  
Reverse lookup: the closest value on the table is 0.7995, i.e. a z-score of 0.84.

21 / 23

## Back to Ann and Tom

- ▶ Since Ann had a z-score of 1.0, her %'ile is 0.8413. (1.0 row, 0.00 column)  
i.e. She did better than 84.13% of SAT test takers.
- ▶ Since Tom had a z-score of 0.6, his %'ile is 0.7257. (0.6 row, 0.00 column)  
i.e. He did better than 72.57% of ACT test takers

22 / 23

## Next Time

Next time we will:

- ▶ Re-iterate the motivation for the normal curve.
- ▶ Go over examples using z-scores.
- ▶ Evaluating the normal approximation.

## Lecture 9: Normal Approximation

### Chapter 3.2

1 / 15

## Goals for Today

- ▶ Discuss how to find %'iles for negative values of  $z$
- ▶ Examples
- ▶ Evaluating how “normal” certain data are.

2 / 15

## Solving Normal Questions

Whenever solving questions of this sort **ALWAYS** draw a rough picture first and keep in mind:

1. The normal distribution/curve is symmetric
2. The total area under the curve is 1

3 / 15

## Normal Probability Tables

Alternatively, whereas

- ▶ table on P.409 gives areas to the left of positive values of  $z$ .
- ▶ table on P.408 gives areas to the left of negative values of  $z$ .

I'm only going to give you P.409 table for exams.

4 / 15

## Speeding on I-5

The distribution of passenger vehicle speeds traveling on Interstate 5 Freeway (I-5) in California is nearly normal with a mean of 72.6 mph and a standard deviation of 4.78 mph.

- a) What percent of passenger vehicles travel slower than 80 mph?
- b) What percent of passenger vehicles travel between 60 and 80 mph?
- c) How fast to do the fastest 5% of passenger vehicles travel?
- d) The speed limit on this stretch of the I-5 is 70 mph.

Approximate what percentage of the passenger vehicles travel above the speed limit on this stretch of the I-5.

## Speeding on I-5

- a) What percent of passenger vehicles travel slower than 80 mph?

## Speeding on I-5

- b) What percent of passenger vehicles travel between 60 and 80 mph?

7 / 15

## Speeding on I-5

- c) How fast do the fastest 5% of passenger vehicles travel?

8 / 15

## Speeding on I-5

- d) The speed limit on this stretch of the I-5 is 70 mph.  
Approximate what percentage of the passenger vehicles travel above the speed limit on this stretch of the I-5.

9 / 15

## Switching Gears: Normal Approximation

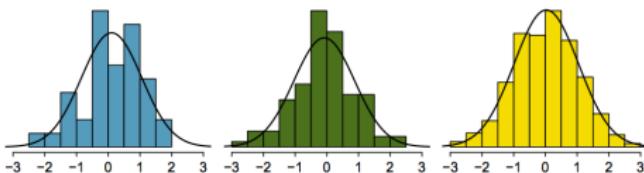
Although we stated that many processes in the physical world look bell-shaped, i.e. roughly normal, we must keep in mind that this is an **approximation**.

**Question:** How do we verify normality?

10 / 15

## Normal Approximation

What about these ones? How well do the histograms fit to the normal curve?



11 / 15

## Normal Probability Plots

Normal probability plots (AKA quantile-quantile plots AKA QQ-plots) are a method for visually displaying how well data fit a normal curve.

The  $k^{\text{th}}$  *q-quantile* is the value such that proportion  $\frac{k}{q}$  of the observations fall below it. So

- ▶ The 4-quantiles are the *quartiles*.
- ▶ The 100-quantiles are the *percentiles*.

12 / 15

## Normal Probability Plots

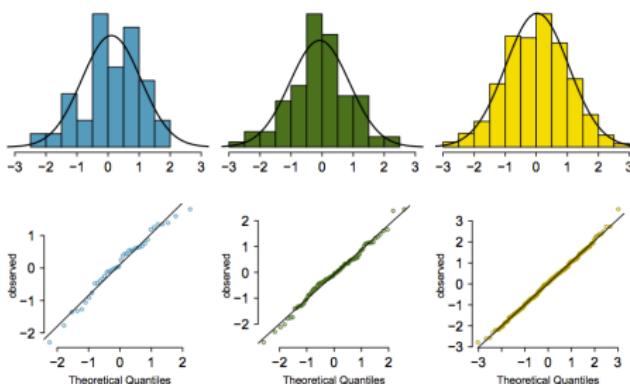
A normal probability plot compares:

- ▶ The **observed** quantiles of a data set (on the **y-axis**)
- ▶ The **theoretical** quantiles that are **exactly** normal (on the **x-axis**)

The more “normal” the data is, the better the fit.

13 / 15

## Normal Probability Plots



14 / 15

## Next Time

- ▶ Introduce some of the more useful other distributions:  
Bernoulli, Geometric, Binomial, and Poisson

## Lecture 10: Bernoulli and Geometric Random Variables

Chapter 3.3-3.5

1 / 1

### Goals for Today

Define

- ▶ Bernoulli random variables
- ▶ Geometric random variables

2 / 1

## Mathematical Definition of a Bernoulli Random Variable

A random variable  $X$  is a random process or variable with a numerical outcome.

Random variables are described in terms of their distribution.

3 / 1

## Bernoulli Distribution

Say we have an experiment where we define each trial (or instance) to have two possible outcomes of interest. Examples

- ▶ Coin flips: heads vs tails
- ▶ Medical test (for a disease): positive vs negative
- ▶ Rolling a die and getting a 6 vs not getting a 6

In each case we can define the outcomes to be success vs failure.  
No moral judgement; just labels.

4 / 1

## Bernoulli Distribution

Say we have trials where we have two outcomes: either a “success” or a “failure”. Classic example: coin flips have  $p = 0.5$  of heads, if we define heads as the success.

- ▶ probability  $p$  of a “success.” Denote successes with a “1.”
- ▶ probability  $1 - p$  of a “failure.” Denote failures with a “0.”

5 / 1

## Definition of a Bernoulli Random Variable

If  $X$  is a random variable that takes value

- ▶ 1 with probability of success  $p$
- ▶ 0 with probability of failure  $1 - p$

then  $X$  is a [Bernoulli random variable](#) with mean and standard deviation:

$$\begin{aligned}\mu &= p \\ \sigma &= \sqrt{p(1-p)}\end{aligned}$$

6 / 1

## Intuition Behind $\sigma$

7 / 1

## Sample Proportion

Say you repeat  $n$  instances of a Bernoulli random variable. You end up with a sample  $x_1, \dots, x_n$

The sample proportion  $\hat{p}$  (p-hat) is the sample mean of these observations. i.e.

$$\hat{p} = \frac{\text{\# of successes}}{\text{\# of trials}} = \frac{1}{n} \sum_{i=1}^n x_i$$

8 / 1

## Example of Bernoulli Distribution

- ▶ A success as rolling a 6.  
So  $P(X = 1) = P(\text{success}) = p = \frac{1}{6}$ .
- ▶ A failure as rolling anything else.  
So  $P(X = 0) = P(\text{failure}) = 1 - p = \frac{5}{6}$ .

9 / 1

## Back to Lecture 3.1: Population vs Sample Values

	True Population Value	Sample Value
Mean	$\mu$	$\bar{x}$
Variance	$\sigma^2$	$s^2$
Standard Deviation	$\sigma$	$s$
Proportion	$p$	$\hat{p}$

The sample proportion  $\hat{p}$  is a specific kind of sample mean for Bernoulli random variables, which estimates  $p$ , a specific kind of population mean.

10 / 1

## Scenario

Question: Say

- ▶ the San Francisco Giants have equal probability  $p = 0.6$  of winning any game
- ▶ games are independent

It's the beginning of the season. What is the probability that they don't win their first game until the 5th game of the season?

For this to happen, there must be 4 loses in the first 4 games AND a win in the 5th game:

$$\begin{aligned} P(\text{1st W in 5th game}) &= P(4 \text{ loses}) \times P(\text{win}) \\ &= (P(\text{loss}))^4 \times P(\text{win}) \\ &= (1 - p)^4 \times p \\ &= 0.4^4 \times 0.6 = 0.01536. \end{aligned}$$

11 / 1

## Geometric Random Variables

**Geometric Distribution:** If the probability of a success in any trial is  $p$ , the trials are independent, then the probability of finding the first success on the  $n^{\text{th}}$  trial is given by

$$(1 - p)^{n-1} p$$

Also

$$\begin{aligned} \mu &= \frac{1}{p} \\ \sigma^2 &= \frac{1-p}{p^2} \\ \sigma &= \frac{\sqrt{1-p}}{p} \end{aligned}$$

12 / 1

## Intuition Behind $\mu$

Think about  $\mu$ :  $\frac{1}{p}$  is the average number of trials we need until the first success.

So compare:

- ▶ Say  $p = 0.5$ . Then  $\mu = \frac{1}{0.5} = 2$
- ▶ Say  $p = 0.001$ . Then  $\mu = \frac{1}{0.001} = 1000$

In the first case, the probability of a success is [lower](#), so we expect on average it will take more trials until the [first](#) success.

13 / 1

## Yesterday's Quiz: Placebos

[Question 1](#): Was Dr. Irving Kirsch arguing that anti-depressants are no better than placebos for everyone with depression?

[Solution](#): No, while he argued that anti-depressants were no better than placebo for those with mild to moderate depression, he is of the opinion that there is clinical benefit for those who are severely depressed.

14 / 1

## Yesterday's Quiz: Placebos

**Question 2:** What is Dr. Walter Brown's (bald guy from Yale) criticism of the way the FDA approves anti-depressants?

**Solution:** That all that is required are two clinical trials where the drug performs better than placebo, regardless of the number of trials with "negative results." Ex: say a drug performs better than placebo in 2 trials, but fails in 998 trials, it will still be approved by the FDA.

## Lecture 11: Binomial and Poisson Random Variables

Chapter 3.3-3.5

1 / 15

### Goals for Today

Define

- ▶ Binomial random variables
- ▶ Poisson random variables

2 / 15

## Binomial Distribution

So say now, instead of  $P(\text{1st W in 5th game}) = P(\text{LLLLW})$ , we want the probability that they win **exactly one** out of the five games. Five ways:

Pattern	Probability	Equals
WLLLL	$p \times (1-p)^4$	$= p \times (1-p)^4$
LWLLL	$(1-p) \times p \times (1-p)^3$	$= p \times (1-p)^4$
LLWLL	$(1-p)^2 \times p \times (1-p)^2$	$= p \times (1-p)^4$
LLLWL	$(1-p)^3 \times p \times (1-p)$	$= p \times (1-p)^4$
LLLLW	$(1-p)^4 \times p$	$= p \times (1-p)^4$

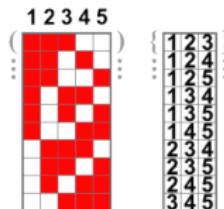
Each pattern (book calls it scenario) has the same probability regardless of order by independence, and there are 5 ways to **choose** the pattern.

So  $P(\text{win exactly one out of five})$  is  
 $5 \times p \times (1-p)^4 = 5 \times 0.4^4 \times 0.6 = 0.0768$

3 / 15

## Step Back... Example of $n$ choose $k$

Say I give you  $n = 5$  balls labeled 1 thru 5. How many different ways can you choose  $k = 3$  of them?



As we see, 10 ways.

4 / 15

## Step Back... $n$ choose $k$ in General

Say I give you  $n$  balls labeled 1 thru  $n$ . How many different ways can you choose  $k$  of them?

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

This is read  $n$  choose  $k$ .

In example:  $n = 5$  and  $k = 3$

$$\binom{5}{3} = \frac{5!}{3!(5-3)!} = \frac{5 \times 4 \times 3 \times 2 \times 1}{(3 \times 2 \times 1)(2 \times 1)} = \frac{120}{12} = 10$$

Note that  $0! = 1$

5 / 15

## Binomial Distribution

Suppose the probability of a single trial being a success is  $p$ . Then the probability of observing exactly  $k$  successes in  $n$  independent trials is given by:

$$\begin{aligned} P(\text{exactly } k \text{ successes}) &= \binom{n}{k} p^k (1-p)^{n-k} \\ &= \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k} \end{aligned}$$

The mean, variance, and SD are:

$$\mu = np \quad \sigma^2 = np(1-p) \quad \sigma = \sqrt{np(1-p)}$$

6 / 15

## Conditions for Binomial Distribution

1. The trials are independent.
2. The number of trials  $n$  is fixed
3. Each trial outcome can be classified as a failure or a success
4. The probability of a success  $p$  is the same for each trial

7 / 15

## Back to Soccer Example

The Portland Timbers have equal probability  $p = 0.6$  of winning any particular soccer game. We want the probability that they win exactly one out of the five games. Five ways:

Pattern	Probability	Equals
WLPLL	$p \times (1 - p)^4$	$= p \times (1 - p)^4$
LWLPL	$(1 - p) \times p \times (1 - p)^3$	$= p \times (1 - p)^4$
LLWLL	$(1 - p)^2 \times p \times (1 - p)^2$	$= p \times (1 - p)^4$
LLLWL	$(1 - p)^3 \times p \times (1 - p)$	$= p \times (1 - p)^4$
LLLLW	$(1 - p)^4 \times p$	$= p \times (1 - p)^4$

Letting a win be a “success”:

$$\begin{aligned}P(k = 1 \text{ win}) &= \binom{n}{k} p^k (1-p)^{n-k} = \frac{5!}{1! \times 4!} 0.6 \times 0.4^4 \\&= 5 \times 0.6 \times 0.4^4 = 0.0768\end{aligned}$$

8 / 15

## Back to Soccer Example

What about the probability that they win all their games! i.e.

$k = 5$ :

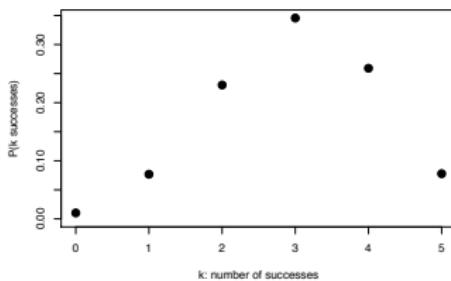
$$\begin{aligned} P(k = 5 \text{ wins}) &= \binom{n}{k} p^k (1-p)^{n-k} = \binom{5}{5} 0.6^5 (1-0.6)^0 \\ &= \frac{5!}{5! \times 0!} 0.6^5 \times 1 = 0.08 \end{aligned}$$

What about the probability that they at least one game?

$$\begin{aligned} P(\text{at least } k = 1 \text{ wins}) &= P(k = 1 \text{ win}) + \dots + P(k = 5 \text{ wins}) \\ &= 1 - P(k=0 \text{ wins}) \\ &= 1 - \frac{5!}{0! \times 5!} 0.6^0 \times 0.4^5 = 1 - 0.01024 \\ &= 0.98976 \end{aligned}$$

9 / 15

## Back to Soccer Example



10 / 15

## Poisson Distribution

Say you want to count the number of rare events in a large population over a unit of time. Examples:

- ▶ the number of car accidents at a particular intersection on a given week
- ▶ the number of ambulance calls on any given day in Portland
- ▶ the number of soldiers in the Prussian army killed accidentally by horse kick from 1875 to 1894

The Poisson distribution helps us describe the number of such events that will occur in a short unit of time for a fixed population if the individuals within the population are independent.

11 / 15

## Poisson Distribution

Suppose we are watching for rare events and the number of observed events follows a Poisson distribution with rate  $\lambda$

$$P(\text{observe } k \text{ rare events}) = \frac{\lambda^k e^{-\lambda}}{k!}$$

where  $k$  may take a value 0, 1, 2, ... where  $e \approx 2.718$ .

The mean and SD are  $\lambda$  and  $\sqrt{\lambda}$ .

12 / 15

## Conditions for Poisson Distribution

A random variable **may** be Poisson distributed if

1. The event in question is rare
2. The population is large
3. The events occur independently of each other

13 / 15

## Exercise 3.47 on Page 158

A coffee shop serves an average of 75 customers per hour during the morning rush. What is the probability that the coffee shop serves 70 customers in one hour during this time of the day?

In this case,  $\lambda = 75$  is the rate

$$P(k = 70) = \frac{75^{70} e^{-75}}{70!} = 0.040$$

Type `dpois(x=70, lambda=75)` in R

14 / 15

## Next Time

### Chapter 4: Foundations for Inference

- ▶ Variability in estimates  $\bar{x}$ ,  $\hat{p}$ , etc.
- ▶ In fact, we can associate a [distribution](#) to these estimates

## Lecture 12: Sampling Distributions & Standard Errors

### Chapter 4.1

1 / 19

## Goals for Today

Start Chapter 4: Arguably the most important chapter as it goes to the heart of what statistical inference is. Three important definitions today:

1. point estimate
2. sampling distribution
3. standard error

2 / 19

## Point Estimates

**Definition 1:** Point estimates are functions of a random sample of  $n$  observations  $x_1, \dots, x_n$ . They estimate the value of some unknown population parameter.

Ex: the sample mean

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + \dots + x_n}{n}$$

is a point estimate of the true population mean  $\mu$

3 / 19

## Behavior of Point Estimates

Ex: Say we draw a random sample of size  $n = 100$  from a large population that is normally distributed with  $\mu = 5$  and  $\sigma = 2$ .

**Two Important Questions:**

1. Is  $\bar{x}$  going to be exactly 5?
2. Say we get  $\bar{x} = 5.025$ . If we repeat this procedure: i.e. generate a new sample of size  $n = 100$  and compute  $\bar{x}$ ), will we get  $\bar{x} = 5.025$ ?

We need to characterize this random error.

4 / 19

## Behavior of Point Estimates

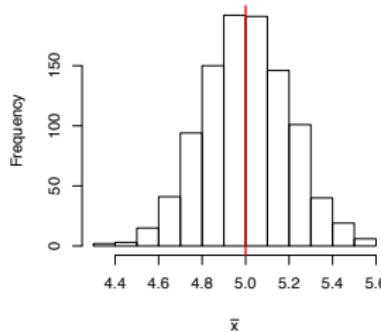
Let's repeat this procedure, say, 1000 times:

1st time	We get $\bar{x} = 4.831$
2nd time	We get $\bar{x} = 5.104$
3rd time	We get $\bar{x} = 4.965$
...	
1000th time	We get $\bar{x} = 4.957$

5 / 19

## Sampling Distribution

This histogram is the 1000 instances of  $\bar{x}$ , where each  $\bar{x}$  is based on a sample of  $n = 100$ . This is the [sampling distribution](#) of  $\bar{x}$ :



6 / 19

## Sampling Distributions

**Definition 2:** the **sampling distribution** is the distribution of point estimates based on samples of fixed size  $n$ .

Every instance of a point estimate can be thought of as a draw from the sampling distribution.

If the sampling is **representative** (unbiased) then the sampling distribution will be centered around the true population parameter (in our case  $\mu$ ).

7 / 19

## Sampling Distributions

We can define the sampling distributions for **any** point estimate, not just  $\bar{x}$ :

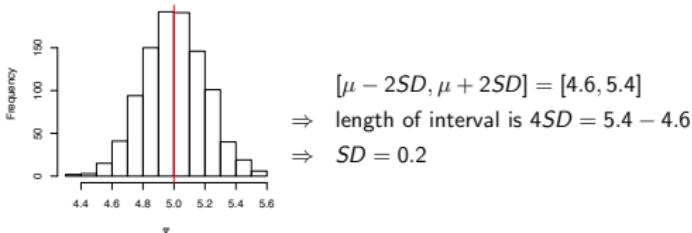
- ▶  $s$
- ▶ the sample median
- ▶ etc.

We will only focus on sample means, including the sample proportion  $\hat{p}$ .

8 / 19

## Measure of Spread

What about spread?  $[4.6, 5.4]$  contains roughly 95% of the data.



9 / 19

## Standard Errors

**Definition 3:** The **standard error** is the standard deviation of the sampling distribution of a point estimate.

It describes the uncertainty/variability associated with the point estimate. In other words, the “typical” error.

**Confusing:** the **standard error** is a specific kind of standard deviation.

10 / 19

## Standard Error of $\bar{x}$

Given  $n$  independent observations from a population with standard deviation  $\sigma$ , the standard error of the sample mean is

$$SE = \frac{\sigma}{\sqrt{n}}$$

**Rule of thumb for independence:** You need a simple random sample consisting of less than 10% of the population.

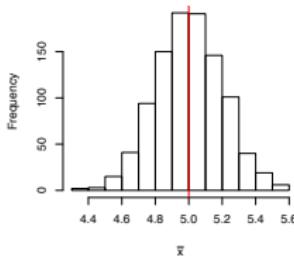
**Notice:**  $\sqrt{n}$  in the denominator: as  $n$  increases, SE decreases! This is why sample size matters.

11 / 19

## Back to Histogram

Samples were of size  $n = 100$  with  $\sigma = 2$ . We estimated that the SD of the sampling distribution was 0.2. Using the formula:

$$SE = \frac{\sigma}{\sqrt{n}} = \frac{2}{\sqrt{100}} = \frac{2}{10} = 0.2$$

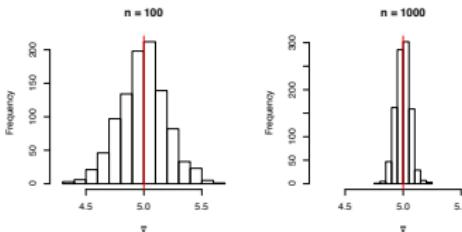


12 / 19

## Standard Error of the Sample Mean $\bar{x}$

Compare 1000 instances of  $\bar{x}$  when

- ▶  $n = 100$ .  $SE = \frac{\sigma}{\sqrt{n}} = \frac{2}{\sqrt{100}} = 0.2$
- ▶  $n = 1000$ .  $SE = \frac{\sigma}{\sqrt{n}} = \frac{2}{\sqrt{1000}} = 0.0632$ . **Smaller!**



Both are “accurate”, but the estimates on the right are “more precise.”

13 / 19

## Repeated Sampling

**Popular question:** What's up with this “1000” instances? Why would you take 1000 different samples of size  $n$ ?

**Answer:** No, in practice you would **not** sample repeatedly: you do this only **once** for the largest  $n$  possible.

Rather the 1000 instances of  $\bar{x}$  is a theoretical exercise to illustrate that  $\bar{x}$ 's are random and we characterize its randomness by its sampling distribution and its standard error.

14 / 19

## Standard Error of the Sample Mean

In this example we knew  $\sigma$ ; typically we won't. However, when

- ▶  $n \geq 30$
- ▶ the distribution of the population is **not** strongly skewed

we can use the point estimate of  $\sigma$ . i.e. plug in  $s$  in place of  $\sigma$ :

$$SE = \frac{s}{\sqrt{n}}$$

15 / 19

## Example

Say in you take a simple random sample of 100 runners in a race and you are interested in their ages:

- ▶  $\bar{x} = 35.05$
- ▶  $s = 8.97$

Assuming that the 100 runners consist of less than 10% of the population, the standard error of  $\bar{x}$  is

$$SE = \frac{s}{\sqrt{100}} = \frac{8.97}{10} = 0.897$$

16 / 19

## Population Distribution vs Sampling Distribution

17 / 19

## Recap

- ▶ **Point estimates** are based on a sample  $x_1, \dots, x_n$  and are used to estimate population parameters.
- ▶ The **sampling distribution** characterizes the (random) behavior of point estimates.
- ▶ The standard deviation of a sampling distribution is the **standard error**: it quantifies the uncertainty/variability of point estimates.

18 / 19

## Next Time

- ▶ Confidence Intervals
- ▶ When quoting survey results, what does: "the results of this survey are estimated to be accurate within 3.1 percentage points, 19 times out of 20" mean?
- ▶ **Big One:** Central Limit Theorem

## Lecture 13: Central Limit Theorem + Confidence Intervals

Chapter 4.4 + 4.2

1 / 25

### Goals for Today

- ▶ Discuss the Central Limit Theorem
- ▶ Introduce confidence intervals
- ▶ Interpretation

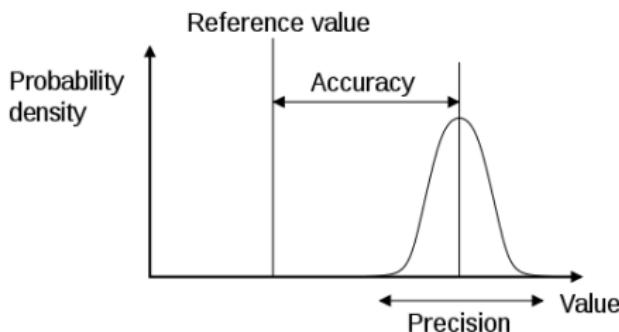
2 / 25

## Recap

- ▶ Point estimates are based on a sample  $x_1, \dots, x_n$  and are used to estimate population parameters.
- ▶ The sampling distribution characterizes the (random) behavior of point estimates (like  $\bar{x}$ ).
- ▶ The standard deviation of a sampling distribution is the standard error: it quantifies the uncertainty/variability of point estimates.

3 / 25

## Illustrative Image of Sampling Distribution



4 / 25

## Central Limit Theorem

### Central Limit Theorem



The averages of samples have approximately normal distributions

Sample size → Bigger  
Distribution of Averages → more normal and narrower

5 / 25

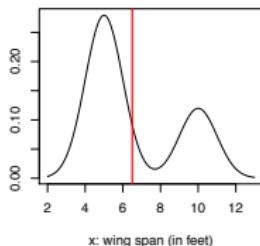
## Central Limit Theorem

Question: Why do we care about the CLT?

Answer: We want the sampling distribution of  $\bar{x}$  to be Normal regardless of the shape of population distribution.

Example: The bimodal (population) distribution of dragon wing spans has a mean of 6.5:

Population Dist'n



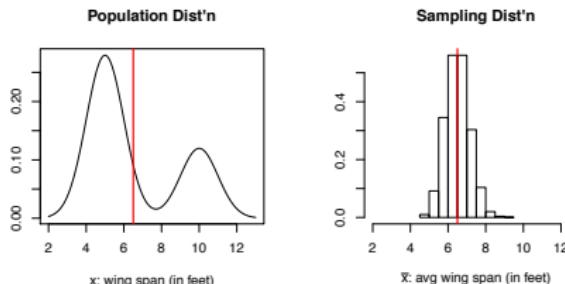
6 / 25

## Central Limit Theorem

**Question:** Why do we care about the CLT?

**Answer:** We want the sampling distribution of  $\bar{x}$  to be Normal regardless of the shape of population distribution.

**Example:** The bimodal (population) distribution of dragon wing spans has a mean of 6.5:



7 / 25

## Central Limit Theorem

**Question:** Why do we care that the sampling distribution of  $\bar{x}$  is Normal?

**Answer:** So we can use the Normal table on p.409 of the book to calculate areas/percentiles/probabilities! We call this using the **normal model**.

Z	Second decimal place of Z									
	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6369	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8434	0.8461	0.8484	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8664	0.8688	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
:	:	:	:	:	:	:	:	:	:	

8 / 25

## Definition

For a sample  $x_1, \dots, x_n$  of independent observations, if  $n$  is "large" enough to counteract the skew of the population distribution, then the sampling distribution of  $\bar{x}$  is approximately Normal with

- ▶ mean  $\mu$
- ▶ SD equal to the  $SE = \frac{\sigma}{\sqrt{n}}$

**Key:** this holds for any population distribution, not just a normally distributed population.

**Recall:** If we don't know  $\sigma$ , we can plug in its point estimate  $s$  if the two conditions are satisfied.

9 / 25

## Conditions for the Normal Model

This translates to the following conditions to verify to be able to use the Normal model with  $s$  in place of  $\sigma$ , as stated in the book:

1.  $n \leq 10\%$  of the population size.

Comment: To ensure independence.

2.  $n \geq 30$ .

Comment: This is a **rule of thumb** that works for most cases.  
You might need less, you might need more.

3. The population distribution is not strongly skewed.

Comment: This is related 2. The larger the  $n$ , the more lenient we can be with the skew assumption.

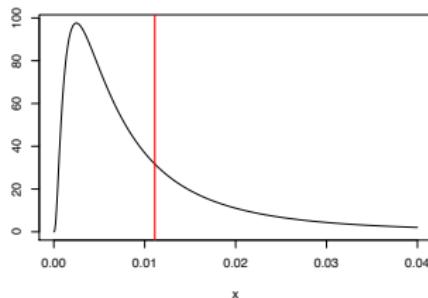
To verify this we can either:

- ▶ Look at the histogram of the sample  $x_1, \dots, x_n$
- ▶ Assume this based on knowledge/previous research

10 / 25

## Example of Skew vs $n$

Let's say your observations come from the following very skewed population distribution with mean  $\mu = 0.011109$ .

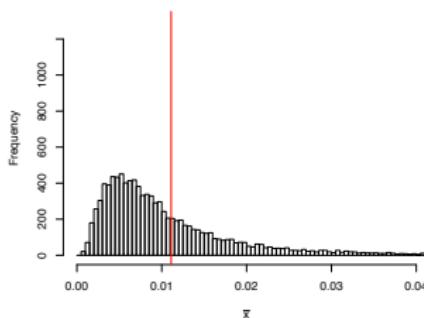


This is where your individual observations  $x_i$  come from. Now compare 10000 values of  $\bar{x}$ 's based on different  $n$ : 2, 10, 30, 75.

11 / 25

## Example of Skew vs $n$

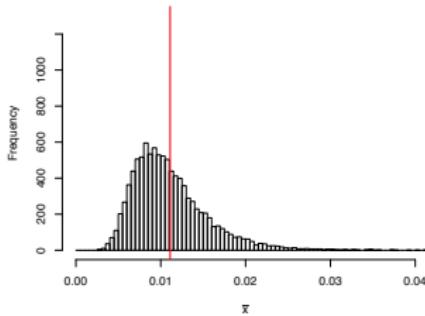
For 10000 values of  $\bar{x}$  based on samples of size  $n = 2$ , the sampling distribution is:



12 / 25

## Example of Skew vs $n$

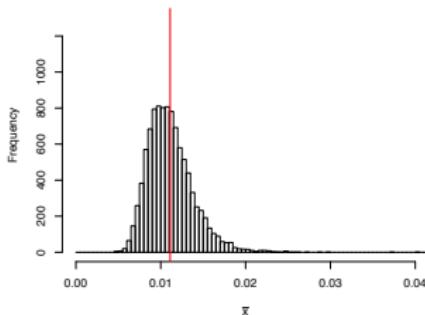
For 10000 values of  $\bar{X}$  based on samples of size  $n = 10$ , the sampling distribution is:



13 / 25

## Example of Skew vs $n$

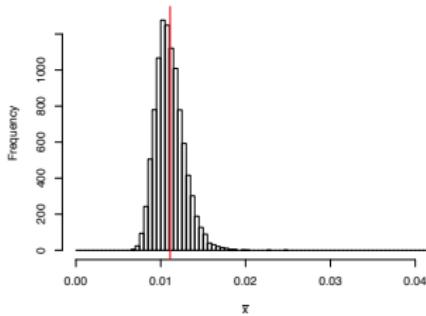
For 10000 values of  $\bar{X}$  based on samples of size  $n = 30$ , the sampling distribution is:



14 / 25

## Example of Skew vs $n$

For 10000 values of  $\bar{x}$  based on samples of size  $n = 75$ , the sampling distribution is:



i.e. more normal and more narrow

15 / 25

## Intuition of a Confidence Interval

**Our Goal:** we want estimate a population parameter (e.g.  $\mu$ ).  
Analogy: imagine  $\mu$  is a fish in a murky river that we want to capture:

Using just the point estimate:



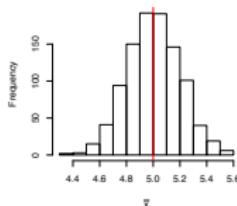
Using a confidence interval:



16 / 25

## Intuition of a Confidence Interval

Recall the example of 1000 instances of  $\bar{x}$  based on  $n = 100$ . Each observation came from a population distribution that was Normal with  $\mu = 5$  &  $\sigma = 2$ .



We observed the sampling distribution

- ▶ is centered at  $\mu$
- ▶ has spread  $SE = \frac{\sigma}{\sqrt{n}} = \frac{2}{\sqrt{100}} = 0.2$

17 / 25

## Intuition of a Confidence Interval

A plausible range of values for the population parameter is called a **confidence interval (CI)**. Since

- ▶ the SE is the standard deviation of the sampling distribution
- ▶ roughly 95% of the time  $\bar{x}$  will be within 2 SE of  $\mu$  **if the sampling distribution is normal**

If the interval spreads out 2 SE from  $\bar{x}$ , we can be roughly "95% confident" that we have captured the true parameter  $\mu$ .

18 / 25

## Intuition of a Confidence Interval

A 95% confidence interval for  $\mu$  is (no more using rule of thumb  $2 \times SD$ ):

$$\begin{aligned}\bar{x} \pm 1.96SE &= [\bar{x} - 1.96SE, \bar{x} + 1.96SE] \\ &= \left[ \bar{x} - 1.96 \frac{\sigma}{\sqrt{n}}, \bar{x} + 1.96 \frac{\sigma}{\sqrt{n}} \right]\end{aligned}$$

If we don't know  $\sigma$ , assuming the conditions hold, plug in  $s$

$$\bar{x} \pm 1.96 \frac{s}{\sqrt{n}} = \left[ \bar{x} - 1.96 \frac{s}{\sqrt{n}}, \bar{x} + 1.96 \frac{s}{\sqrt{n}} \right]$$

19 / 25

## Confidence Intervals

In general a confidence interval for  $\mu$  will be

$$\bar{x} \pm z^*SE = [\bar{x} - z^*SE, \bar{x} + z^*SE]$$

where the **critical value**  $z^*$  is chosen to achieve the desired confidence.

Ex: For 95% confidence  $z^* = 1.96$ . For 99% confidence  $z^* = 2.58$

20 / 25

## Crucial: How to Interpret a Confidence Interval

The confidence interval has nothing to say about any particular calculated interval; it only pertains to the **method** used to construct the interval:

- ▶ **Wrong, yet common, interpretation:** There is a 95% chance that the C.I. captures the true population mean  $\mu$ . The probability is 0 or 1: either it does or it doesn't.
- ▶ **Correct, interpretation:** If we were to repeat this sampling procedure 100 times, we expect 95 (i.e. 95%) of calculated C.I.'s to capture the true  $\mu$

21 / 25

## Illustration: How to Interpret a Confidence Interval

In Chapter 4 there is an example of finish times (in minutes) from the 2012 Cherry Blossom 10 mile run with  $n = 16,924$  participants. In this case, we can compute the **true** population mean  $\mu = 94.52$ .

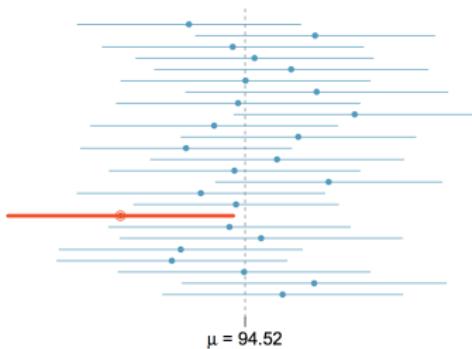
Say we take 25 (random) samples of size  $n = 100$  and for each sample we compute:

- ▶  $\bar{x}$
- ▶  $s$
- ▶ and hence the 95% CI:  $\left[ \bar{x} - 1.96 \times \frac{s}{\sqrt{n}}, \bar{x} + 1.96 \times \frac{s}{\sqrt{n}} \right]$

22 / 25

## How to Interpret a Confidence Interval

Of the 25 CI's based on 25 different samples of size  $n = 100$ , one of them (in red) did not capture the true population mean  $\mu$ :



23 / 25

## Political Polls

We polled the electorate and found that 45% of voters plan to vote for candidate X. The margin of error for this poll is  $\pm 3.4$  percentage points 19 times out of 20.

What does this mean?

- ▶ "19 times out of 20" indicates 95%
- ▶ The margin of error of  $\pm 3.4\%$  indicates that 95% C.I. is:

$$45 \pm 3.4\% = [41.6, 48.4]$$

**Intrepretation:** the interpretation is not that there is a 95% chance that [41.6, 48.4] captures the true %'age. Rather, that if we were to take 20 such polls, 19 of them would capture the true %'age.

24 / 25

## Next Time

Hypothesis Testing: we can perform [statistical tests](#) on population parameters such as  $\mu$ :

Define:

- ▶ Null and alternative hypotheses.
- ▶ Testing hypotheses using confidence intervals.
- ▶ Types of errors

## Lecture 14: Hypothesis Testing Part I

### Chapter 4.3

1 / 17

## Goals for Today

- ▶ Introduce Hypothesis Testing Framework
- ▶ Testing Hypotheses Using Confidence Intervals
- ▶ Types of Errors
- ▶ Testing Hypotheses Using p-Values

2 / 17

## Statistical Hypothesis Testing

(For now) A **hypothesis** is a claim about a population parameter.

A **hypothesis test** is a method for using sample data to decide between two competing hypotheses about the population parameter:

- ▶ A **null hypothesis**  $H_0$ .  
i.e. the **status quo** that is initially assumed to be true, but will be tested.
- ▶ An **alternative hypothesis**  $H_A$ .  
i.e. the **challenger**.

3 / 17

## Example

We flip a coin many times and start to suspect that it is biased:

- ▶  $H_0$ : the coin is fair. i.e. the probability of heads is  $p = 0.5$
- ▶  $H_A$ : the coin is not fair. i.e.  $p \neq 0.5$

4 / 17

## Crucial Concept: Conclusions of Hypothesis Tests

There are two potential outcomes of a hypothesis test. Either we

- ▶ reject  $H_0$  in favor of  $H_A$
- ▶ fail to reject  $H_0$

Note the difference between accepting  $H_0$  & failing to reject  $H_0$

- ▶ “accepting  $H_0$ ” is saying we are sure  $H_0$  is true
- ▶ “failing to reject  $H_0$ ” is saying something not as strong: we do not have enough evidence to reject  $H_0$ .

5 / 17

## Analogy: US Criminal Justice System

In the criminal justice system, the jury's verdict does NOT make any statement about the defendant being **innocent**, rather that there was not enough evidence to prove beyond a reasonable doubt that they were guilty.

6 / 17

## Analogy: US Criminal Justice System

Let's compare criminal trials to hypothesis tests:

### Truth:

- ▶ Truth about the defendant: innocent vs guilty
- ▶ Truth about the hypothesis:  $H_0$  or  $H_A$

### Decision:

- ▶ Verdict: not guilty vs guilty
- ▶ Test outcome: "Do not reject  $H_0$ " vs "Reject  $H_0$ "

7 / 17

## Testing Hypotheses Using Confidence Intervals

Example on page 173: The average 10 mile run time for the Cherry Blossom Run in 2006  $\mu_{2006}$  was 93.29 min. Researchers suspect  $\mu_{2012}$  was different:

- ▶  $H_0$ : average time was the same. i.e.  $\mu_{2012} = 93.29$
- ▶  $H_A$ : average time was different. i.e.  $\mu_{2012} \neq 93.29$

8 / 17

## Testing Hypotheses Using Confidence Intervals

9 / 17

## Decision Errors

Hypothesis tests will get things right sometimes and wrong sometimes:

		Test conclusion	
		do not reject $H_0$	reject $H_0$ in favor of $H_A$
Truth	$H_0$ true	OK	Type I Error
	$H_A$ true	Type II Error	OK

Two kinds of errors:

- ▶ Type I Error: a false positive
- ▶ Type II Error: a false negative

10 / 17

## Decision Errors

- ▶ Trade-off between these two error rates
  - ▶ procedures with lower type I error rates typically have higher type II error rates
  - ▶ vice-versa
- ▶ In other words, there is almost never a procedure that makes no type I errors and no type II errors. Some sort of balance between the two is required

11 / 17

## Next Time

- ▶ More Hypothesis Testing

12 / 17

## Lecture 15: Hypothesis Testing Part II

### Chapter 4.3

1 / 21

## Previously... Statistical Hypothesis Testing

A **hypothesis test** is a method for using sample data to decide between two competing hypotheses about the population parameter:

- ▶ A **null hypothesis**  $H_0$ .  
i.e. the **status quo** that is initially assumed to be true, but will be tested.
- ▶ An **alternative hypothesis**  $H_A$ . i.e. the **challenger**.

There are two potential outcomes of a hypothesis test. Either we

- ▶ reject  $H_0$
- ▶ fail to reject  $H_0$

2 / 21

## Previously... Decision Errors

Hypothesis tests will get things right sometimes and wrong sometimes:

		Test conclusion	
		do not reject $H_0$	reject $H_0$ in favor of $H_A$
Truth	$H_0$ true	OK	Type I Error
	$H_A$ true	Type II Error	OK

Two kinds of errors:

- ▶ Type I Error: a false positive (test result)
- ▶ Type II Error: a false negative (test result)

3 / 21

## Type I Errors: US Criminal Justice System

Defendants must be proven “guilty beyond a reasonable doubt”: in theory they would rather let a guilty person go free, than put an innocent person in jail.

- ▶  $H_0$ : the defendant is innocent
- ▶  $H_A$ : the defendant is guilty

thus “rejecting  $H_0$ ” is a guilty verdict  $\Rightarrow$  putting them in jail

In this case:

- ▶ Type I error is putting an innocent person in jail (considered worse)
- ▶ Type II error is letting a guilty person go free.

4 / 21

## Type II Errors: Airport Screening

An example of where Type II errors are more serious: [airport screening](#).

$H_0$  : passenger X does not have a weapon

$H_A$  : passenger X has a weapon

Failing to reject  $H_0$  when  $H_A$  is true is not “patting down” passenger X when they have a weapon.

Hence the long lines at airport security.

5 / 21

## Goals for Today

- ▶ Define significance level
- ▶ Tie-in p-Values with sampling distributions
- ▶ Example

6 / 21

## Significance Level

Hypothesis testing is built around rejecting or failing to reject the null hypothesis.

i.e. we do not reject  $H_0$  unless we have **strong evidence**.

As a rule of thumb, when  $H_0$  is true, we do not want to incorrectly reject  $H_0$  more than 5% of the time.

i.e.  $\alpha = 0.05 = 5\%$  is the **significance level**.

With 95% confidence intervals from earlier, we expect it to miss the true population parameter 5% of the time. This corresponds to  $\alpha = 0.05$ .

7 / 21

## Thought experiment: p-Values

Say you flip a coin you think is fair 1000 times. Say you observe

- ▶ 501 heads? Do you think the coin is biased?
- ▶ 525 heads? Do you think the coin is biased?
- ▶ 900 heads? Do you think the coin is biased?

8 / 21

## Thought experiment: p-Values

Intuitively, a **p-value** quantifies how **extreme** an observation is given the null hypothesis.

The smaller the p-value, the more **extreme** the observation, where the meaning of extreme depends on the context.

Note the p-value is different than the population proportion  $p$  (bad historical choice).

9 / 21

## p-Values

Definition: The **p-value** or **observed significance level** is the probability of observing a test statistic as extreme or more extreme (in favor of the alternative) as the one observed, assuming  $H_0$  is true.

It is **NOT** the probability of  $H_0$  being true. This is the most common misinterpretation of the p-value.

10 / 21

## Recall our Coin Example

You have a coin that test for fairness with  $n = 1000$  flips. Set  $p_0 = 0.5$  and define a “success” as getting heads. i.e.

$$H_0 : p = p_0 \text{ i.e. coin is fair}$$

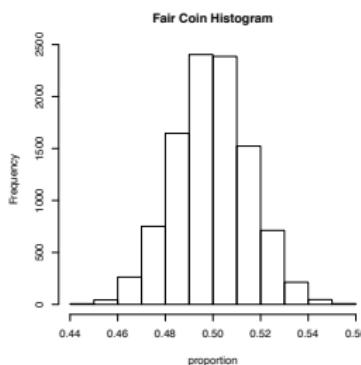
vs       $H_A : p \neq p_0$

- ▶ The point estimate  $\hat{p}$  of  $p$  is  $\frac{\# \text{ of successes}}{\# \text{ of trials}}$ .
- ▶ Since it is based on a sample,  $\hat{p}$  has a sampling distribution
- ▶ The standard error is  $\sqrt{\frac{p(1-p)}{n}}$  (Chapter 6).
- ▶ Furthermore, since conditions hold, the sampling distribution is Normal (CLT)

11 / 21

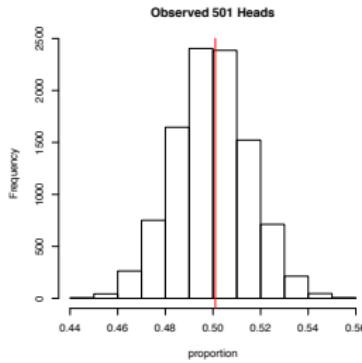
## Sampling Distribution of $\hat{p}$

Under  $H_0$  that the coin is fair i.e.  $p = p_0 = 0.5$ , the sampling distribution of  $\hat{p}$  when  $n = 1000$  is:



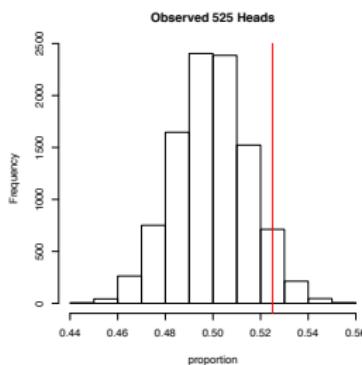
12 / 21

Say we observe...



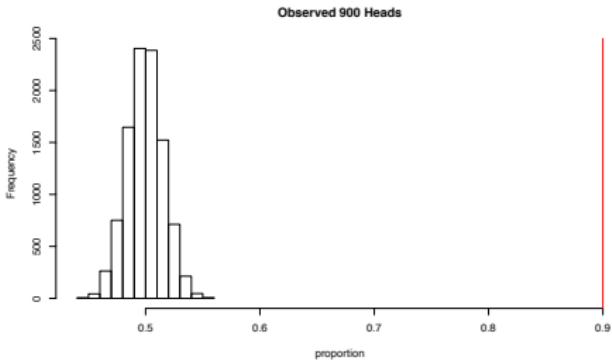
13 / 21

Say we observe...



14 / 21

Say we observe...



15 / 21

## p-Value Definition

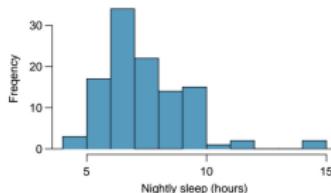
The **p-value** or *observed significance level* is the probability of observing a test statistic as extreme or more extreme (in favor of the alternative) as the one observed, assuming  $H_0$  is true.

It is **NOT** the probability of  $H_0$  being true. This is the most common misinterpretation of the p-value.

16 / 21

## Example about Sleep Habits

A poll found that college students sleep about 7 hours a night. Researchers suspect that Reedies sleep more. They want to investigate this claim at a pre-specified  $\alpha = 0.05$  level. They sample  $n = 110$  Reedies and find that  $\bar{x} = 7.42$  and  $s = 1.75$  and the histogram looks like:

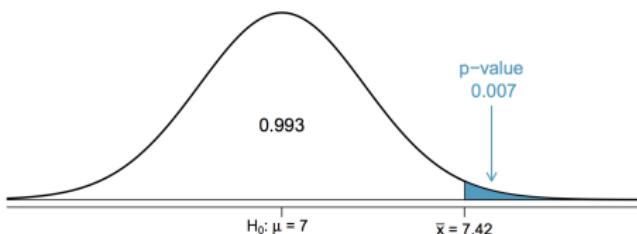


17 / 21

## Example about Sleep Habits

In our case, since  $H_A : \mu > 7$ , more extreme means to the right of  $z = 2.47$ .

Hence, the p-value is 0.007:



18 / 21

## Example about Sleep Habits

Since the p-value  $0.007 < 0.05 = \alpha$ , the pre-specified significance level, it has a high degree of extremeness, and thus we reject  $H_0$ .

**Interpretation:** we reject (at the  $\alpha = 0.05$  significance level) the hypothesis that the average # of hours of Reedies sleep is 7, in favor of the hypothesis that sleep more.

19 / 21

## Example about Sleep Habits

**Correct interpretation of the p-value:** If the null hypothesis is true ( $\mu = 7$ ), the probability of observing a sample mean  $\bar{x} = 7.42$  or greater is 0.007.

**Incorrect interpretation of the p-value:** The probability that the null hypothesis ( $\mu = 7$ ) is true is 0.007.

20 / 21

## Next Time

- ▶ How big a sample size do I need? i.e. power calculations
- ▶ Statistical vs practical significance

## Lecture 16: Sample Size and Power

### Chapter 4.6

1 / 19

## Last Time: Reddie Sleep Example

Tested number of hours of sleep:

- ▶  $H_0 : \mu = 7$
- ▶  $H_A : \mu > 7$

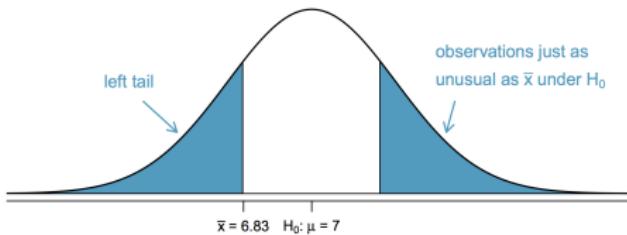
2 / 19

## Two-Sided Alternative Hypothesis

Say instead we had a two-sided alternative hypothesis:

- ▶  $H_0 : \mu = 7$
- ▶  $H_A : \mu \neq 7$

The the p-value would be double:  $2 \times 0.007 = 0.014$ . Picture:



3 / 19

## Setting $\alpha$

Say Dr. Quack is conducting a hypothesis tests. They start with  $\alpha = 0.05$ .

They conduct the test and get  $p\text{-value} = 0.09$ . They then declare "having used an  $\alpha = 0.10$ , we reject the null hypothesis and declare our results to be significant."

What's not honest about this approach?

Ronald Fisher, the creator of p-values, never intended for them to be used this way: <http://en.wikipedia.org/wiki/P-value#Criticisms>

4 / 19

## Goals for Today

- ▶ More in depth discussion of
  - ▶ 10% sampling rule
  - ▶ Skew condition to check to use the normal model
- ▶ How big a sample size do I need?
- ▶ Statistical power
- ▶ Statistical vs practical significance

5 / 19

## 10% Sampling Rule

**Question:** Why do we set  $n$  to be less than 10% of the population size  $N$ ?

**Intuition:** Shouldn't we always sample as many people as we can?

**Answer:** Yes, if we only care about the mean. If we also care about the SE, then we need to be careful.

**Explanation:** Recall from HW5 Q1, sampling without replacement from a rooms that are half male/female but with  $N = 10$  and  $N = 10000$ .

6 / 19

## Finite Population Correction

The finite population correction (FPC) to the SE accounts for the sampling without replacement:

$$SE = \frac{\sigma}{\sqrt{n}} \times \sqrt{\frac{N-n}{N-1}} = \frac{\sigma}{\sqrt{n}} \times FPC$$

Say we have  $N = 10000$ .

- ▶ Let  $n = 100$  (1%), then

$$FPC = \sqrt{\frac{10000 - 100}{10000 - 1}} = 0.995$$

- ▶ Let  $n = 5000$  (50%), then

$$FPC = \sqrt{\frac{10000 - 5000}{10000 - 1}} = 0.707$$

7 / 19

## Finite Population Correction

$$SE = \frac{\sigma}{\sqrt{n}} \times \sqrt{\frac{N-n}{N-1}} = \frac{\sigma}{\sqrt{n}} \times FPC$$

We've been ignoring the FPC. So when

- ▶  $n$  is relatively small, the FPC  $\approx 1$ , so not a problem.
- ▶  $n$  is relatively large, the FPC  $\rightarrow 0$ .  
i.e.  $\frac{\sigma}{\sqrt{n}}$  is not the true SE.

Conclusion: By capping  $n \leq 10\%$  of  $N$ , we have a rule of thumb for keeping the FPC "close" to 1.

8 / 19

## Sampling

We can tie the **conceptual** and **mathematical** notions of sampling:

**Conceptual:** If we sample everybody, we know the true  $\mu$ .

and

**Mathematical:** If  $n = N$  then  $FPC = \sqrt{\frac{N-n}{N-1}} = 0$  then

$$SE = \frac{\sigma}{\sqrt{n}} \times FPC = 0$$

i.e.

- ▶ the sampling distribution is just one point: the true  $\mu$ .
- ▶ if we repeat this procedure many times, we get the same value each time: 0 variability.

9 / 19

## Sampling and the SE

**Question:** Why do we care that our SE is correct?

**Answer:** If not

- ▶ the  $SE$  in confidence intervals is off
- ▶ the z-scores of  $\bar{X}$  have the wrong denominator

10 / 19

## Skew Condition to Check to Use Normal Model

Throughout the book, they talk about the condition for  $\bar{x}$  being nearly normal and using  $s$  in place of  $\sigma$  in  $SE = \frac{\sigma}{\sqrt{n}}$ :

- ▶ On page 164: the population distribution is not strongly skewed
- ▶ On page 167: the data are not strongly skewed
- ▶ On page 168: the distribution of sample observations is not strongly skewed
- ▶ On page 185: the population data are not strongly skewed

11 / 19

## Skew Condition to Check to Use Normal Model

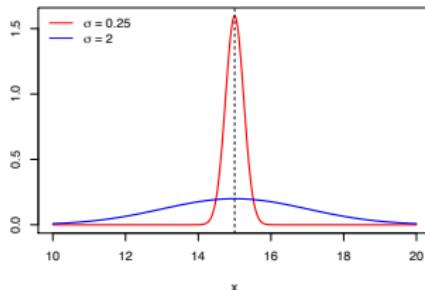
However, they all mean the same thing:

1. The **true population** distribution from which you are drawing your sample observations/data  $x_1, \dots, x_n$  is not too skewed.
2. The **histogram** (visual estimate) of the sample observations/data  $x_1, \dots, x_n$  is not too skewed.

12 / 19

## Sample Size: Thought Experiment

Say you have two distributions with  $\mu = 15$  but different  $\sigma$ .



Which of the two distributions do you think will require a bigger  $n$  to estimate  $\mu$  "well"?

13 / 19

## Margin of Error

Recall our formula for a 95% confidence interval:

$$\left[ \bar{x} - 1.96 \frac{s}{\sqrt{n}}, \bar{x} + 1.96 \frac{s}{\sqrt{n}} \right]$$

The **margin of error** is half the width of the CI.

Say we knew the **true** standard deviation  $\sigma$ , then

$$\text{Margin of Error} = 1.96 \frac{\sigma}{\sqrt{n}}$$

14 / 19

## Identify $n$ for a Desired Margin of Error

To estimate the necessary sample size  $n$  for a maximum desired margin of error  $m$ , we set

$$m \geq z^* \frac{\sigma}{\sqrt{n}}$$

and solve for  $n$ .

15 / 19

## Identify $n$ for a Desired Margin of Error

Since

$$m \geq z^* \frac{\sigma}{\sqrt{n}}$$

$$\sqrt{n} \geq z^* \frac{\sigma}{m}$$

$$n \geq \left( z^* \frac{\sigma}{m} \right)^2$$

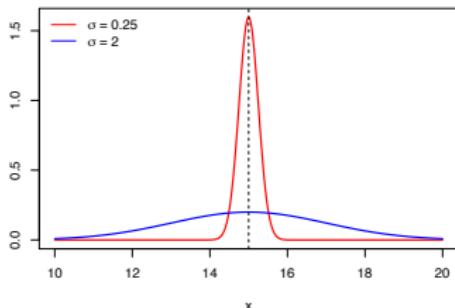
So

- ▶ As  $\sigma$  goes up, you need more  $n$
- ▶ As  $z^*$  goes up, i.e. higher confidence level, you need more  $n$
- ▶ As the desired margin of error goes down, you need more  $n$

16 / 19

## Back to Thought Experiment

For the same desired maximal margin of error  $m$  and same confidence level, we need a larger  $n$  to estimate the mean of the blue curve:



17 / 19

## Type II Error Rate and Power

For a hypothesis test:

- ▶ The significance level  $\alpha$  is the **type I error rate**: the rate at which we reject  $H_0$  when it is true.
- ▶ The **type II error rate  $\beta$**  is the rate at which we fail to reject  $H_0$  when  $H_A$  is true.
- ▶  $1 - \beta$  is called the **statistical power**: the rate at which we reject  $H_0$  when  $H_A$  is true.

18 / 19

## Type II Error Rate and Power

Say we are conducting  $N = A + B + C + D$  hypothesis tests.

		Test conclusion	
		do not reject $H_0$	reject $H_0$ in favor of $H_A$
Truth	$H_0$ true	A	B
	$H_A$ true	C	D

- ▶ The Type I Error rate is  $\alpha = \frac{B}{A+B}$ : rate at which B occurs given  $H_0$  is true.
- ▶ The Type II Error is  $\beta = \frac{C}{C+D}$ : rate at which C occurs given  $H_A$  is true.
- ▶ The power is  $1 - \beta = 1 - \frac{C}{C+D} = \frac{D}{C+D}$ : rate at which D occurs given  $H_A$  is true.

## Lecture 17: Paired Data and Difference of Two Means

Chapter 5.2, 5.1

1 / 18

### Goals for Today

- ▶ Difference of means
- ▶ Note on Practical vs Statistical Significance
- ▶ Paired differences of means

2 / 18

## 6 Types of Questions

Here are the 6 broad types of questions about **population parameters** we'll be answering with statistical methods: confidence intervals and hypothesis tests

1. What is the mean value  $\mu$ ?
2. Are the means  $\mu_1$  and  $\mu_2$  of two groups different?
3. What is the mean paired difference  $\mu_{diff}$ ?
4. What is the proportion  $p$  of "successes"?
5. Are the proportions of "successes"  $p_1$  and  $p_2$  of two groups different?
6. Are the means  $\mu_1, \dots, \mu_k$  of  $k$  groups different?

Today we look at 3 and 2.

3 / 18

## General Outline

We now generalize what we did in Chapter 4:

1. Define the population parameter and determine its point estimate
2. Show that the sampling distribution of the point estimate is Normal
  - ▶ Verify CLT & any additional conditions
  - ▶ Find the SE

Then we either:

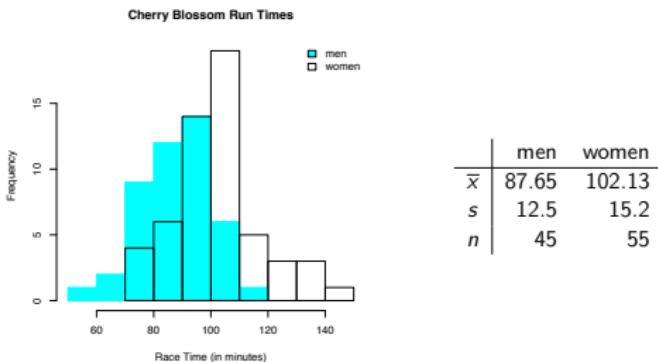
- ▶ Build a confidence interval: point estimate  $\pm z^*SE$
- ▶ Conduct a hypothesis test with test statistic: z-score of the point estimate

$$z = \frac{\text{point estimate} - \text{null value}}{SE}$$

4 / 18

## Chapter 5.2: Are Two Means $\mu_1$ & $\mu_2$ Different?

We randomly sample 45 men (of 7192) and 55 women (of 9732) runners in the 2012 Cherry Blossom Run. Did men run faster than women?



5 / 18

## Difference in Means

We want the difference of two population means:

- ▶  $\mu_w$ : mean time for women
- ▶  $\mu_m$ : mean time for men

Thus:

- ▶ Population parameter:  $\mu_w - \mu_m$ .  
i.e. if men run faster, this is positive
- ▶ Point estimate:  $\bar{x}_w - \bar{x}_m = 102.13 - 87.65 = 14.48$   
i.e. difference of sample means

6 / 18

## Normality of Sampling Distribution

If two sample means  $\bar{x}_1$  and  $\bar{x}_2$

- ▶ each satisfy the 3 CLT conditions
- ▶ Additionally: the two samples are independent from each other

Then the sampling distribution of  $\bar{x}_1 - \bar{x}_2$  will be approximately normal with

- ▶ mean  $\mu_1 - \mu_2$
- ▶ estimated standard error

$$SE_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \approx \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

7 / 18

## Normality of Sampling Distribution

We verify the conditions:

1. Each sample consists of  $\leq 10\%$  of their respective populations.
2. Both histograms don't look too skewed.
3. Each sample has at least 30 observations (rule of thumb).
4. Additionally: the samples are independent (not paired or linked in any way).

Thus the sampling distribution is Normal with mean =  $\mu_w - \mu_m$  and

$$SE_{\bar{x}_w - \bar{x}_m} = \sqrt{\frac{15.2^2}{55} + \frac{12.5^2}{45}} = 2.77$$

8 / 18

## Confidence Interval

A 95% confidence interval for  $\mu_1 - \mu_2$  is

$$\begin{aligned} (\text{point estimate for } \mu_1 - \mu_2) &\pm z^* \times SE \\ (\bar{x}_1 - \bar{x}_2) &\pm 1.96 \times SE_{\bar{x}_1 - \bar{x}_2} \end{aligned}$$

For the Cherry Blossom Run data, a 95% CI for  $\mu_w - \mu_m$  is:

$$14.48 \pm 1.96 \times 2.77 = [9.05, 19.91]$$

9 / 18

## Hypothesis Test

For  $\alpha = 0.001$  (i.e. we want reject with high confidence) we test

- ▶  $H_0 : \mu_w - \mu_m = 0$
- ▶  $H_A : \mu_w - \mu_m > 0$

Test statistic: z-score of  $\bar{x}_w - \bar{x}_m$  under  $H_0$ :

$$\begin{aligned} \frac{\text{point estimate} - \text{null value}}{SE} &= \frac{(\bar{x}_w - \bar{x}_m) - \text{null value}}{SE_{\bar{x}_1 - \bar{x}_2}} \\ &= \frac{14.48 - 0}{2.77} = 5.23 \end{aligned}$$

The p-value is 0, hence we reject  $H_0$  and declare that men ran significantly faster than women.

10 / 18

## Practical vs Statistical Significance

When rejecting  $H_0$ , we call this a **statistically significant** result. But statistically significant results aren't always **practically significant**.

Say for **very** large  $n_M$  &  $n_F$  we observe  $\bar{x}_M = 87.65$  and  $\bar{x}_F = 87.651$  and reject  $H_0$ .

The point estimate of the difference  $\bar{x}_M - \bar{x}_F = 0.001$ . Near negligible!

However, the 95% CI might be:

$$[0.0005, 0.0015]$$

11 / 18

## Practical vs Statistical Significance

Moral of the story

- ▶ Hypothesis tests with “rejections of  $H_0$ ” focus almost entirely on **statistical significance**.
- ▶ Confidence intervals allow you to also focus on **practical significance**.

12 / 18

## Hypothesis Test

13 / 18

## Chapter 5.1: Paired Data

Two sets of observations are **paired** if each observation in one set has a special correspondence or connection with exactly one observation in the other data set.

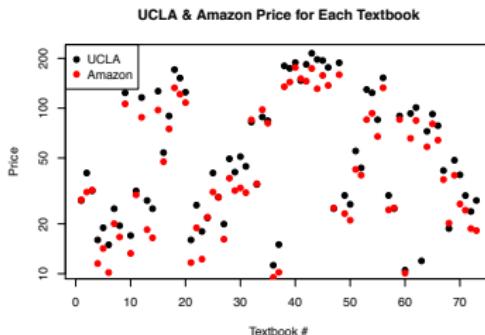
Examples:

- ▶ Cholesterol levels before and after some intervention for the same person
- ▶ Disease rates amongst pairs of twins
- ▶ In the text: price of the same textbook at the UCLA bookstore vs Amazon

14 / 18

## Paired Differences

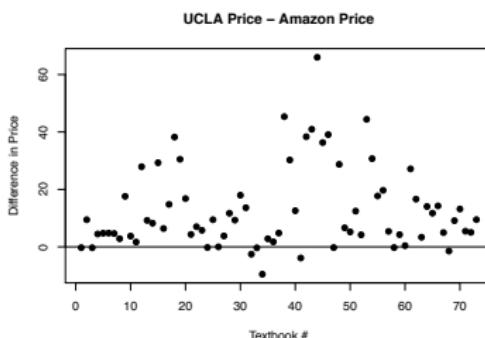
The methodology for paired data remains the same, except our **observations** are the difference in pairs. Example, for the UCLA Bookstore vs Amazon book price example in the text



15 / 18

## Paired Differences

The methodology for paired data remains the same, except our **observations** are the difference in pairs. Example, for the UCLA Bookstore vs Amazon book price example in the text



16 / 18

## Paired Differences

We have

- ▶ population parameter is  $\mu_{\text{diff}}$  with point estimate  $\bar{x}_{\text{diff}}$
- ▶ Check the conditions not on the original observations, but rather the differences.
- ▶ If met,  $\bar{x}_{\text{diff}}$  has a normal sampling distribution
  - ▶ mean  $\mu_{\text{diff}}$
  - ▶  $SE_{\text{diff}} = \frac{\sigma_{\text{diff}}}{\sqrt{n_{\text{diff}}}} \approx \frac{s_{\text{diff}}}{\sqrt{n_{\text{diff}}}}$

17 / 18

## Next Time

- ▶ t-test

18 / 18

## Lecture 18: One-Sample Means With the *t*-Distribution

Chapter 5.3

1 / 17

### Goals for Today

- ▶ What do we do when  $n$  is small?

2 / 17

## Sample Size $n$

We need a **large  $n$**  for two reasons:

1. CLT so sampling distribution of  $\bar{x}$  is normal regardless of the true population distribution.
2. ensure  $s$  is a good point estimate of  $\sigma$ , so  $SE = \frac{\sigma}{\sqrt{n}} \approx \frac{s}{\sqrt{n}}$

What if  $n$  is small? We're stuck **except when**:

1. the sample observations are independent
2. the observations are from a population distribution that is normal

the sampling distribution of  $\bar{x}$  is nearly normal **regardless** of  $n$ .

3 / 17

## Verifying Normality of Population Distribution

Be cautious when verifying the normality condition for small  $n$ . It is important to not only examine the data but also think about where the data come from. For example, ask:

- ▶ Would I expect this distribution to be symmetric?
- ▶ Am I confident that outliers are rare?

4 / 17

## *t* Distribution

Let  $x_1, \dots, x_n$  be a random sample from a **normal** population distribution where  $\sigma$  is unknown. Then the **t-statistic**

$$t = \frac{\bar{x} - \mu}{SE} = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

has probability distribution called a ***t* distribution** with  $n - 1$  degrees of freedom (*df*).

5 / 17

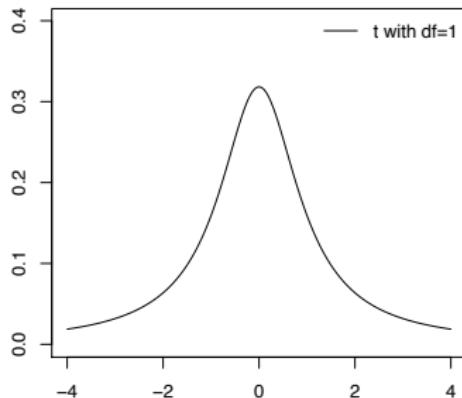
## *t* Distribution

Properties of the *t* distribution:

- ▶ a *t*-distribution has only one parameter: the degrees of freedom *df*.
- ▶ It is bell-shaped and centered at 0
- ▶ Any *t* curve is more spread out than a *z* curve.  
i.e. it has **fatter tails**
- ▶ As the *df* goes to  $\infty$ , the *t* curve approaches the *z* curve.

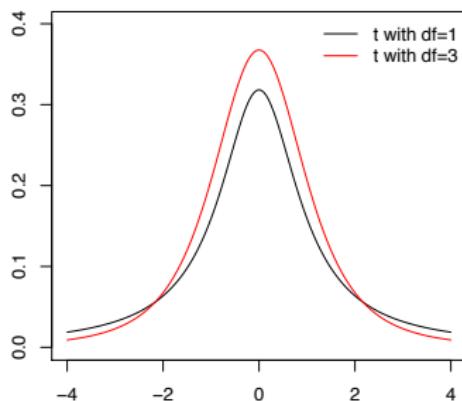
6 / 17

## *t* Distribution Examples



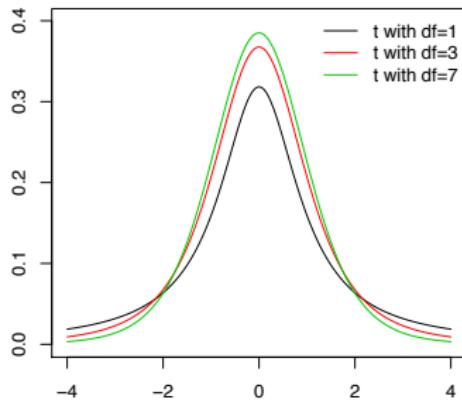
7 / 17

## *t* Distribution Examples



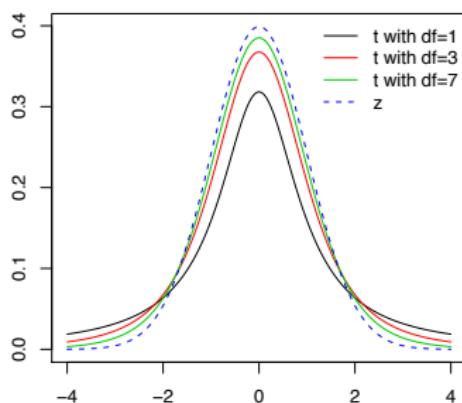
7 / 17

## *t* Distribution Examples



7 / 17

## *t* Distribution Examples



7 / 17

## Conditions for Using t Distribution

We use the  $t$  distribution when you have

- ▶  $n$  is small. E.g. 3, 5, 10, 15.
- ▶ **Independence:**  $n \leq 10\%$  rule
- ▶ **Observations come from a nearly normal distribution:**
  - ▶ Look at a histogram of the data (difficult when  $n$  is small)
  - ▶ Consider whether any previous experiences alert us that the data may be normal

8 / 17

## $t$ -Tables

If  $n = 11$ , we use  $df = 11 - 1 = 10$  and do a look up on the  $t$ -table on page 410:

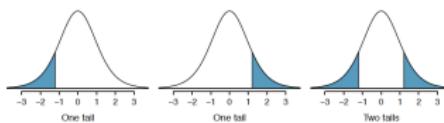


Figure B.1: Three  $t$  distributions.

	one tail	0.100	0.050	0.025	0.010	0.005
	two tails	0.200	0.100	0.050	0.020	0.010
df		3.08	6.31	12.71	31.82	63.66
1		1.44	1.32	2.45	3.14	3.71
2		1.41	1.89	2.95	3.65	3.90
3		1.40	1.86	2.31	2.99	3.39
4		1.38	1.83	2.26	2.83	3.25
5		1.37	1.81	2.23	2.76	3.17
6		1.36	1.80	2.20	2.72	3.11
7		1.36	1.79	2.18	2.68	3.05
8		1.36	1.78	2.17	2.65	3.00
9		1.36	1.78	2.16	2.63	2.97
10		1.37	1.81	2.23	2.76	3.17
11		1.36	1.80	2.20	2.72	3.11

9 / 17

## Confidence Intervals

Confidence intervals: Use  $t_{df}^*$  instead of  $z^*$

$$\bar{x} \pm t_{df}^* SE = \left[ \bar{x} - t_{df}^* \times \frac{s}{\sqrt{n}}, \bar{x} + t_{df}^* \times \frac{s}{\sqrt{n}} \right]$$

So for example, to get a 95% C.I. based on

$$n = 11 \Rightarrow df = 10 \Rightarrow t_{10}^* = 2.23$$

10 / 17

## *t*-Test Example

Example 5.19 on page 252: A random sample of 25 New Yorkers were asked how much sleep they get per night. Does the data below provide strong evidence that New Yorkers sleep less than 8 hours a night on average? Set  $\alpha = 0.05$

$n$	$\bar{x}$	$s$	$\min$	$\max$
25	7.73	0.77	6.17	9.78

11 / 17

## *t*-Test

Conditions:

- ▶ **Independence:** 25 is obviously less than 10% of the population of NYC
- ▶ **Normality:** Not an exact science. The halfway point of the min and max is 7.975, which is fairly close to  $\bar{x} = 7.73$ . So symmetric enough?

The test statistic is the *t*-statistic:

$$t = \frac{\bar{x} - \text{null value}}{SE} = \frac{\bar{x} - \text{null value}}{\frac{s}{\sqrt{n}}} = \frac{7.73 - 8}{\frac{0.77}{\sqrt{25}}} = -1.75$$

Since  $n = 25$ ,  $df = 25 - 1 = 24$ .

12 / 17

## *t*-Test

p-Value: we use the *t* distribution i.e. the *t*-table on page 410:

one-tail	0.100	0.050	0.025	0.010	0.005
two-tail	0.200	0.100	0.050	0.020	0.010
df = 24	1.32	1.71	2.06	2.49	2.80

Since

- ▶ 1.75 is in [1.71, 2.06]
- ▶ by symmetry -1.75 is in [-2.06, -1.71]
- ▶ the one-sided p-value is in between [0.025, 0.05]

Decision: Since the p-value  $< \alpha = 0.05$ , we reject  $H_0$  that NY'ers sleep 8 hours a night at the  $\alpha = 0.05$  significance level in favor of the hypothesis they sleep more.

13 / 17

## History of *t* Distribution

The *t* distribution was derived by William Sealy Gosset in 1908, a chemist/statistician at the Guinness Brewery in Dublin, Ireland.



14 / 17

## History of *t* Distribution

Gosset was concerned with **small-sample statistics** about barley given that brewers are limited in the number of batches of beer they can brew.

Guinness prohibited its employees from publishing. So Gosset had to use the pseudonym "Student" to conceal his identity.

The *t*-test's complete name is the **(Student's) *t*-test**.

15 / 17

## History of *t* Distribution

In fact if you go to the Guinness Brewery at St James's Gate in Dublin, Ireland...



16 / 17

## History of *t* Distribution



17 / 17

## Lecture 19: ANOVA Part I

### Chapter 5.5

1 / 25

## Previously: Conditions for Using t Distribution

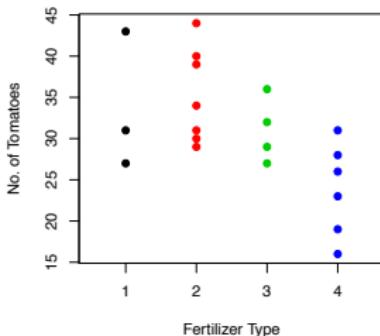
We use the  $t$  distribution when you have

- ▶  $n$  is small.
- ▶ **Independence:**  $n \leq 10\%$  rule
- ▶ **Observations come from a nearly normal distribution:**
  - ▶ Look at a histogram of the data (difficult when  $n$  is small)
  - ▶ Consider whether any previous experiences alert us that the data may be normal

2 / 25

## Analysis of Variance (ANOVA)

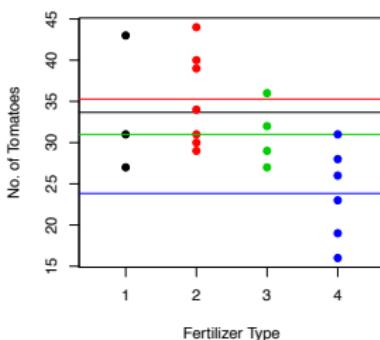
A farmer has the choice of four tomato fertilizers and wants to compare their performance in terms of crop yield.



3 / 25

## Analysis of Variance (ANOVA)

A farmer has the choice of four tomato fertilizers and wants to compare their performance in terms of crop yield.



4 / 25

## Analysis of Variance (ANOVA)

We have  $k = 4$  groups AKA **levels of a factor**: the 4 types of fertilizer.

- ▶  $n_i$  plants assigned to each of the  $k = 4$  fertilizers:

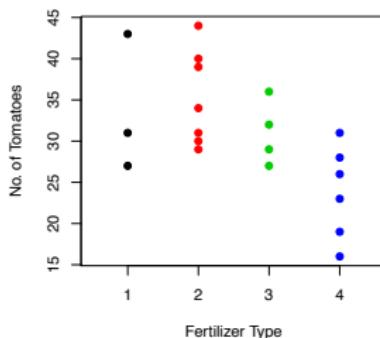
$n_1$	$n_2$	$n_3$	$n_4$	total $n$
3	7	4	6	20

- ▶ Count the number of tomatoes on each plant

5 / 25

## Tomato Fertilizer

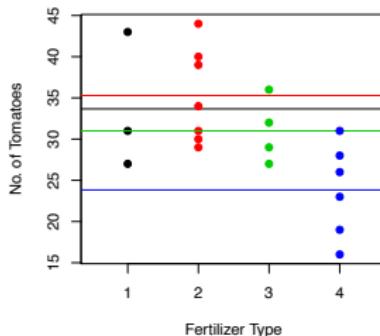
We observe the following, where each point is one tomato plant.



6 / 25

## Tomato Fertilizer

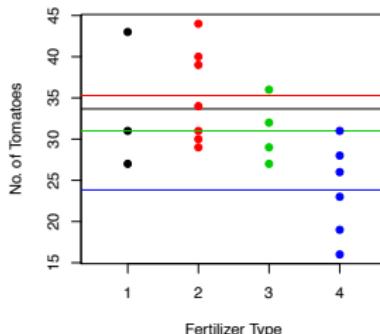
We observe the following, where each point is one tomato plant.  
Plot the sample mean of each level.



6 / 25

## Tomato Fertilizer

We observe the following, where each point is one tomato plant.  
Plot the sample mean of each level. [Question: are the mean tomato yields different?](#)



6 / 25

## Analysis of Variance

Say we have  $k$  groups and want to compare the  $k$  means:

$$\mu_1, \mu_2, \dots, \mu_k$$

We could do  $\binom{k}{2}$  individual two-sample tests.

Ex. for groups 1 & 2:

$$\begin{aligned} H_0 : \quad & \mu_1 = \mu_2 \\ \text{vs. } H_a : \quad & \mu_1 \neq \mu_2 \end{aligned}$$

7 / 25

## Analysis of Variance

Or we do a single overall test via Analysis of Variance ANOVA:

The hypothesis test is:

$$\begin{aligned} H_0 : \quad & \mu_1 = \mu_2 = \dots = \mu_k \\ \text{vs. } H_a : \quad & \text{at least one of the } \mu_i \text{'s are different} \end{aligned}$$

8 / 25

## Analysis of Variance

ANOVA asks: where is the overall variability of the observations originate from?

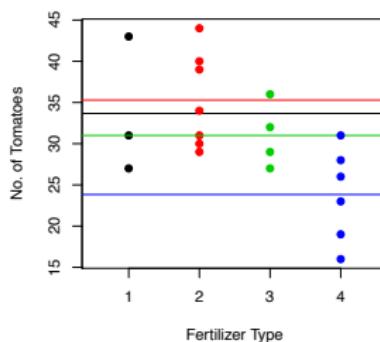
The test statistic used to compute a  $p$ -value is now the F-statistic:

$$F = \frac{\text{measure of between-group variability}}{\text{measure of within-group variability}}$$

9 / 25

## Tomato Fertilizer Example

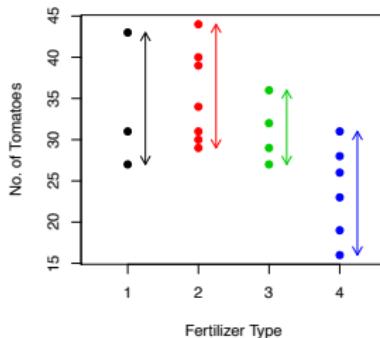
Numerator: the between-group variation refers to the variability between the levels (the 4 horizontal lines):



10 / 25

## Tomato Fertilizer Example

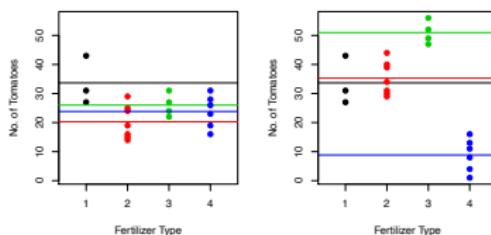
Denominator: the **within-group variation** refers to the variability **within** each level (the 4 vertical arrows):



11 / 25

## Tomato Fertilizer Example

Now compare the following two plots. Which has “more different” means?



12 / 25

## Tomato Fertilizer Example

- ▶ They have the same within-group variability. Call this value  $W$
- ▶ The right plot has higher between group variability b/c the 4 means are more different. Call these values  $B_{left}$  and  $B_{right}$  with  $B_{left} < B_{right}$
- ▶ Recall  $F = \frac{\text{measure of between-group variability}}{\text{measure of within-group variability}}$
- ▶ Since  $\frac{B_{left}}{W} < \frac{B_{right}}{W}$ , thus  $F_{left} < F_{right}$  The right plot as a larger  $F$ -statistic

13 / 25

## $F$ Distributions

Assuming  $H_0$  is true (that  $\mu_1 = \mu_2 = \dots = \mu_k$ ), the  $F$ -statistic

$$F = \frac{\text{measure of between-group variability}}{\text{measure of within-group variability}}$$

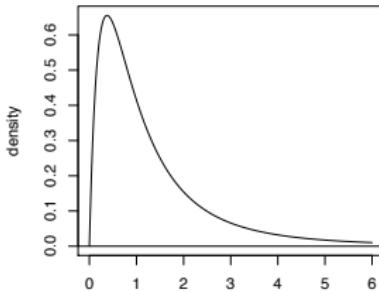
follows the  $F$  distribution with  $df_1 = k - 1$  and  $df_2 = n - k$  degrees of freedom where

- ▶  $n$  = total number of observations
- ▶  $k$  = number of groups

14 / 25

## *F* Distributions

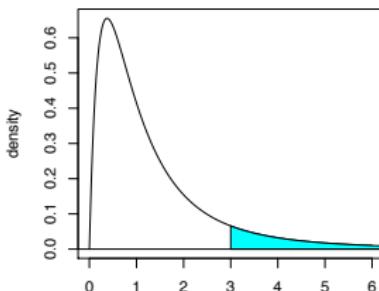
For  $df_1 = 4$  and  $df_2 = 6$ , the *F* distribution looks like:



15 / 25

## *F* Distributions

*p*-values are computed where “more extreme” means **larger**. Say the  $F = 3$ , the *p*-value is the [area to the right of 3](#) and is computed in R: `pf(3,df1=4,df2=6,lower.tail=FALSE)`



16 / 25

## Conducting An F-Test

The results are typically summarized in an [ANOVA table](#):

Source of Variation	df	SS	MS	F	p-value
Between groups	$k - 1$	$SSTr$	$MSTr = \frac{SSTr}{k-1}$	$\frac{MSTr}{MSE}$	$p$
Within groups	$n - k$	$SSE$	$MSE = \frac{SSE}{n-k}$		
Total	$n - 1$	$SST$			

17 / 25

## Conditions

1. The observations have to be [independent](#). 10% rule.
2. Trade off of  $n$  and [normality](#) of observations within each group.
3. Each of the groups has [constant variance](#)  $\sigma_1^2 = \dots = \sigma_k^2 = \sigma^2$ .  
Check via:
  - ▶ boxplots
  - ▶ comparing the sample standard deviations  $s_1, \dots, s_k$

18 / 25

## Discussion of Quiz

Question 1: Why did  $\frac{1}{20}$  studies yield a positive/significant result i.e. that there is a link between jelly beans and acne?

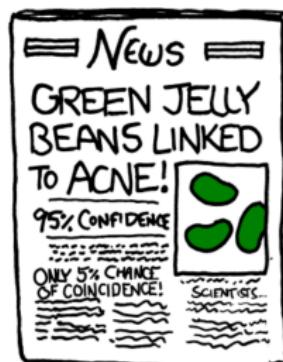
Not that the p-value is 0.05, rather that  $\alpha = 0.05$ :

- ▶ significance level AKA
- ▶ type I error rate AKA
- ▶ false positive rate

i.e. we expect 1 out of 20 results to be significant even if there is no effect.

19 / 25

## Publication Bias



20 / 25

## Publication Bias

**Publication bias:** people only highlight significant/positive results.  
From Wikipedia: "Publication bias occurs when the publication of research results depends on their nature and direction."

To counter this, some prominent medical journals including

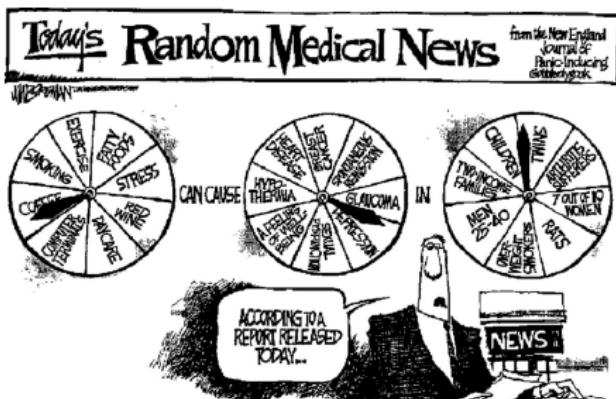
- ▶ New England Journal of Medicine
- ▶ The Lancet
- ▶ Journal of the American Medical Association

require registration of a trial **before** it starts so that unfavorable results are not withheld from publication.

Journal of Negative Results: <http://www.jnrbm.com/>

21 / 25

## Publication Bias



From: Sterne JA, Davey Smith G (2001) Sifting the evidence - What's wrong with significance tests. BMJ 322: 226231.

22 / 25

## What $\alpha$ to Use?

Should I use  $\alpha = 0.05$  as my significance level? Before using it, put some thought into the balance between:

- ▶ **Type I errors.** Setting a smaller  $\alpha$  yields a more **conservative** procedure: all things being equal, you will reject  $H_0$  less often.
- ▶ **Type II errors.** Setting a bigger  $\alpha$  yields a more **liberal** procedure: all things being equal, you will reject  $H_0$  more often.

23 / 25

## Multiple Testing

A related issue is the statistical concept of **multiple testing**.

Say we are conducting many experiments, and  $H_0$  is true for all of them.

If you repeat experiments many times, you're bound to get a significant result eventually just by **chance alone**.

24 / 25

## Multiple Testing

What do people do? Make the  $\alpha$  stricter! i.e.

- ▶ make the  $\alpha$  smaller
- ▶ i.e. less chance the p-value is smaller than  $\alpha$
- ▶ i.e. less chance of incorrectly rejecting  $H_0$  when it is true

Use the [Bonferroni correction](#) to  $\alpha$ : If you are conducting  $n$  tests, use  $\alpha^* = \frac{\alpha}{n}$ . You'll study its properties in HW8.

## Lecture 20: Single Proportion Test

### Chapter 6.1

1 / 17

## Quiz 8

**Question 1:** According to the article, why do scientists even bother with correlational/observational studies, when no notions of causality can be established?

**Answer:** One reason is that correlational studies are excellent starting points for deciding which hypotheses to evaluate with the more rigorous randomized controlled experiment.

2 / 17

## Quiz 8

**Question 2:** The article argues that various scientific disciplines should set professional labeling standards for material discussed in the media... Rank the four possible labels in order of how much credence the public should give them, from lowest to highest.

**Answer:**

3. preliminary result
1. large-scale observational study
4. large-sample randomized controlled test
2. well-established scientific law that we know how to apply in a wide range of conditions

3 / 17

## Causality

What is causality? How do we establish it?

- ▶ [http://nfs.unipv.it/nfs/minf/dispense/patgen/lectures/files/disease\\_causality.html](http://nfs.unipv.it/nfs/minf/dispense/patgen/lectures/files/disease_causality.html)
- ▶ <http://bayes.cs.ucla.edu/BOOK-2K/>

4 / 17

## Question for Today

According to a (representatively sampled) poll done by the New York Times/CBS News in June 2012, only about 44% of the American public approved of the Supreme Court's performance.

The sample proportion  $\hat{p} = 0.44$  is point estimate of  $p$ : the true (population) proportion of the American public who approves.

What are some next things to ask?

- ▶ What was  $n$ ?
- ▶ What is the SE of  $\hat{p} = 44\% = 0.44$ ?
- ▶ What is the sampling distribution of  $\hat{p}$ ?

5 / 17

## Question for Today

Just like with  $\bar{x}$ , if we want to use the normal model to

- ▶ build confidence intervals via  $z^*$
- ▶ conduct hypothesis tests via the normal tables

we need the sampling distribution of  $\hat{p}$  to be nearly normal.

This happens when the population distribution of 0's and 1's is not too strongly skewed. As the sample size  $n \rightarrow \infty$ , this is less of an issue by the CLT.

Note:

$$\hat{p} = \frac{x_1 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

where each of the  $x_i$ 's are 0/1 success/failure Bernoulli random variables.

6 / 17

## Conditions for Sampling Dist'n of $\hat{p}$ Being Nearly Normal

The sampling distribution of the sample proportion  $\hat{p}$  based on sample size  $n$  is nearly normal when

- ▶ The observations are independent: the 10% rule
- ▶ We expect to see at least 10 successes and 10 failures in our sample. This is called the **success-failure condition**:
  - ▶  $np \geq 10$
  - ▶  $n(1 - p) \geq 10$

7 / 17

## Conditions for Sampling Dist'n of $\hat{p}$ Being Nearly Normal

If conditions are met, then the sampling distribution of  $\hat{p}$  is nearly normal with

- ▶ mean  $p$  (the true population proportion)
- ▶ standard error

$$SE_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$$

Note the similarity of the previous formula for the sample mean  $\bar{x}$ :

$$SE_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \sqrt{\frac{\sigma^2}{n}}$$

8 / 17

## What $p$ to use?

But we **don't know** what  $p$  is. So what  $p$  do we use

- ▶ to check the success/failure condition?
- ▶ for the  $SE_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$ ?

For

- ▶ Confidence intervals: plug in the **point estimate**  $\hat{p}$  of  $p$
- ▶ Hypothesis tests: plug in the **null value**  $p_0$  from  $H_0 : p = p_0$

9 / 17

## Confidence Intervals

Going back to the poll:  $\hat{p} = 0.44$  based on  $n = 976$ . What is a 95% confidence interval?

Check the conditions and find SE **using  $p = \hat{p}$**

- ▶  $976 < 10\%$  of 313 million  $\Rightarrow$  independence
- ▶ Defining a success as a person approving of the job done by the Supreme Court:
  - ▶  $976 \times \hat{p} = 976 \times .44 = 429$  successes  $\geq 10$
  - ▶  $976 \times (1 - \hat{p}) = 976 \times .56 = 547$  failures  $\geq 10$

$$SE_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}} \approx \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{0.44(1-0.44)}{976}} = 0.016$$

10 / 17

## Confidence Intervals

A 95% confidence interval using the normal model has  $z^* = 1.96$ , thus:

$$\text{point estimate} \pm 1.96 \times SE$$

In our case

$$\hat{p} \pm 1.96 \times SE_{\hat{p}} = 0.44 \pm 1.96 \times 0.016 = (0.409, 0.471)$$

11 / 17

## Hypothesis Tests

Thomas Carcetti is running for mayor of Baltimore. His campaign manager claims he has more than 50% support of the electorate.

The Baltimore Sun collects a random sample of  $n = 500$  likely voters and finds that 52% support him. Does this provide convincing evidence for the claim of Carcetti's manager at the 5% significance level?

12 / 17

## Hypothesis Tests

The hypothesis test is, with the null value  $p_0 = 0.5$

$$H_0 : p = p_0 \\ \text{vs} \\ H_A : p > p_0$$

Check the conditions and find SE using  $p = p_0$

- ▶  $500 < 10\%$  of the population of Baltimore  $\Rightarrow$  independence
- ▶ Success-failure condition

- ▶  $np_0 = 500 \times 0.5 = 250 \geq 10$
- ▶  $n(1 - p_0) = 500 \times (1 - 0.5) = 250 \geq 10$

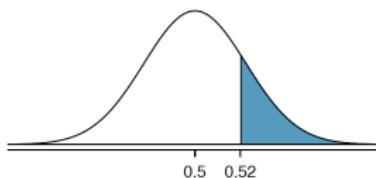
$$SE_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}} \approx \sqrt{\frac{p_0(1-p_0)}{n}} = \sqrt{\frac{0.5(1-0.5)}{500}} = 0.022$$

13 / 17

## Hypothesis Tests

$$z = \frac{\text{point estimate } \hat{p} - \text{null value } p_0}{SE_{\hat{p}}} = \frac{0.52 - 0.50}{0.022} = 0.89$$

p-value is 0.1867. In the original %'age scale:



Hence we do **not** reject the null hypothesis, and we do not find convincing evidence to support the campaign manager's claim.

14 / 17

## Next Time

Same as with the jump from

$$\mu \text{ to } \mu_1 - \mu_2$$

i.e. from one to two-sample tests for means, we make the jump from

$$p \text{ to } p_1 - p_2$$

i.e. from one to two-sample tests for proportions.

## Lecture 21: Difference of two proportions

Chapter 6.2

1 / 23

### Question for today

How do we infer about a difference in proportions  $p_1 - p_2$ ?

2 / 23

## Surveys

The way a question is phrased in survey can influence a person's response. Ex on p.269: the Pew Research Center conducted a survey with the following question:

*By 2014 all Americans will be required to have health insurance. X while Y. Do you approve or disapprove of this policy?*

where X and Y were randomly ordered between

- ▶ People who do not buy insurance will pay a penalty
- ▶ People who cannot afford it will receive financial help from the government

Let's infer about the difference in proportion of people who **approve**. Any guesses which is higher?

3 / 23

## Example from Text

	Sample size $n_i$	Approve (%)	Disapprove (%)	Other (%)
people who do not buy it will pay a penalty given first	771	47	49	3
people who cannot afford it will receive financial help from the gov't given first	732	34	63	3

4 / 23

## Example from Text

	Sample size $n_i$	Approve (%)	Don't Approve (%)
people who do not buy it will pay a penalty given first	771	47	53
people who cannot afford it will receive financial help from the gov't given first	732	34	66

So  $\hat{p}_1 - \hat{p}_2 = 0.47 - 0.34 > 0$ : people are more likely to support Obamacare in the first scenario.

5 / 23

## Conditions...

When

- ▶ Both sample proportions  $\hat{p}_1$  and  $\hat{p}_2$  are approximately **normal**:
  - ▶ independence
  - ▶ success/failure condition: at least 10 successes and failures
- ▶ the two samples are independent from each other

6 / 23

## ... for Sampling Dist'n of $\hat{p}_1 - \hat{p}_2$ Being Normal

The sampling distribution of  $\hat{p}_1 - \hat{p}_2$  is approximately Normal with

- ▶ mean  $p_1 - p_2$
- ▶ standard error

$$SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{SE_{\hat{p}_1}^2 + SE_{\hat{p}_2}^2} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

7 / 23

## Standard Error

Recall we showed that the SE for  $\bar{x}_1 - \bar{x}_2$  was

$$SE_{\bar{x}_1 - \bar{x}_2} = \sqrt{SE_{\bar{x}_1}^2 + SE_{\bar{x}_2}^2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Compare this to

$$SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{SE_{\hat{p}_1}^2 + SE_{\hat{p}_2}^2} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

8 / 23

## What $p_1$ & $p_2$ ?

What  $p_1$  &  $p_2$  do we

- ▶ Use to check success/failure condition?
- ▶ Use in  $SE_{\hat{p}_1 - \hat{p}_2}$ ?

For

- ▶ Confidence intervals: plug in  $\hat{p}_1$  and  $\hat{p}_2$
- ▶ Hypothesis tests: plug in **pooled estimate**  $\hat{p}$

9 / 23

## Confidence Intervals

What is a 90% confidence interval for the difference in proportions?

Check the conditions:

- ▶ Normality for each group
  - ▶ Independence: both groups  $\leq 10\%$  of respective populations
  - ▶ The success/failure condition for **both** groups:
    - ▶ Group 1: 362 successes and  $771 - 362 = 409$  failures
    - ▶ Group 2: 249 successes and 483 failures
- ▶ We assume both groups were sampled independently.

10 / 23

## Confidence Intervals

- ▶ Point estimate is  $\hat{p}_1 - \hat{p}_2 = 0.47 - 0.34 = 0.13$
- ▶ Plug in  $\hat{p}_1$  and  $\hat{p}_2$  into SE:

$$SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}} = \dots = 0.025$$

- ▶ A 90% confidence interval for  $p_1 - p_2$  is:

$$\text{point estimate} \pm z^* \times SE = 0.13 \pm 1.65 \times 0.025 = (0.09, 0.17)$$

11 / 23

## Interpretation

Two key observations:

- ▶ (9%, 17%) does not contain 0, suggestive of a true difference.
- ▶ The sign of the difference:  $\hat{p}_1 - \hat{p}_2 = 0.13 > 0$

More support Obamacare if stated as follows:

*People who do not buy it will pay a penalty while people who cannot afford it will receive financial help from the government.*

12 / 23

## Hypothesis Tests

Now we are interested in testing the difference of two proportions:

$$\begin{aligned} H_0 : p_1 - p_2 &= 0 \\ \text{vs } H_1 : p_1 - p_2 &\neq 0 \end{aligned}$$

Note this can be re-expressed as:

$$\begin{aligned} H_0 : p_1 &= p_2 \\ \text{vs } H_1 : p_1 &\neq p_2 \end{aligned}$$

i.e. under  $H_0$  the two proportions are both equal to some value  $p$ :

$$p_1 = p_2 = p$$

13 / 23

## Hypothesis Tests

So to

- ▶ Verify the success-failure condition
- ▶ Compute the standard SE

we use a [pooled estimate](#)  $\hat{p}$  of the proportion  $p$ . i.e. as if there were [no difference](#) between them, so we can combine them:

$$\hat{p} = \frac{\text{total \# of successes}}{\text{total \# of cases}}$$

The SE to use is:

$$SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n_1} + \frac{\hat{p}(1 - \hat{p})}{n_2}} = \sqrt{\hat{p}(1 - \hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

14 / 23

## Exercise 6.31 on Page 305

A 2010 survey asked 827 randomly sample voters in California "How do you feel about drilling for oil and natural gas off the coast of California?"

	College Grad	
	Yes	No
Support	154	132
Oppose	180	126
Don't Know	104	131
Total	438	389

Test at the  $\alpha = 0.10$  significance level if the proportion of college graduates who support off-shore drilling is different than that of non-college graduates.

15 / 23

## Exercise 6.31 on Page 305

The pooled estimate is  $\hat{p} = \frac{154+132}{438+389} = 0.346$ . Check the conditions:

1. Normality of both point estimates
  - ▶ Independence
    - ▶  $n_1 = 438 \leq 10\%$  of pop. of CA college grads
    - ▶  $n_2 = 389 \leq 10\%$  of pop. of CA non college grads
  - ▶ Success/failure: both groups have at least 10 successes and 10 failures.
2. We assume that both groups are sampled independently.

16 / 23

## Exercise 6.31 on Page 305

- ▶ Point estimate  $\hat{p}_1 - \hat{p}_2 = 0.352 - 0.339 = 0.013$
- ▶  $SE_{\hat{p}_1 - \hat{p}_2} \approx \sqrt{\hat{p}(1 - \hat{p}) \left( \frac{1}{n_1} + \frac{1}{n_2} \right)} = 0.033$
- ▶ Test statistic: z-score of  $\hat{p}_1 - \hat{p}_2$  under  $H_0 : p_1 - p_2 = 0$ 
$$z = \frac{\text{point estimate} - \text{null value}}{SE} = \frac{0.013 - 0}{0.033} = 0.392$$
- ▶  $p$ -value: 0.6922. i.e. we fail to reject  $H_0$ . We don't have strong evidence of a difference in support.

17 / 23

## Jury Selection

Preview of next lecture: In many trials a big issue is the racial makeup of the jury.

Question: is there a way to figure out if there is a racial bias in jury selection?

18 / 23

## Jury Selection

Say we have a juror pool (registered voters) where the racial breakdown is:

Race	White	Black	Hispanic	Other	Total
Registered Voters	72%	7%	12%	9%	100%

19 / 23

## Jury Selection

If we pick  $n = 100$  jurors at random (i.e. unbiasedly), we expect the breakdown of counts to be:

Race	White	Black	Hispanic	Other	Total
Registered Voters	72%	7%	12%	9%	100%
Representation	72	7	12	9	$n = 100$

20 / 23

## Jury Selection

Say we *observe* the following counts:

Race	White	Black	Hispanic	Other	Total
Registered Voters	72%	7%	12%	9%	100%
Representation	0	0	100	0	$n = 100$

Fairly obvious bias in juror selection!

21 / 23

## Jury Selection

But what about the following? Is there a bias? i.e. a non-random mechanism at play?

Race	White	Black	Hispanic	Other	Total
Registered Voters	72%	7%	12%	9%	100%
Representation	75	6	11	8	$n = 100$

22 / 23

## Next Two Lectures

Chi-square tests are used to compare **expected** counts with **observed** counts.

Two tests we'll see:

- ▶ Goodness-of-fit tests: for frequency tables
- ▶ Tests for independence: for contingency/two-way tables

## Lecture 22: Chi-Square Tests for Goodness-of-Fit

### Chapter 6.3

1 / 18

### Question for Today

Say we had  $n = 100$  people picked as jurors, we expect the breakdown to be:

Race	White	Black	Hispanic	Other	Total
Registered Voters	72%	7%	12%	9%	100%
Representation	72	7	12	9	$n = 100$

2 / 18

## Question for Today

Say we observe the following. Is there a bias? i.e. a non-random mechanism?

Race	White	Black	Hispanic	Other	Total
Registered Voters	72%	7%	12%	9%	100%
Representation	75	6	11	8	$n = 100$

3 / 18

## Chi-Square Tests

Chi-square  $\chi^2$  tests allow us to compare

- ▶ Observed counts
- ▶ Expected counts

i.e. What is the “goodness” of the fit of the observed counts to the expected counts?

4 / 18

## The Data

Let's use  $n = 275$  people. Assuming the same proportions as above, we compute the **expected** counts. Ex:  $198 = 275 \times 0.72$ .

Race	White	Black	Hispanic	Other	Total
Expected Counts	198	19.25	33	24.75	275

5 / 18

## The Data

Let's use  $n = 275$  people. Assuming the same proportions as above, we compute the **expected** counts. Ex:  $198 = 275 \times 0.72$ . Now say we observe the following counts:

Race	White	Black	Hispanic	Other	Total
Expected Counts	198	19.25	33	24.75	275
Observed Counts	205	26	25	19	275

6 / 18

## Hypothesis Test in General

$H_0$  : The data are consistent with the specified distribution.  
vs     $H_A$  : The data are not consistent with the specified distribution.

$H_0$  can also be stated: the data are a random sample from the distribution and any differences of observed vs expected reflect natural sampling variation.

7 / 18

## Hypothesis Test in Our Case

$H_0$  : the jurors are randomly sampled i.e. there is no racial bias  
vs     $H_A$  : the jurors are not randomly sampled i.e. there is racial bias

8 / 18

## Null Distributions

To compute p-values we compare the [computed test statistic](#) to a [null distribution](#): the distribution of the test statistic under  $H_0$ .

### 1. means/proportions:

- ▶ test statistic: z-score of  $\bar{x}/\hat{p}$
- ▶ null distribution: normal distribution

### 2. t-test:

- ▶ test statistic: t-statistic
- ▶ null distribution: t-distribution with  $df = n - 1$

### 3. ANOVA:

- ▶ test statistic: F-statistic
- ▶ null distribution: F-distribution with  $df_1 = k - 1$  and  $df_2 = n - k$

### 4. Goodness-of-fit:

- ▶ test statistic:  $\chi^2$ -statistic
- ▶ null distribution:  $\chi^2$  distribution with  $df = k - 1$

## Deviations

Previously, many test statistics had the following form:

$$z = \frac{\text{point estimate} - \text{null value}}{\text{SE of point estimate}}$$

For goodness-of-fit, it's similar. For each of the  $k$  groups compute

$$Z = \frac{\text{observed} - \text{expected}}{\sqrt{\text{expected}}}$$

Note when:

- ▶ observed = expected  $\Rightarrow Z = 0$
- ▶ observed > expected  $\Rightarrow Z > 0$
- ▶ observed < expected  $\Rightarrow Z < 0$

The Z's measure deviations.

## Deviations

Now treat +'ve and -'ve differences as the same by squaring  $Z$ :

$$\begin{aligned} Z^2 &= \left( \frac{\text{observed} - \text{expected}}{\sqrt{\text{expected}}} \right)^2 \\ &= \frac{(\text{observed} - \text{expected})^2}{\text{expected}} \end{aligned}$$

Why square it and not absolute value it? It's easier to do [calculus](#) on  $x^2$  than  $|x|$ .

11 / 18

## Chi-Square Test Statistic

Finally sum all values of  $Z^2$ . This is the [chi-square test statistic for one-way tables](#).

$$\chi^2 = \sum_{i=1}^k Z_i^2 = \sum_{i=1}^k \frac{(\text{observed}_i - \text{expected}_i)^2}{\text{expected}_i}$$

12 / 18

## Chi-Square Test Statistic

In the case of the jury data, we have 4 groups: white, black, hispanic, and other:

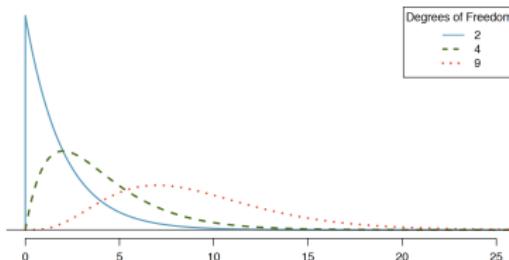
$$\begin{aligned}\chi^2 &= Z_w^2 + Z_b^2 + Z_h^2 + Z_o^2 \\ &= \frac{(205 - 198)^2}{198} + \dots + \dots + \frac{(19 - 24.75)^2}{24.75} \\ &= 5.89\end{aligned}$$

13 / 18

## p-values

We compare the test statistic to a  $\chi^2$  distribution with  $df = k - 1$  degrees of freedom.

Note: not  $df = n - 1$  like with t-test.



14 / 18

## p-values

The *p*-value is the **area to the right** of the test statistic. Use p.412:

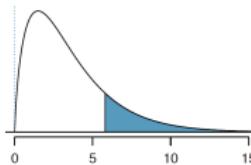


Figure B.2: Areas in the chi-square table always refer to the right tail.

Upper tail	0.3	0.2	0.1	0.05	0.02	0.01	0.005	0.001	
df	2	3.21	3.22	4.61	5.99	7.82	9.21	10.60	13.82
	3	3.66	4.64	6.25	7.81	9.84	11.34	12.84	16.27
	4	4.88	5.99	7.78	9.49	11.67	13.28	14.86	18.47
	5	6.06	7.29	9.24	11.07	13.39	15.09	16.75	20.52

In our case,  $df = k - 1 = 3$ , and  $\chi^2 = 5.89$ , which is in between (4.64, 6.25), so *p*-value is in between (0.1, 0.2). Not overwhelming evidence against  $H_0$ .

15 / 18

## Hypothetical Scenarios

Say we have two hypothetical scenarios of observed counts:

Race	White	Black	Hispanic	Other	Total
Expected Counts	198	19.25	33	24.75	275
Observed Counts					275

- ▶ For all 4 groups, say observed = expected, then

$$\chi^2 = 0 + 0 + 0 + 0 = 0$$

hence *p*-value = 1.

- ▶ Say we observed 275 others and 0 for the rest, then

$$\chi^2 = 2786.11 + 84.46 + 189.57 + 12588.11 = 15648.25$$

hence *p*-value = 0.

16 / 18

## Assumptions for Chi-Square Test

1. **Independence:** Each case is independent of the other
2. **Sample size:** Similarly like with proportions, we need at least 5 cases in each scenario (each cell in the table)
3. **Degrees of freedom:** We need at least  $df = 2$ , i.e.  $k \geq 3$

17 / 18

## Next Time

We look at [chi-square tests for two-way tables](#) to test for independence. i.e. are two variables independent from each other?

18 / 18

## Lecture 23: Tests for Independence in Two-Way Tables

Chapter 6.4

1 / 26

### Quiz 9

**Question:** While the results of the controlled experiment suggesting that women are at a disadvantage in science hiring may come as no surprise, what argument is made that this discrimination is not entirely due to overt misogyny? Answer in one sentence.

**Answer:** Women rated women candidates lower as well, suggesting not so much explicit misogyny, but rather manifestation of subtler prejudices internalized from societal stereotypes.

2 / 26

## Quiz 9

**Question:** Knowing nothing else about the problem (sample sizes, SE, etc), what can we conclude about the difference in means between men and women?

**Answer:** No, refer to HW8 Question 7. We had two overlapping CIs, but the CI on the difference did not include 0.

3 / 26

## Conditions for Chi-Square Test for Goodness-of-Fit

1. **Independence:** Each case is independent of the each other
2. **Sample size/distribution:** We need at least 5 cases in each scenario i.e. each cell in the table
3. **Degrees of freedom:** We need at least  $df = 2$ , i.e.  $k \geq 3$

4 / 26

## Today's Example

Google is always tinkering with its search ranking algorithm. Say we want to compare the following 3 algorithms:

1. the current version
2. test algorithm 1
3. test algorithm 2

5 / 26

## Today's Example

They measure user satisfaction with the results for a particular search with the `new_search` variable:

- ▶ no new search: User clicked on a result. Suggests user is satisfied with result.
- ▶ new search: User `did not` click on a result and tried a new related search. Suggests user is dissatisfied with result.

6 / 26

## Today's Example

So we have two categorical variables:

- ▶ algorithm: current, test 1, or test 2
- ▶ new search: yes or no

Are they independent? i.e. independent of which algorithm is used, do we have the same levels of new search?

7 / 26

## Today's Example

Say we observed the following results:

		algorithm			Total
		Current	Test 1	Test 2	
new search	No new search	4000	2000	2000	8000
	New search	1000	500	500	2000
Total	5000	2500	2500	10000	

For all 3 algorithms, there is a new search  $\frac{1}{5}$  of the time.

algorithm and new search are **independent**: regardless of which algorithm used, the proportion of new searches stays the same.

8 / 26

## Today's Example

Now say instead we observed the following results:

new search	algorithm			Total
	Current	Test 1	Test 2	
No new search	4000	2500	1500	8000
New search	1000	0	1000	2000
Total	5000	2500	2500	10000

In this case, `algorithm` and `new search` are **not independent**: depending on which `algorithm` used, the proportion of new searches **is different**.

9 / 26

## Hypothesis Test

We test at the  $\alpha = 0.05$  significance level:

$$\begin{aligned} H_0 : & \text{ the algorithms each perform equally well} \\ \text{vs } H_A : & \text{ the algorithms do not perform equally well} \end{aligned}$$

i.e. are the categorial variables `algorithm` and `new search` independent?

10 / 26

## Different Names

The following all refer to the same test:  $\chi^2$  test for

- ▶ two-way tables
- ▶ i.e. contingency tables
- ▶ independence of two categorical variables
- ▶ homogeneity: are the algorithms homogeneous in their performance?

11 / 26

## Example from Textbook

Let's make the values match the example from the textbook on page 284:

		algorithm			Total
		Current	Test 1	Test 2	
new search	No new search	3511	1749	1818	7078
	New search	1489	751	682	2922
Total	5000	2500	2500	10000	

12 / 26

## Example from Textbook

Before we start, let's make each column reflect a proportion and not a count.

new search	algorithm			Total
	Current	Test 1	Test 2	
No new search	0.7022	0.6996	0.7272	0.7078
New search	0.2978	0.3004	0.2728	0.2922
Total	1	1	1	1

If all algorithms performed the same, we'd expect

- ▶ 0.7078 for all 3 values in the top row
- ▶ 0.2922 for all 3 values in the bottom row

Are we observing what we expect? i.e. What is the degree of this deviation?

13 / 26

## What's Expected

We expect:

new search	algorithm			Total
	Current	Test 1	Test 2	
No new search				7078 = 0.7078 × 10000
New search				2922 = 0.2922 × 10000
Total	5000	2500	2500	10000

14 / 26

## What's Expected

We expect:

new search	Current	algorithm		Total
		Test 1	Test 2	
No new search		$1769.5 = 0.7078 \times 2500$		7078
New search		$730.5 = 0.2922 \times 2500$		2922
Total	5000	2500	2500	10000

15 / 26

## What's Expected

We expect:

new search	Current	algorithm		Total
		Test 1	Test 2	
No new search		$1769.5 = 0.7078 \times 2500$	1769.5	7078
New search		$730.5 = 0.2922 \times 2500$	730.5	2922
Total	5000		2500	10000

16 / 26

## What's Expected

We expect:

new search	algorithm				Total
	Current	Test 1	Test 2		
No new search	3539 = 0.7078 × 5000	1769.5	1769.5	7078	
New search	1461 = 0.2922 × 5000	730.5	730.5	2922	
Total	5000	2500	2500	10000	

17 / 26

## Observed vs. Expected

Expected Counts:

new search	algorithm				Total
	Current	Test 1	Test 2		
No new search	3539	1769.5	1769.5	7078	
New search	1461	730.5	730.5	2922	
Total	5000	2500	2500	10000	

Observed Counts:

new search	algorithm				Total
	Current	Test 1	Test 2		
No new search	3511	1749	1818	7078	
New search	1489	751	682	2922	
Total	5000	2500	2500	10000	

18 / 26

## Chi-Square Statistic

We compute  $\chi^2$  test statistic: for all  $i = 1, \dots, 6$  cells

$$\frac{(\text{observed count}_i - \text{expected count}_i)^2}{\text{expected count}_i}$$

$$\text{Row 1, Col 1} = \frac{(3511 - 3539)^2}{3539} = 0.222$$

⋮

⋮

$$\text{Row 2, Col 3} = \frac{(682 - 730.5)^2}{730.5} = 3.220$$

So

$$\begin{aligned}\chi^2 &= 0.222 + 0.237 + \dots + 3.220 \\ &= 6.120\end{aligned}$$

19 / 26

## Chi-Square Distribution

We compare this to a  $\chi^2$  distribution to get the p-value. What are the degrees of freedom?

$$\begin{aligned}df &= (\# \text{ of rows} - 1) \times (\# \text{ of columns} - 1) \\ &= (R - 1) \times (C - 1) \\ &= (2 - 1) \times (3 - 1) = 2 \text{ in our case}\end{aligned}$$

20 / 26

## Chi-Square Distribution

Looking up 6.120 in the  $\chi^2$  table on page 412 on the  $df = 2$  row, it would be between 0.05 and 0.01. Since our  $\alpha = 0.05$ , we reject the null hypothesis and accept the alternative that the algorithms do not perform equally well.

i.e. the algorithm and new search categorical variables are independent.

21 / 26

## Conditions/Assumptions

Nearly identical to conditions/assumptions for  $\chi^2$  tests for goodness-of-fit:

1. **Independence:** Each case is independent of the other
2. **Sample size/distribution:** We need at least 5 cases in each scenario i.e. each cell in the table
3. **Degrees of freedom:** (Different than before) We need  $df = (R - 1) \times (C - 1) \geq 2$ .

22 / 26

## Why Are They Called Degrees of Freedom?

In the case of  $\chi^2$  tests, the degrees of freedom is the number of values needed before you specify **all** values in the cells of the table.

23 / 26

## Why Are They Called Degrees of Freedom? Rows

Each row has  $df = 2$  because if we specify 2 values, all values in the row are specified.

Example:

new search	algorithm			Total
	Current	Test 1	Test 2	
No new search	X	Y		7078
New search				2922
Total	5000	2500	2500	10000

then the missing value is  $7078 - X - Y$ .

i.e. the **wiggle room** we have is  $C - 1$  two cells

24 / 26

## Why Are They Called Degrees of Freedom? Columns

Each column has  $df = 1$  because if we specify 1 value, all values in the column are specified.

Example:

		algorithm			Total
new search		Current	Test 1	Test 2	
No new search	X				7078
	New search				2922
Total		5000	2500	2500	10000

then the missing value is  $5000 - X$ .

i.e. the [wiggle room](#) we have is  $R - 1$  one cell

25 / 26

## Why Are They Called Degrees of Freedom? Columns

So the overall  $df$  is  $(C - 1) \times (R - 1)$ , in our case  $df = 2$ .

		algorithm			Total
new search		Current	Test 1	Test 2	
No new search	X	Y			7078
	New search				2922
Total		5000	2500	2500	10000

i.e. if we know these two values, we can fill the rest of the table.

26 / 26

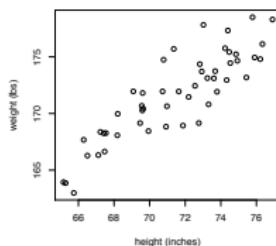
## Lecture 24: Linear Regression Part I

Chapter 7.1-7.2

1 / 23

### Questions for Today

Say we have the height/weight of 50 individuals and we display the scatterplot/bivariate plot of the seemingly linear relationship:



Questions:

- ▶ What is the “best” fitting line through these points?
- ▶ What do we mean by “best”?

2 / 23

## Regression

There are many types of regression, all in order to estimate the relationship between variables. We start by considering simple

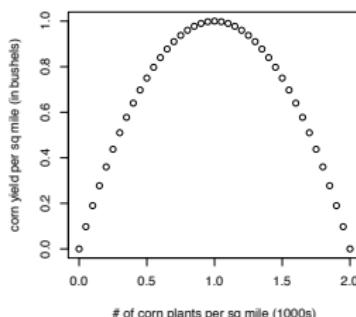
linear regression (SLR):

- ▶ a single explanatory variable / independent variable / predictor variable  $x$
- ▶ an outcome variable / dependent variable  $y$
- ▶ a presumed linear relationship between them

3 / 23

## Example of Non-Linear Relationship

At first as you plant more corn plants, you have higher yield, but past a certain point plants fight for limited resources and they die.



4 / 23

## Modeling $x$ and $y$ Linearly

The **SLR model** assumes that the relationship between  $x$  and  $y$  can be modeled by a line:

$$y = \beta_0 + \beta_1 x$$

where

- ▶  $\beta_0$  is the unknown **intercept parameter**
- ▶  $\beta_1$  is the unknown **slope parameter**

5 / 23

## Procedure

Based on  $n$  pairs of observations  $(x_i, y_i)$

1. Compute **point estimates**
  - ▶  $b_0$  of parameter  $\beta_0$
  - ▶  $b_1$  of parameter  $\beta_1$
2. Associate standard errors  $SE_{b_0}$  and  $SE_{b_1}$
3. For both the intercept and slope
  - ▶ Build confidence intervals
  - ▶ Do hypothesis test

$$\begin{aligned} H_0 : \beta &= 0 \\ \text{vs} \quad H_A : \beta &\neq 0 \end{aligned}$$

The equation

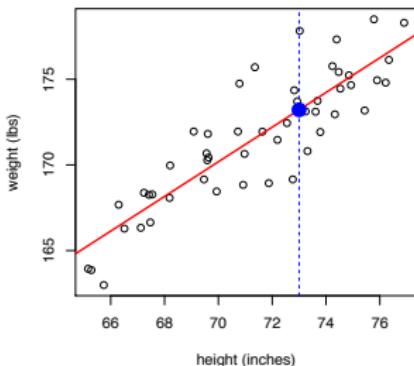
$$\hat{y} = b_0 + b_1 x$$

is called the **least squares line** where  $\hat{y}$  is the **fitted/predicted value**.

6 / 23

## Fitted Value

Here  $\hat{y} = 100 + 0.99x$ . Thus for  $x = 73$ ,  $\hat{y} = 173.22$ :



7 / 23

## Residuals

**Residuals** are what's leftover: leftover variation in the data unexplained by the model:

$$\begin{aligned}\text{Residual} &= \text{Data} - \text{Fit} \\ e_i &= y_i - \hat{y}_i\end{aligned}$$

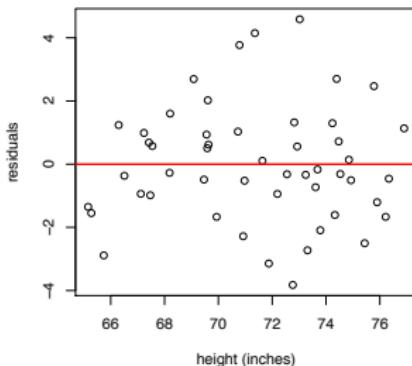
where  $e_i$  is the **residual** of the  $i^{th}$  observation  $(x_i, y_i)$ .

We can think of the  $e_i$ 's as **deviations** from the model. The smaller the deviations, the better the fit.

8 / 23

## Residual Plot

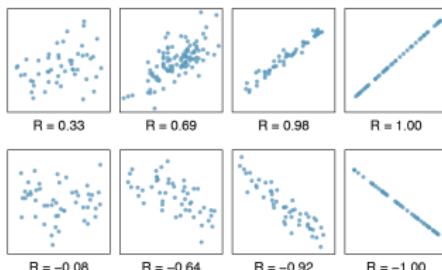
Residual plots: take previous plot and flatten the red line by subtracting  $\hat{y}$  from  $y$ .



9 / 23

## Correlation Coefficient

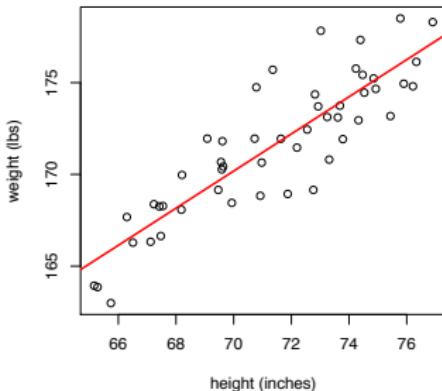
The correlation coefficient  $R$  is a value between  $[-1, 1]$  that measures the strength of the linear relationship between  $x$  and  $y$ .



10 / 23

## Best Fitting Line

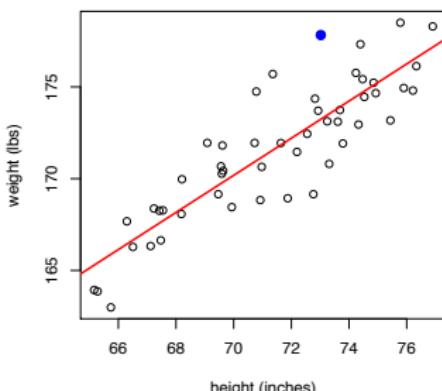
What does “best fitting line” mean?



11 / 23

## Best Fitting Line

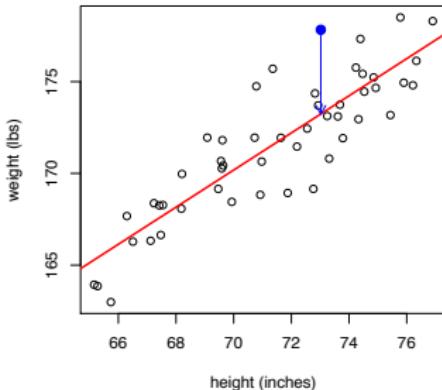
Consider ANY point  $x_i$  for  $i = 1, \dots, 50$  (in blue).



12 / 23

## Best Fitting Line

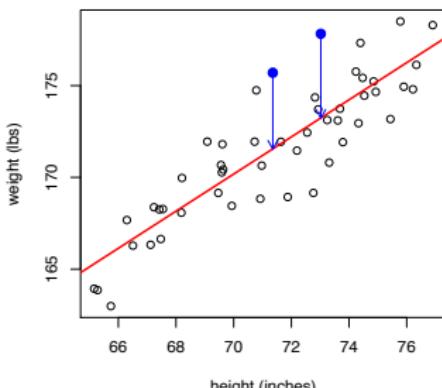
Now consider this point's deviation from the regression line



13 / 23

## Best Fitting Line

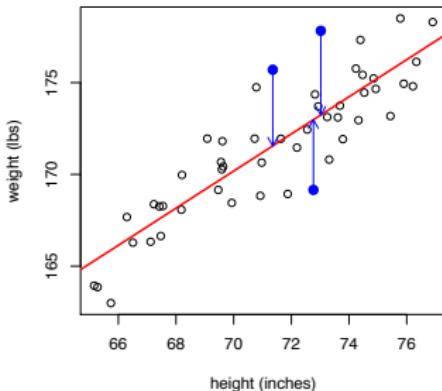
Do this for another point  $x_i \dots$



14 / 23

## Best Fitting Line

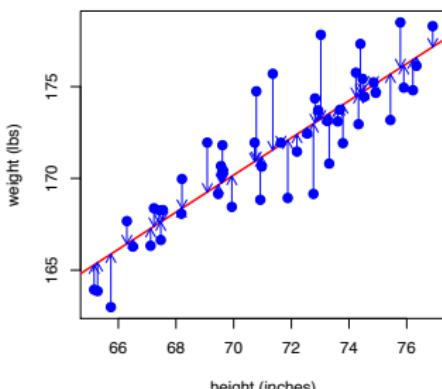
Do this for another point  $x_i \dots$



15 / 23

## Best Fitting Line

The regression line minimizes the sum of the squared arrow lengths.



16 / 23

## Least Squares

i.e. the regression line minimizes:

$$e_1^2 + e_2^2 + \dots + e_n^2$$

This is called **minimizing the least squares criterion**.

17 / 23

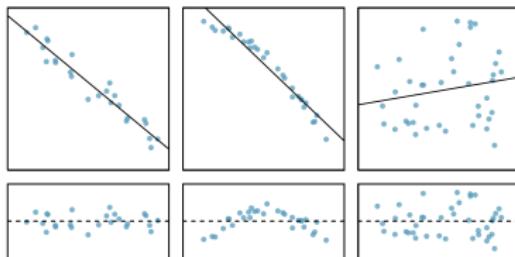
## Conditions for Simple Linear Regression

- ▶ **Linearity:** The data should show a linear trend.
- ▶ **Independence:** The residuals should be independent
- ▶ **Nearly normal residuals:** The residuals  $e_i$  must be nearly normal (verify with QQ-plot) with mean 0.
- ▶ **Constant variability:** The variability of points around the least squares line remains roughly constant (i.e. for all values of  $x$ ).

18 / 23

## Behavior of Residuals: 3 Examples

Sample data + regression on top, residual plots on bottom.



- ▶ Plots 1 and 3 are roughly linear.
- ▶ Plots 1 and 3 have roughly constant variability, but the 3rd plot has higher variability

19 / 23

## Finding the Least Squares Line

To find the least squares line we need to find the point estimates:

- ▶ The point estimate  $b_1$  of the slope  $\beta_1$  is

$$b_1 = \frac{s_y}{s_x} R$$

- ▶ The regression line **always** goes through  $(\bar{x}, \bar{y})$ . We use this fact to find the point estimate of  $b_0$  of the intercept  $\beta_0$ .

20 / 23

## Finding the Point Estimate of the Intercept $b_0$

Given the slope and a point on the line  $(x_0, y_0)$ , the equation for the line can be written as

$$\begin{aligned}\text{slope} &= \frac{\text{rise}}{\text{run}} = \frac{y - y_0}{x - x_0} \\ y - y_0 &= \text{slope} \times (x - x_0)\end{aligned}$$

So

$$\begin{aligned}y - \bar{y} &= b_1(x - \bar{x}) \\ \text{so} \quad y &= (\bar{y} - b_1\bar{x}) + b_1x \\ \text{so} \quad b_0 &= \bar{y} - b_1\bar{x}\end{aligned}$$

21 / 23

## Measuring the Strength of a Fit

If  $R = -1$  or  $R = 1$  we have a perfect linear fit between  $x$  and  $y$ , if  $R = 0$  then there is no fit.

However  $R^2$  is a more commonly used measure of the strength of fit. For SLR, it is correlation coefficient squared, but not for other kinds of regression.

$R^2$  of a linear model describes the proportion of the total variation in  $y$  that is explained by the least squares line.

22 / 23

## Next Time

- ▶ How to interpret regression line parameter estimates
- ▶ Categorical Variable for  $x$ : male vs female, new vs used, etc.
- ▶ Inference for linear regression