

Missing Data

Nicholas Mattei, Tulane University

CMPSS3660 – Introduction to Data Science – Fall 2019

<https://rebrand.ly/TUDatascience>



Many Thanks

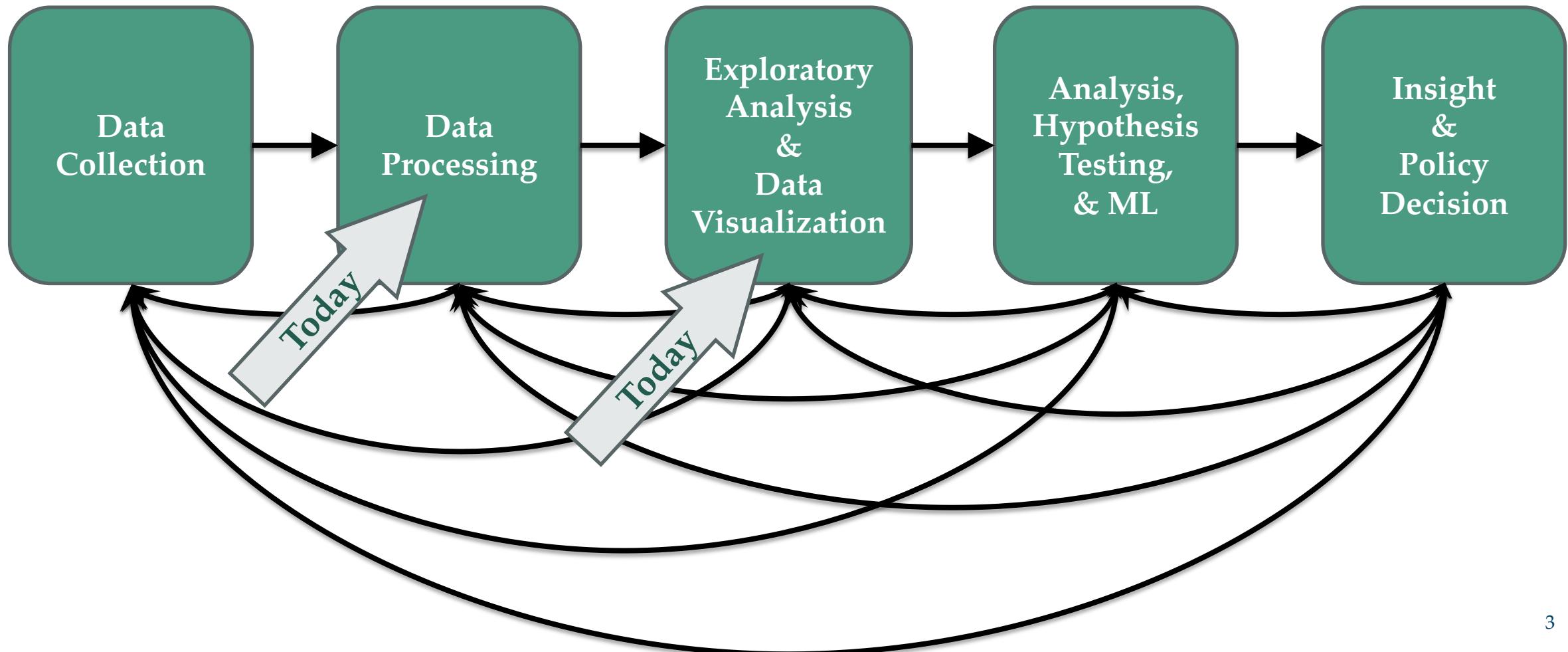
Slides based off Introduction to Data Science from John P. Dickerson -
<https://cmsc320.github.io/>

Announcements

- Peer Observation Tuesday 10/29
- Finish Lab 7
- Review Questions
- Today's Lecture



The Data LifeCycle



Today's Lecture

- Missing Data ...
 - What is it?
 - Simple methods for **imputation**.
 - ... with a tiny taste of Stats/ML lecturers to come.



Thanks to John Atwood and Wenjiang Fu

Missing Data

- Missing data is information that we want to know, but don't!
- It can come in many forms, e.g.:
 - People not answering questions on surveys
 - Inaccurate recordings of the height of plants that need to be discarded
 - Canceled runs in a driving experiment due to rain
- Could also consider missing columns (no collection at all) to be missing data ...



Key Question

- Why is the data missing?
 - What mechanism is it that contributes to, or is associated with, the probability of a data point being absent?
 - Can it be explained by our observed data or not?
- The answers drastically affect what we can ultimately do to compensate for the missing-ness



Complete Case Analysis

- Delete all tuples with any missing values at all, so you are left only with observations with all variables observed

```
# Clean out rows with nil values
df = df.dropna()
```

- Default behavior for libraries for analysis (e.g., regression).
 - We'll talk about this much more during the Stats/ML lectures
- This is the simplest way to handle missing data. In some cases, will work fine; in others, ??????????????:
 - Loss of sample will lead to variance larger than reflected by the size of your data.
 - May bias your sample.

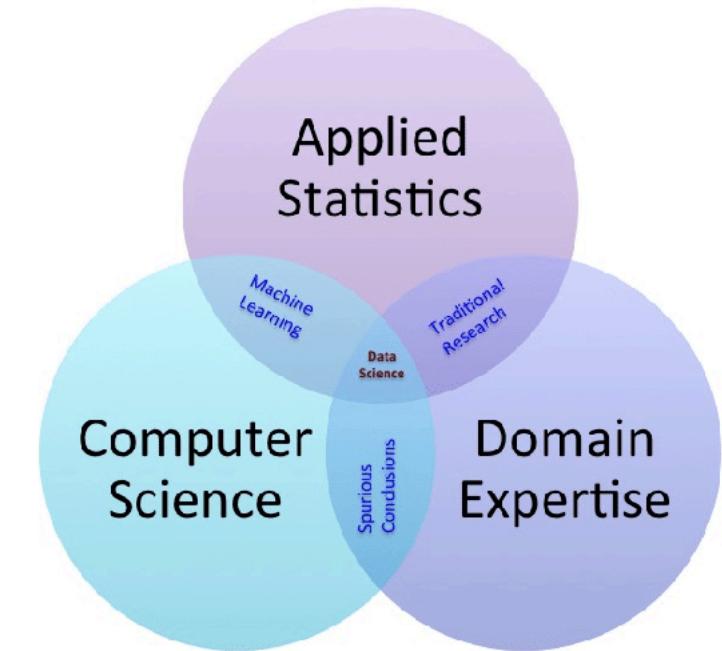


Example

- **Dataset:** Body fat percentage in men, and the circumference of various body parts [Penrose et al., 1985]
 - **Question:** Does the circumference of certain body parts predict body fat percentage?
 - Given **complete** data, how would you answer this ?????????
 - One way to answer is **regression analysis**:
 - One or more independent variables ("predictors")
 - One dependent variables ("outcome")
 - What is the relationship between the predictors and the outcome?
 - What is the conditional expectation of the dependent variable given fixed values for the dependent variables?
-
- Generalized body composition prediction equation for men using simple measurement techniques", K.W. Penrose, A.G. Nelson, A.G. Fisher, FACSM, Human Performance research Center, Brigham Young University, Provo, Utah 84602 as listed in Medicine and Science in Sports and Exercise, vol. 17, no. 2, April 1985, p. 189.
 - <http://staff.pubhealth.ku.dk/~tag/Teaching/share/data/Bodyfat.html>

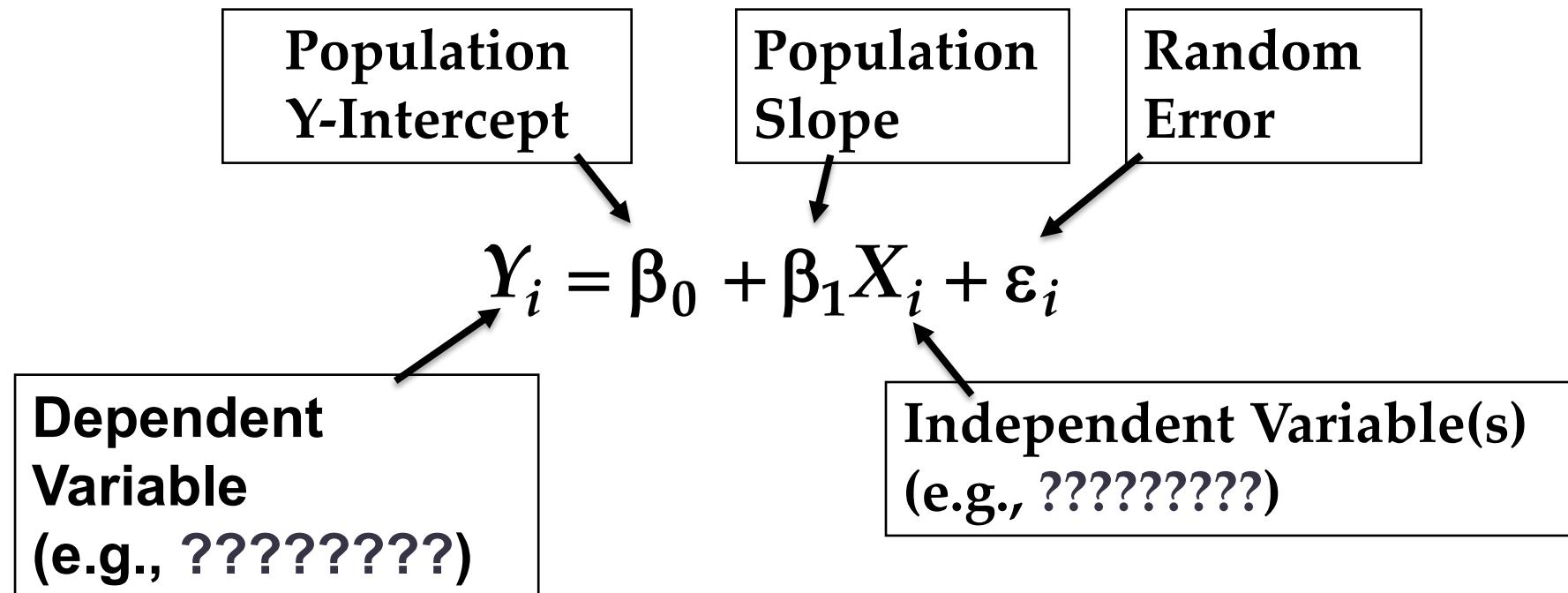
A Side Note On Terms

- For linear regression we have equations of where we want to know the relationship between an **outcome** given some **predictors**.
- If you have a ML background:
 - Get **target, outcome given predictors, observations.**
- If you have a Stats background:
 - Get **endogenous** variables given **exogenous** variables.
- If you are more of a Math person:
 - Get **dependent variable** (y-axis) given one or more **independent variables** (x-axis).



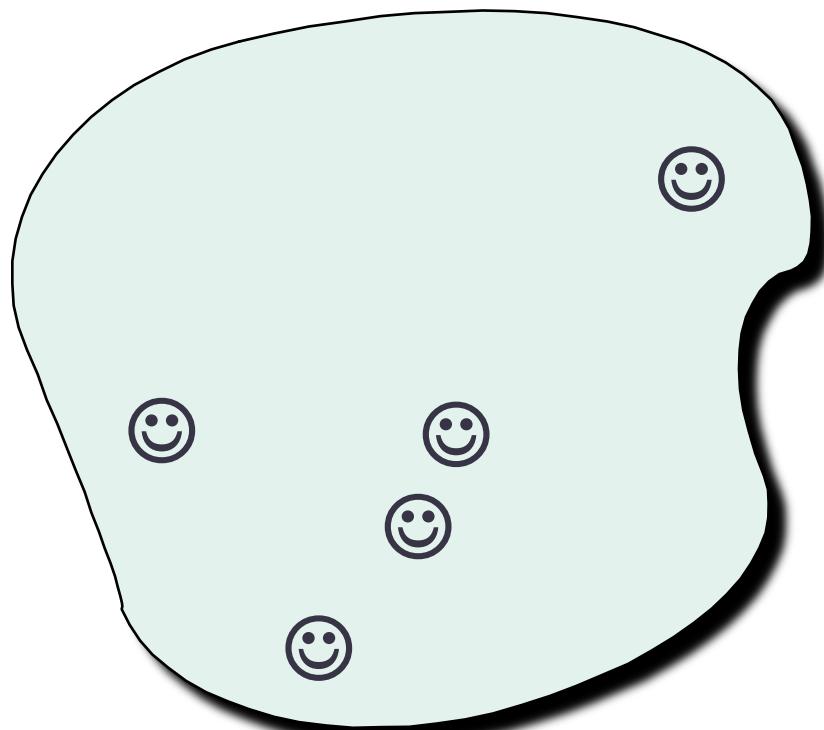
Linear Regression

- Assumption: relationship between variables is **linear**:
 - (We'll relax linearity, study in more depth later.)



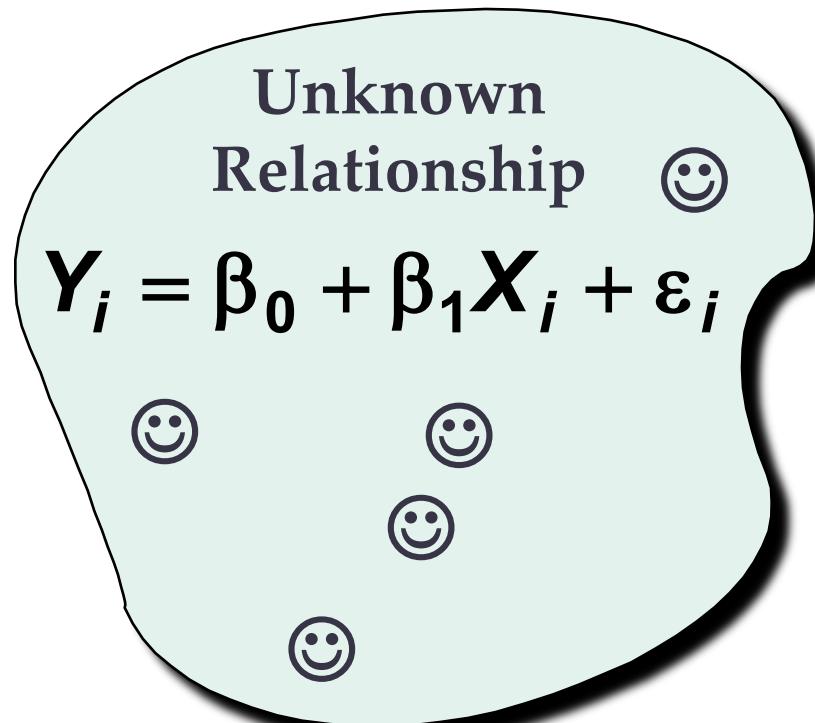
Population & Sample Regression Models

Population

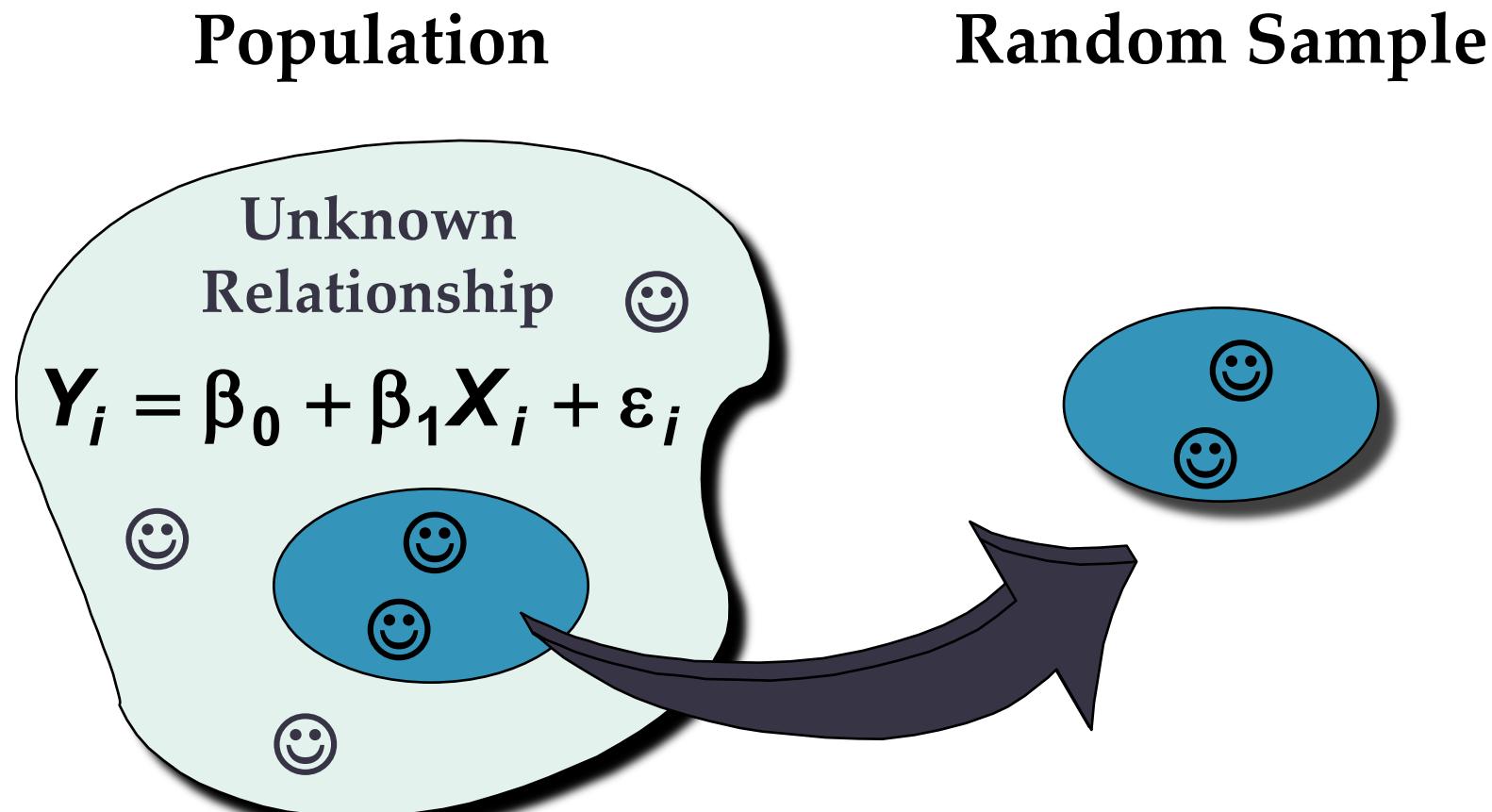


Population & Sample Regression Models

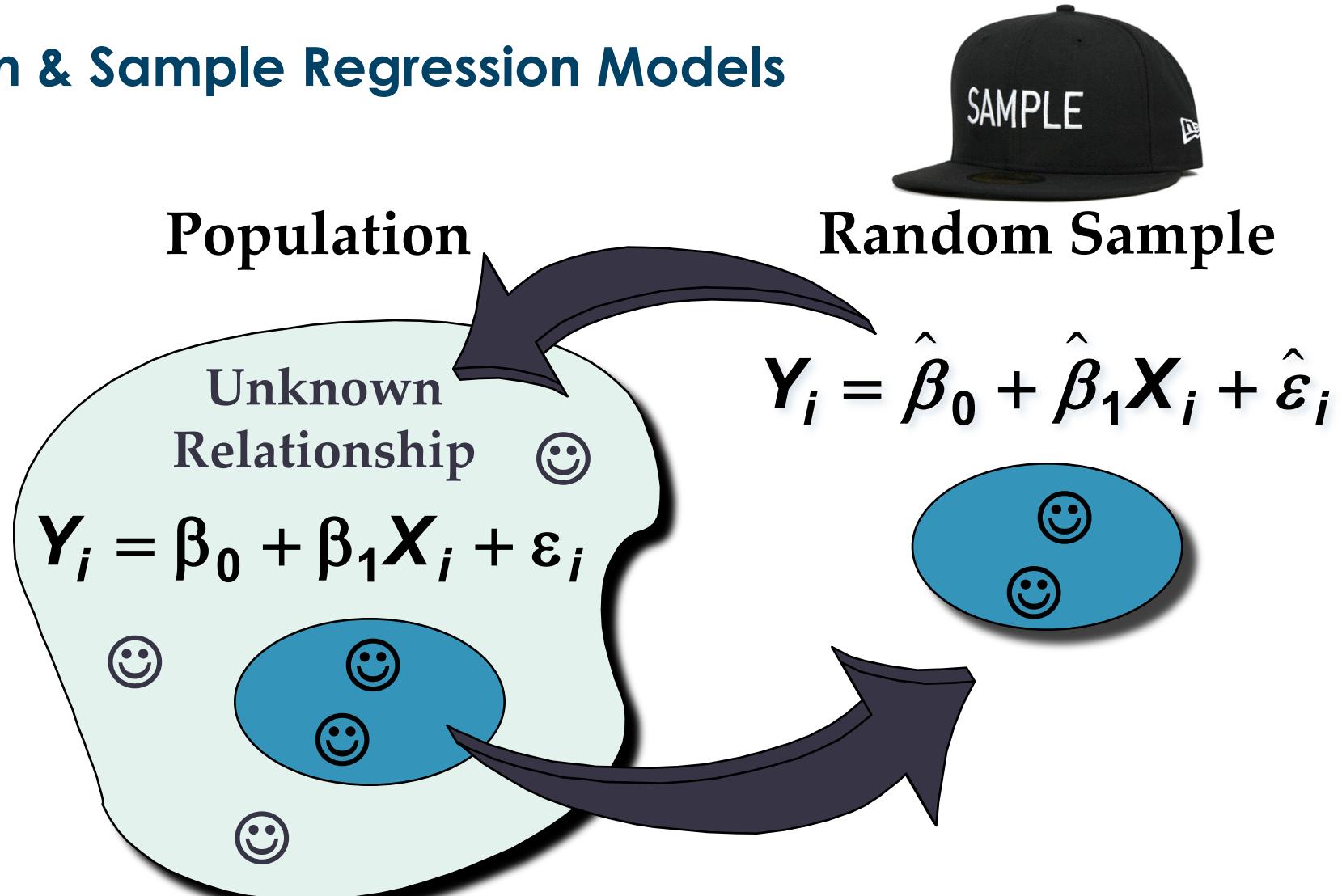
Population



Population & Sample Regression Models



Population & Sample Regression Models

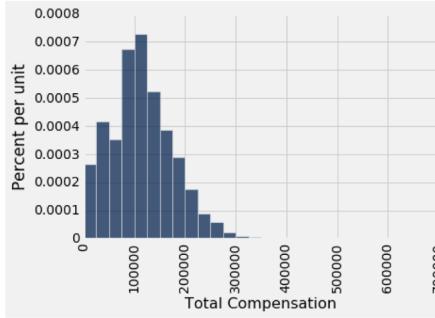


The Bootstrap

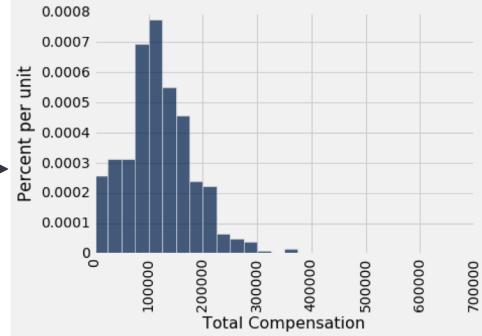
- A technique for simulating repeated random sampling
- All that we have is the original sample
 - ... which is large and random
 - Therefore, it probably resembles the population
- So we sample at random from the original sample!

Why the Bootstrap Works

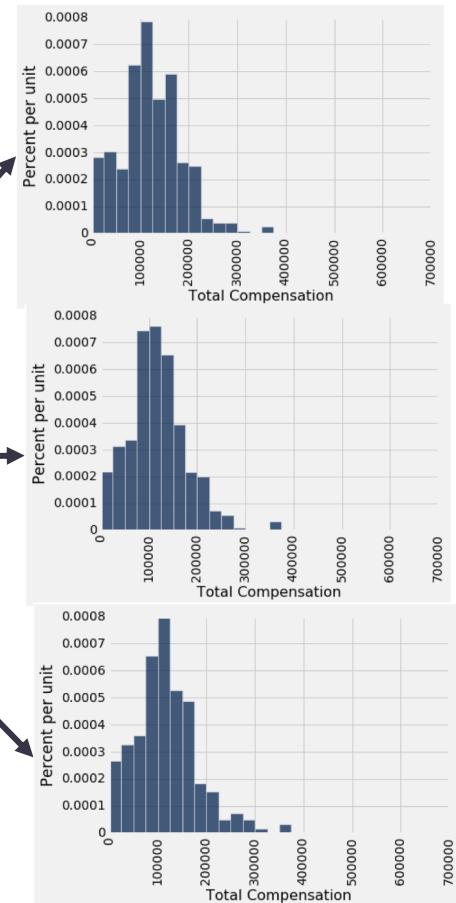
population



sample



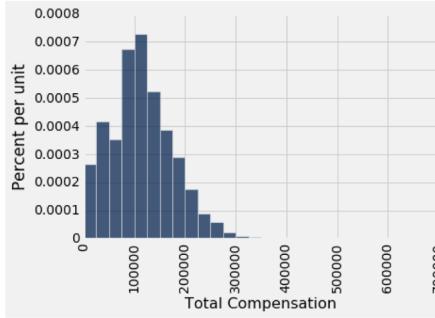
resamples



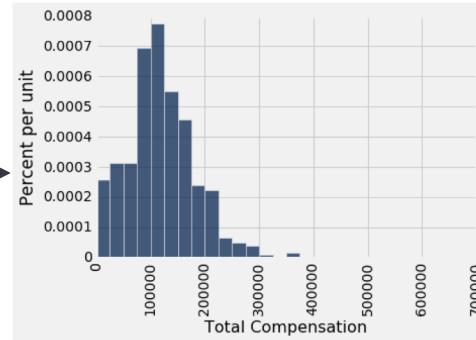
All of these look pretty similar, most likely.

Why We Need the Bootstrap

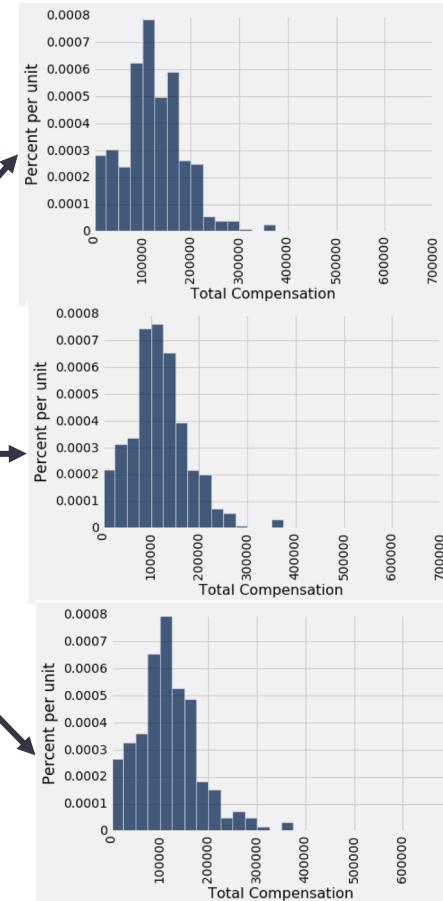
population



sample



resamples



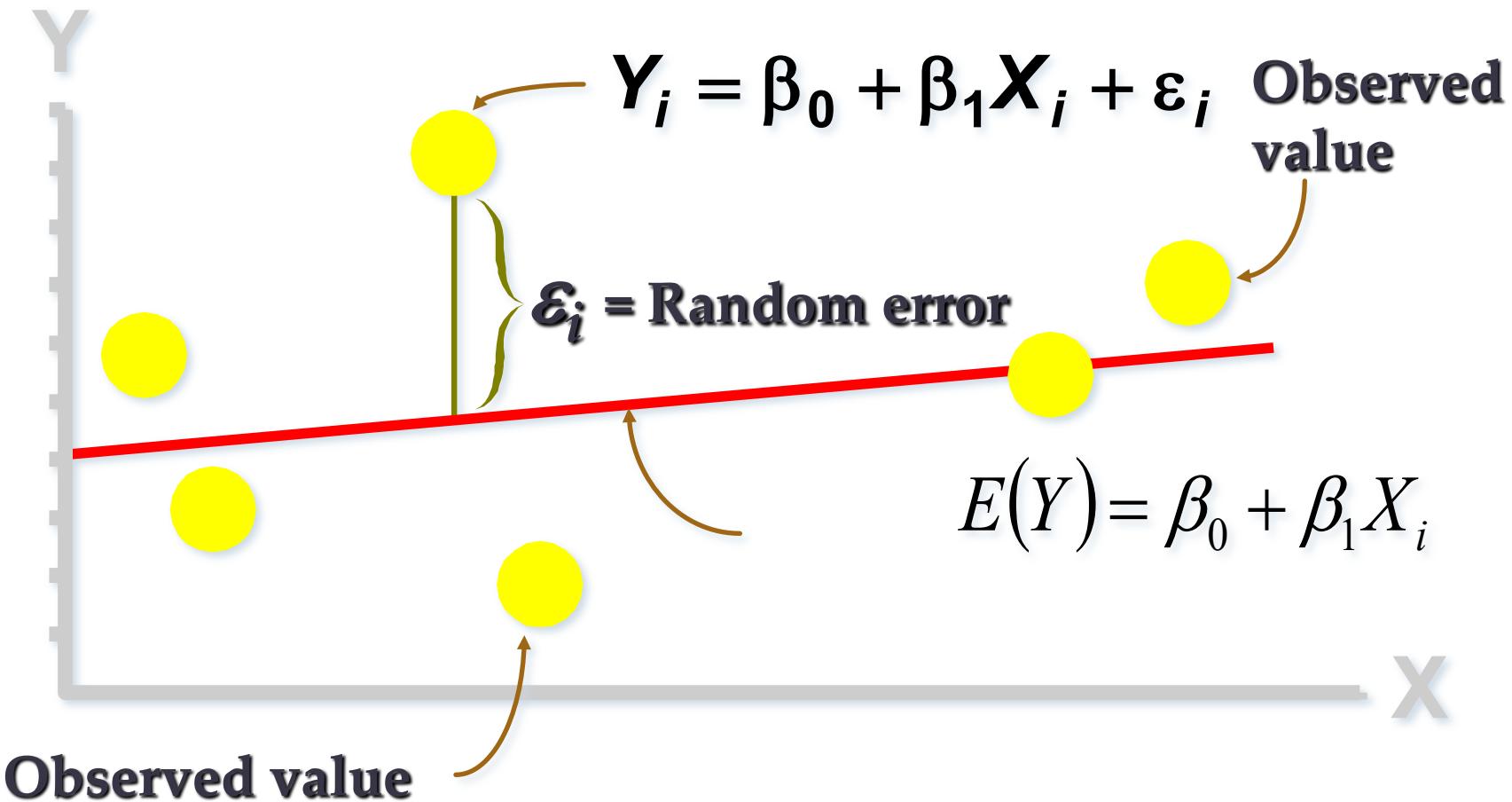
What we wish
we could get

What we
really get

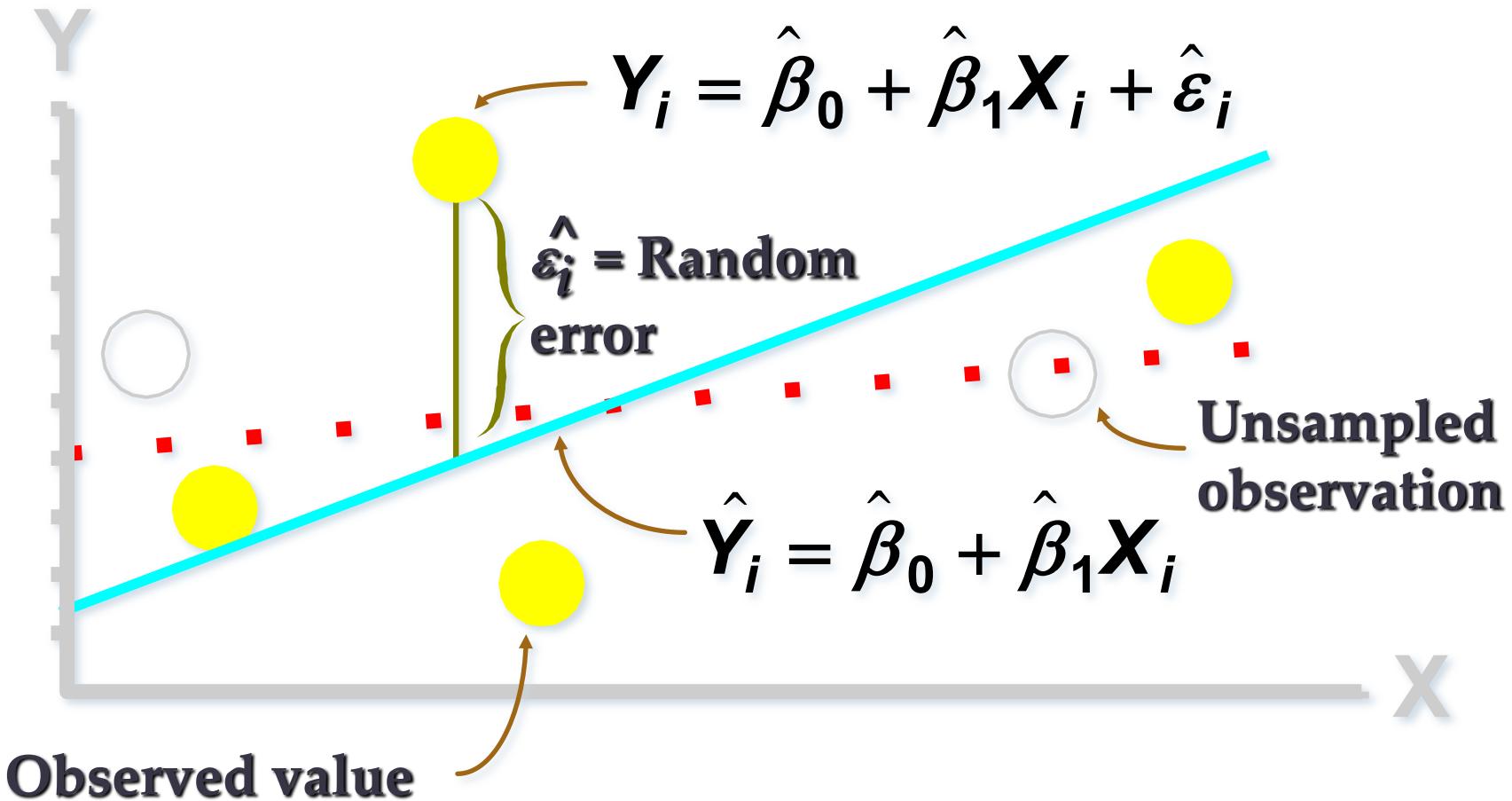
Key to Resampling

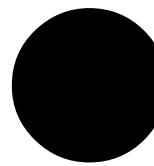
- From the original sample,
 - draw at random
 - with replacement
 - as many values as the original sample contained
- The size of the new sample has to be the same as the original one, so that the two estimates are comparable

Linear Regression



Sample Linear Regression Model

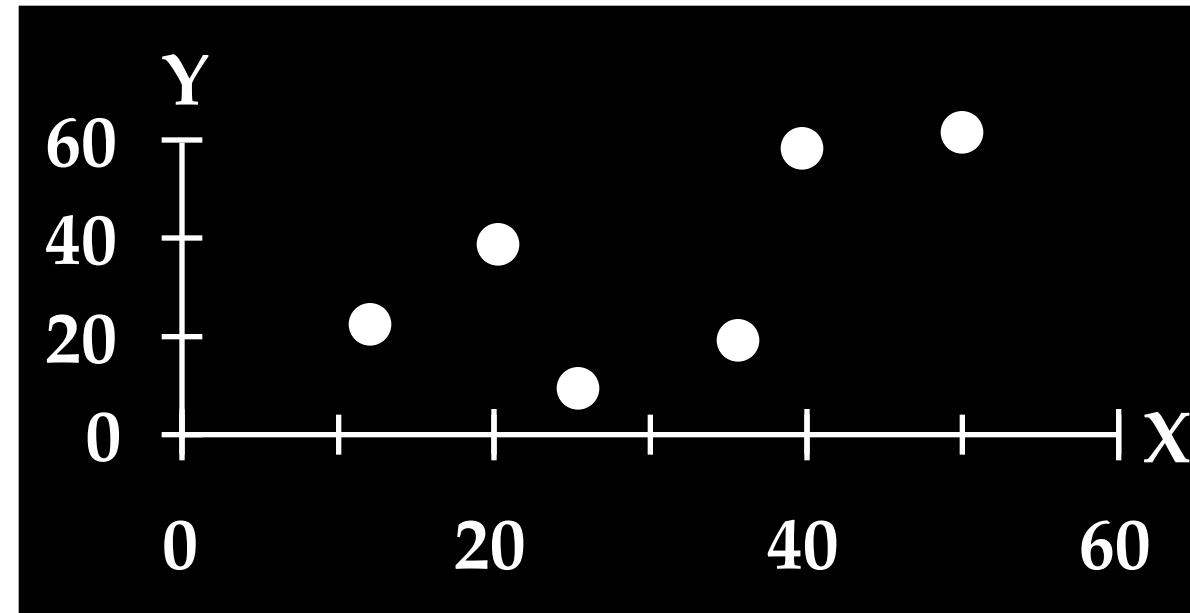




**Estimating Parameters:
Least Squares Method**

Scatter plot

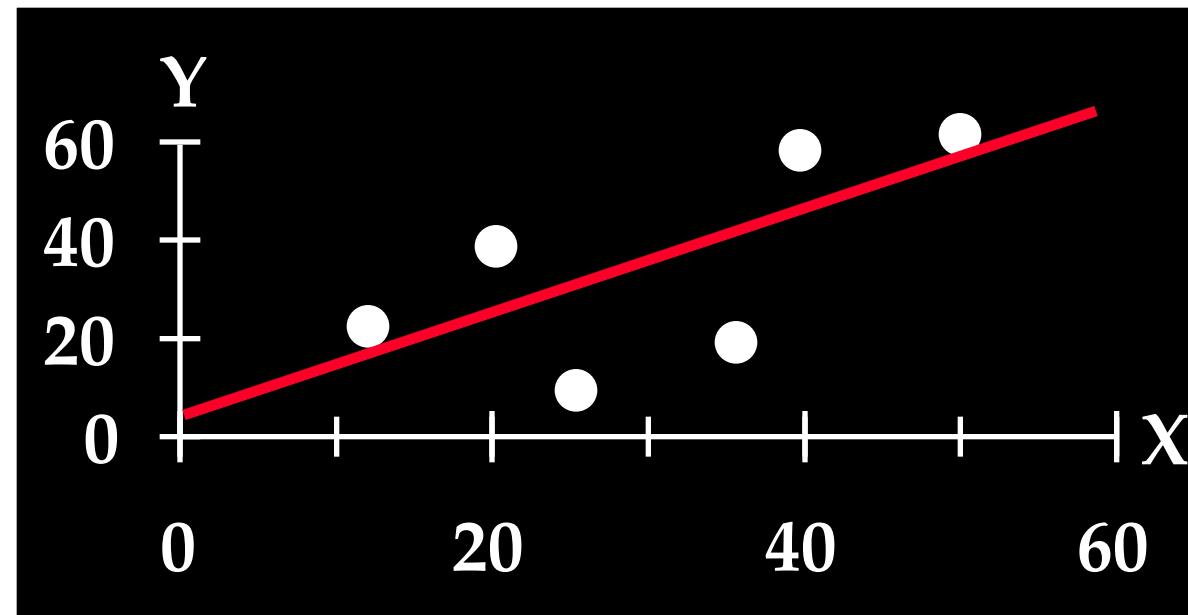
- Plot all (X_i, Y_i) pairs, and plot your learned model
- If you squint, suggests how well the model fits the data



Question

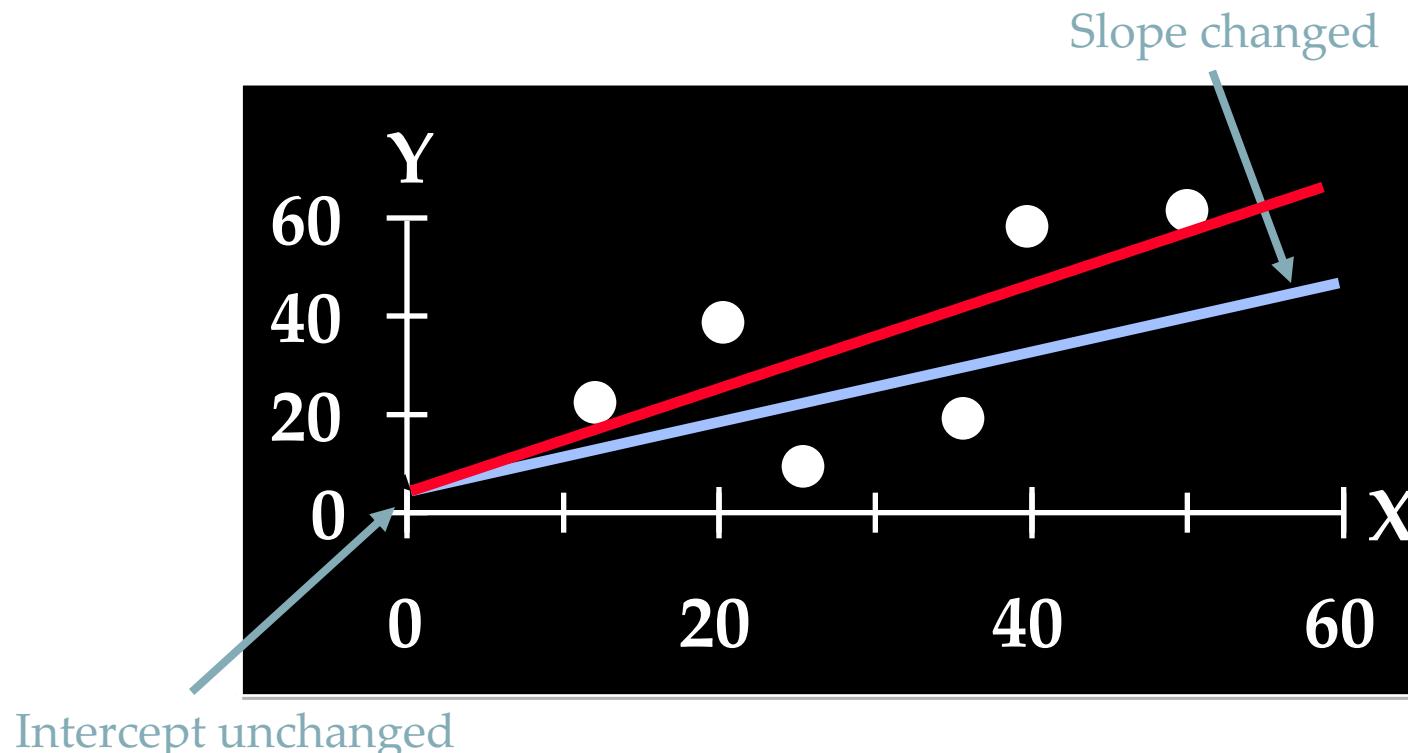
- How would you draw a line through the points?
- How do you determine which line “fits the best” ...?

??????????



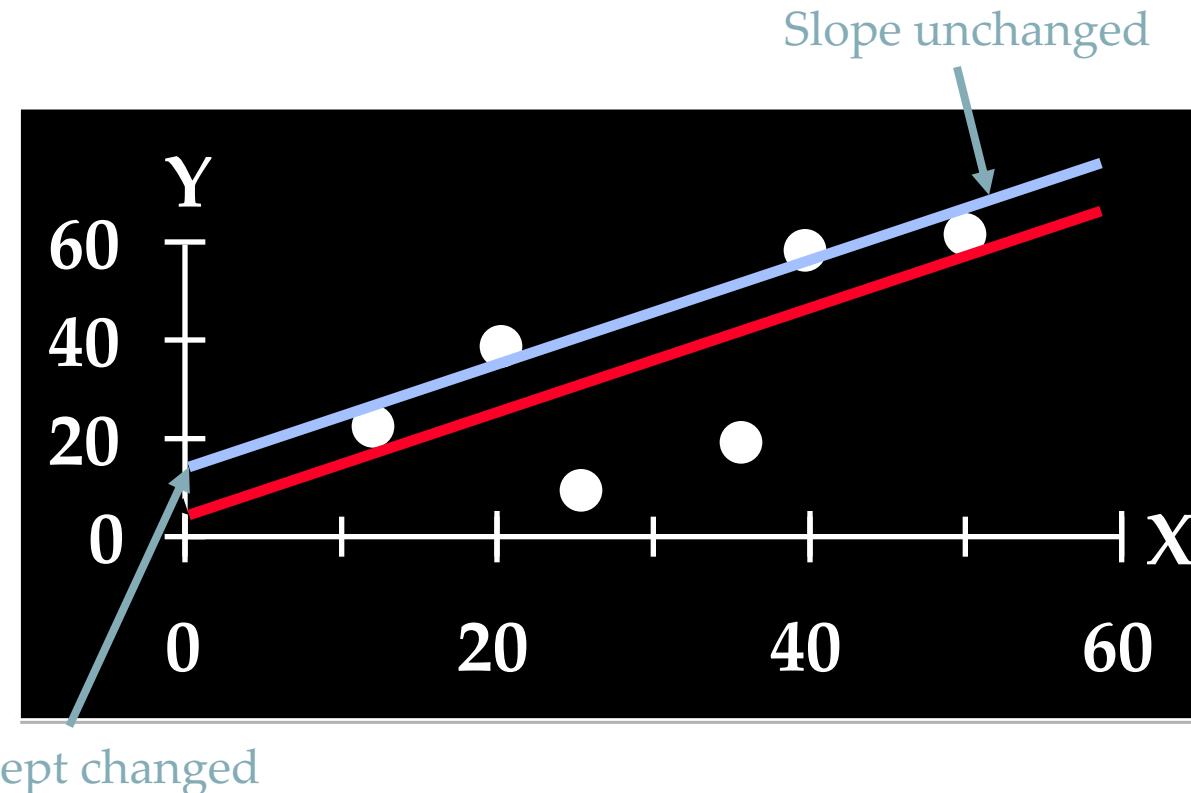
Question

- How would you draw a line through the points?
- How do you determine which line “fits the best” ?????????



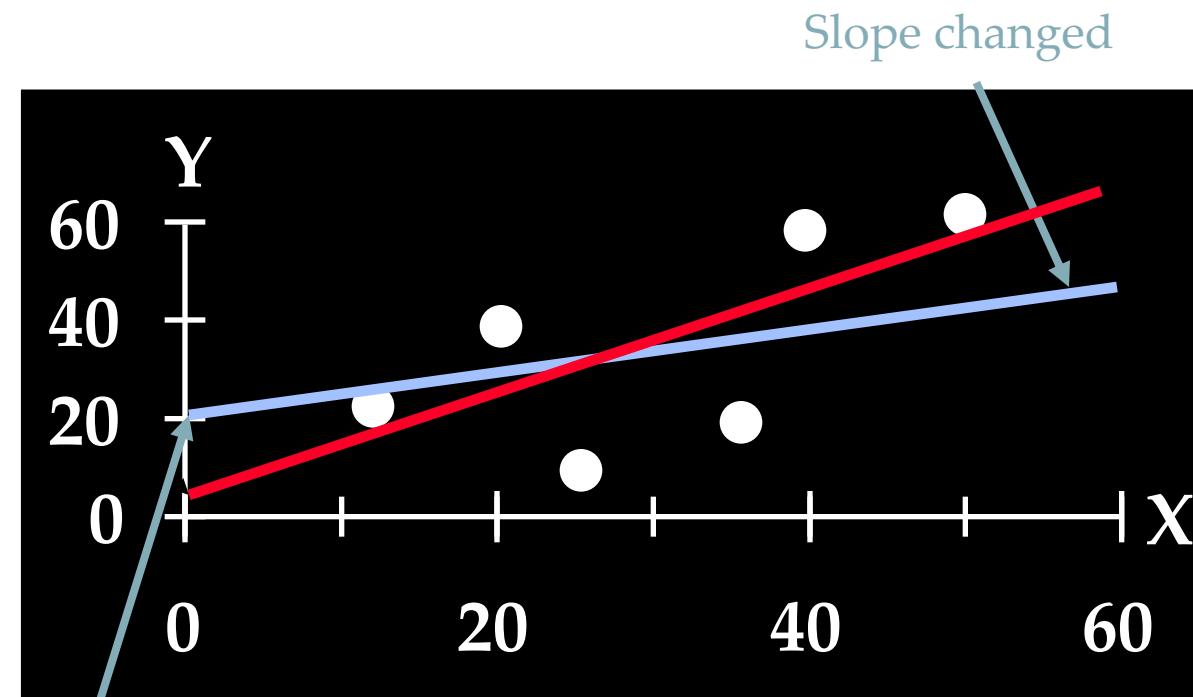
Question

- How would you draw a line through the points?
- How do you determine which line “fits the best” ?????????



Question

- How would you draw a line through the points?
- How do you determine which line “fits the best” ?????????



Least Squares

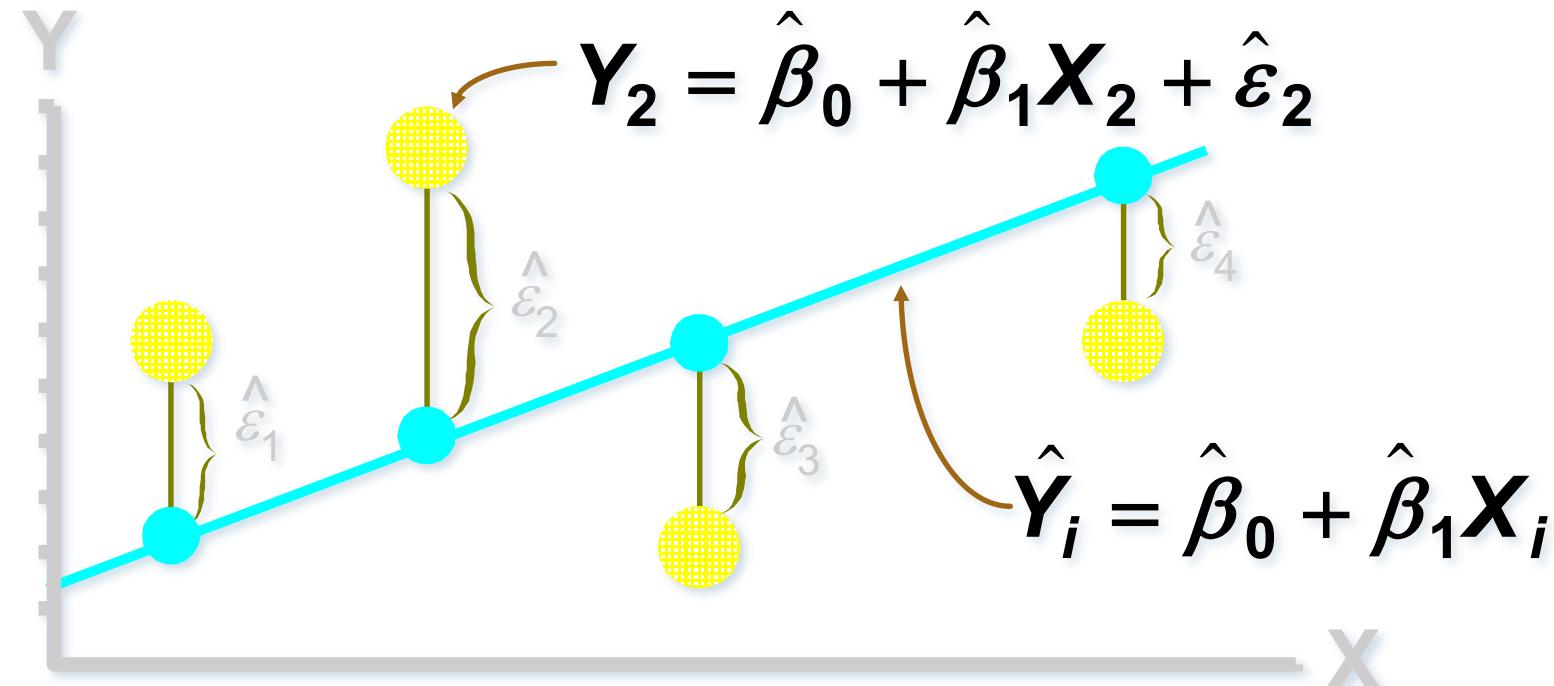
- **Best Fit:** difference between the true Y-values and the estimated Y-values is minimized:
 - Positive errors offset negative errors ...
 - ... square the error!

$$\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n \hat{\epsilon}_i^2$$

- Least squares minimizes the sum of the squared errors
 - Why squared? We'll cover this in more depth in a few weeks.
 - Until then: <http://www.benkuhn.net/squared>

Least Squares, Graphically

LS minimizes $\sum_{i=1}^n \hat{\varepsilon}_i^2 = \hat{\varepsilon}_1^2 + \hat{\varepsilon}_2^2 + \hat{\varepsilon}_3^2 + \hat{\varepsilon}_4^2$

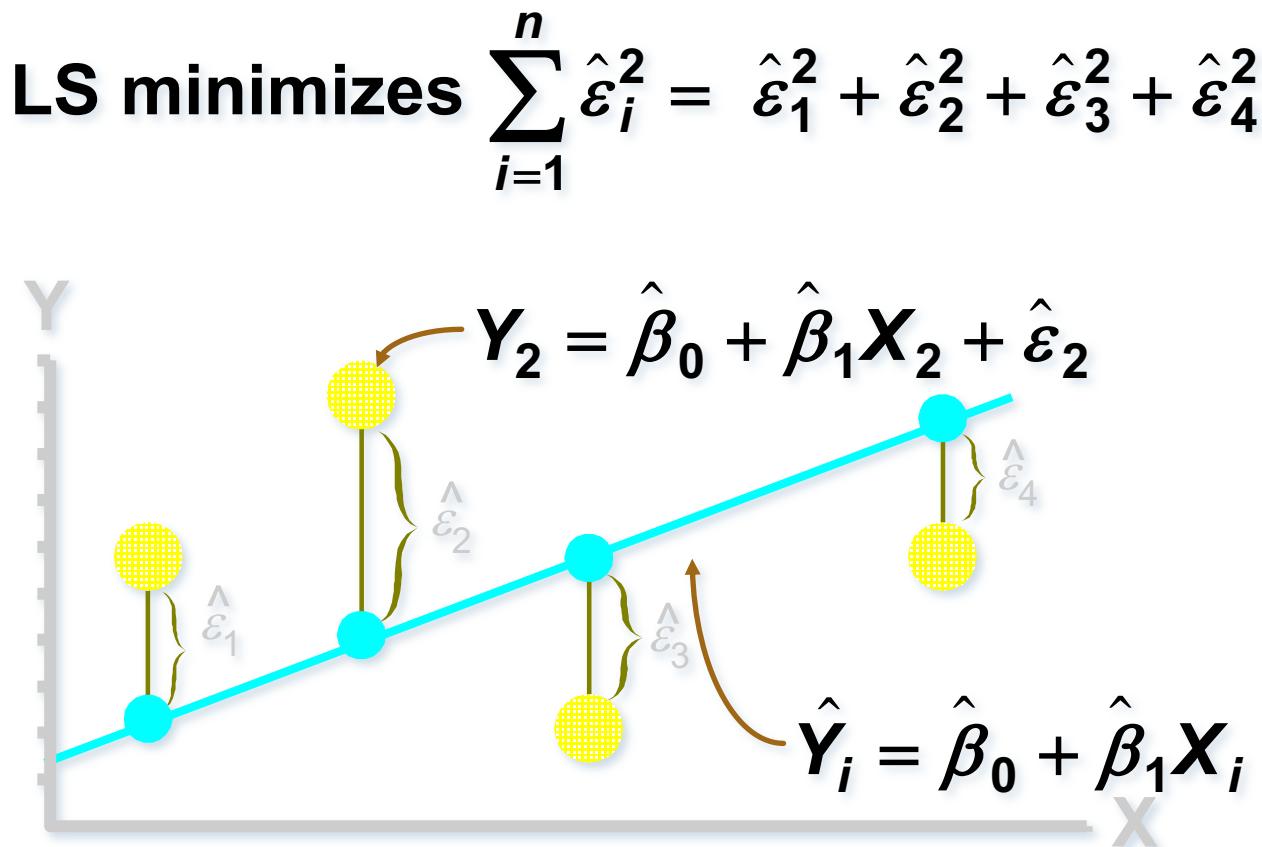


Interpretation of Coefficients

- Slope ($\hat{\beta}_1$):
 - Estimated Y changes by $\hat{\beta}_1$ for each unit increase in X .
 - If $\hat{\beta}_1 = 2$, then Y is expected to increase by 2 for each 1 unit increase in X .
- Y-Intercept ($\hat{\beta}_0$):
 - Average value of Y when $X = 0$.
 - If $\hat{\beta}_0 = 4$, then average Y is expected to be 4 when X is 0.

[WF]

In-depth
analysis to
come!





Now, Back to Missing Data ...

Example

- **Question:** Does the circumference of certain body parts predict BF%?
- **Assumption:** BF% is a linear function of measurements of various body parts and other features ...
- **Analysis:** Results from a regression model with BF% ...

Predictor	Estimate	S.E.	p-value
Age	0.0626	0.0313	0.0463
Neck	-0.4728	0.2294	0.0403
Forearm	0.45315	0.1979	0.0229
Wrist	-1.6181	0.5323	0.0026

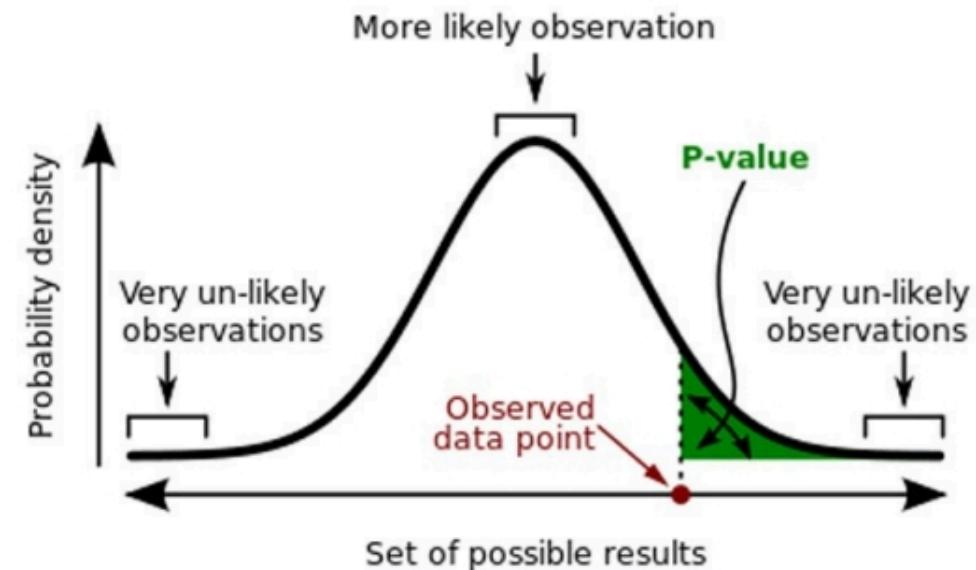
(Interpretation ????????????)

Hypothesis Testing

- One of the core ideas of Data Science – Should be at the core of all the analysis that you do.
 - The Null Hypothesis:
 - H_0 =The feature of interest has no effect on the target
 - The Alternative Hypothesis
 - H_1 =The feature of interest does have an effect on the target
- When you perform a hypothesis test you're testing if H_0 is true or not. If it is, then you say there's no relationship. If it's not, then you say 'we reject H_0 and accept H_1 '. This gets into some tricky logic that statisticians argue over as you actually aren't testing H_1 , but we're going to avoid that discussion here. The key thing is that you should be thinking deeply about what you're actually asking. What are your H_0 and H_1 ? Does your model actually test them? Can your data actually test them? The book explores how to calculate p-values, so

P-values

- In short a **p-value** is the probability of obtaining a result equal to or "more extreme" than what was actually observed, when the null hypothesis is actually true.
- Null Hypothesis:** There is no significant difference between the specified populations. The observed difference is due to sampling or experiment error.
- Typically, when we have a p-value below a certain threshold, say 0.05, then we can reject the null hypothesis and say that there is an effect.



A **p-value** (shaded green area) is the probability of an observed (or more extreme) result assuming that the null hypothesis is true.

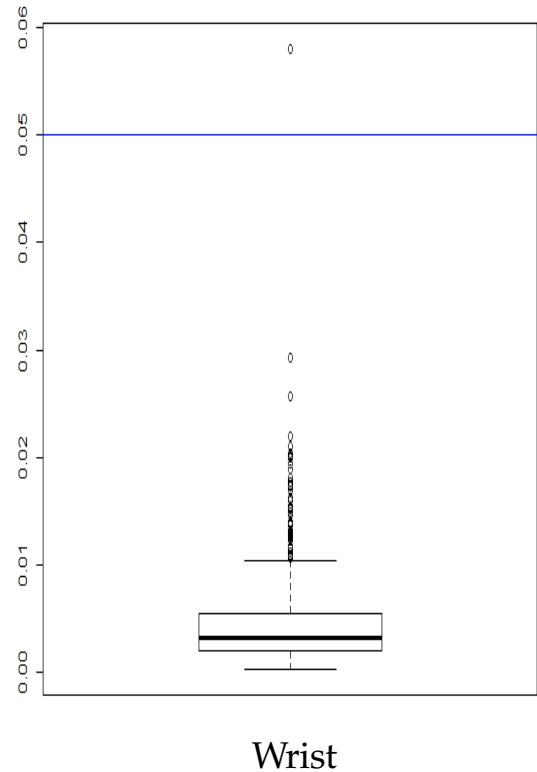
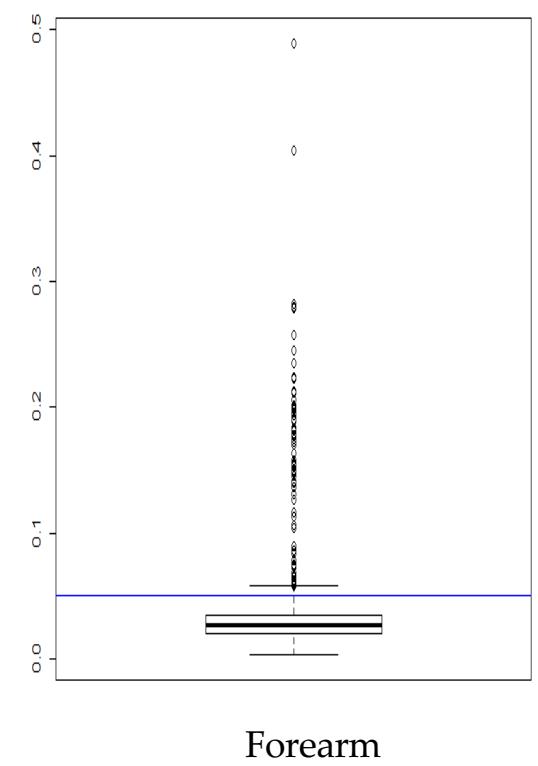
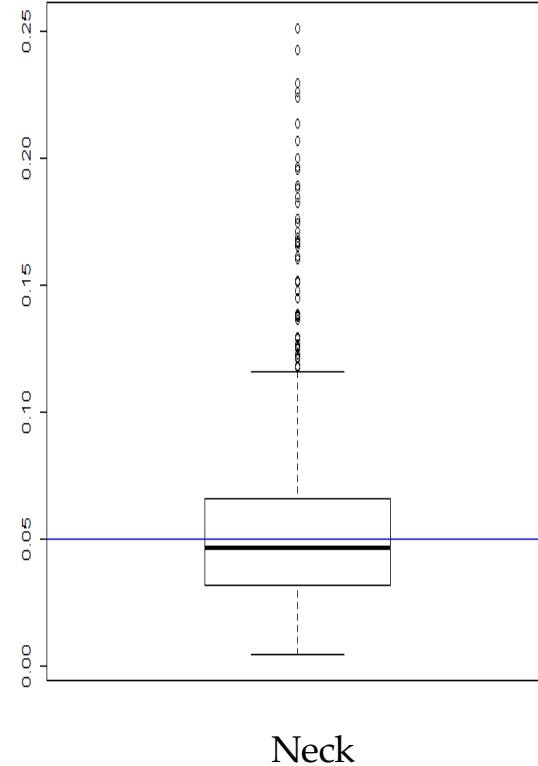
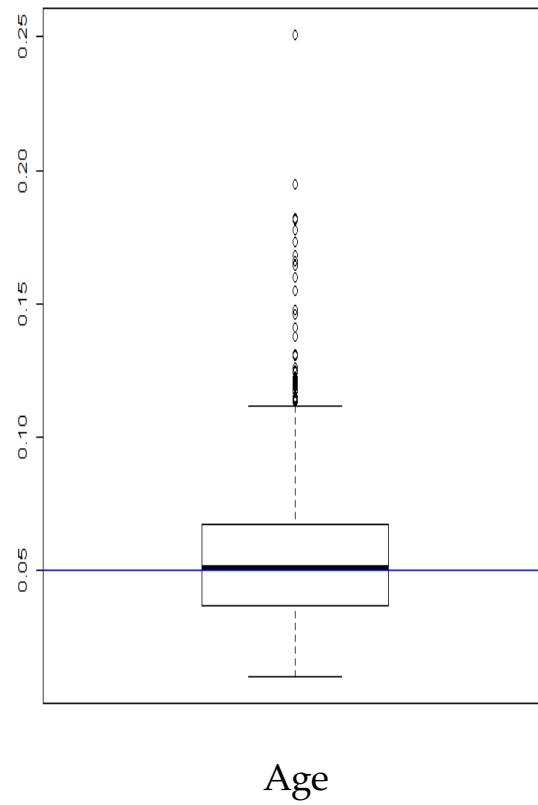
What If Data Were Missing?

- In this case, the dataset is complete:
 - But what if 5 percent of the participants had missing values? 10 percent? 20 percent?
- What if we performed complete case analysis and removed those who had missing values?
- Let's examine the effect if we do this if when the data is **missing completely at random (MCAR)**
 - Removed cases at random, reran analysis, stored the p-values
 - p-value: probability of getting at least as extreme a result as what we observed given that there is no relationship
 - Repeat 1000 times, plot p-values ...
 - Implemented in our notebook but with a different statistical test...



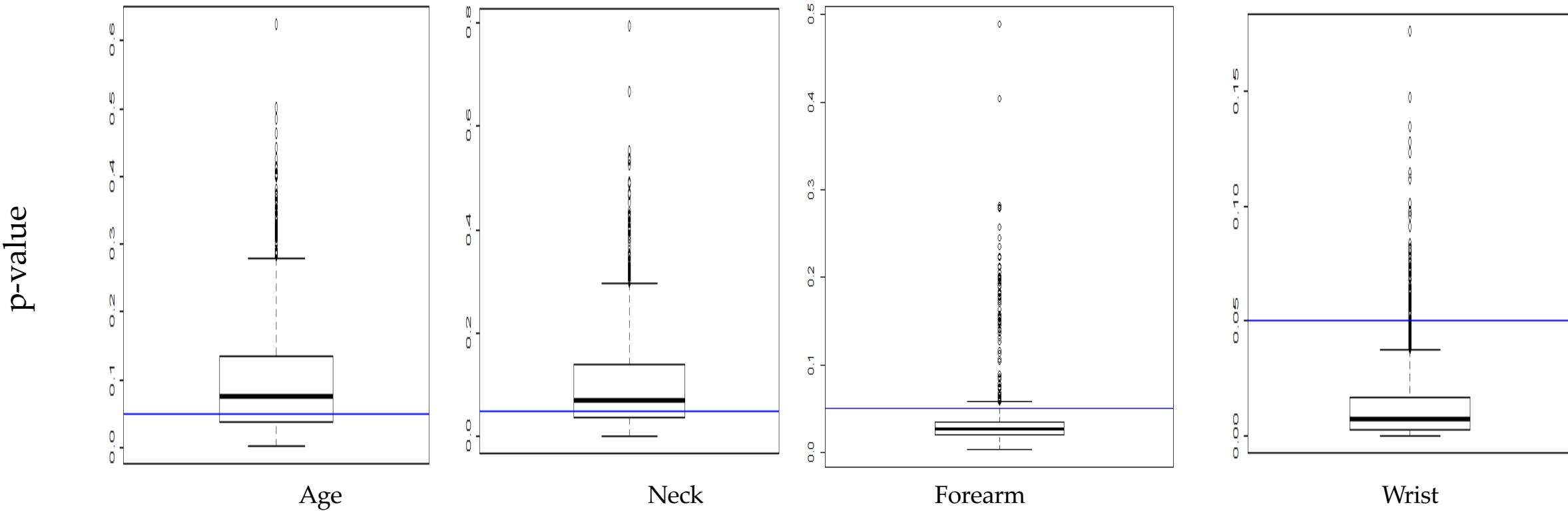
~5% Deleted (N=13)

p-value



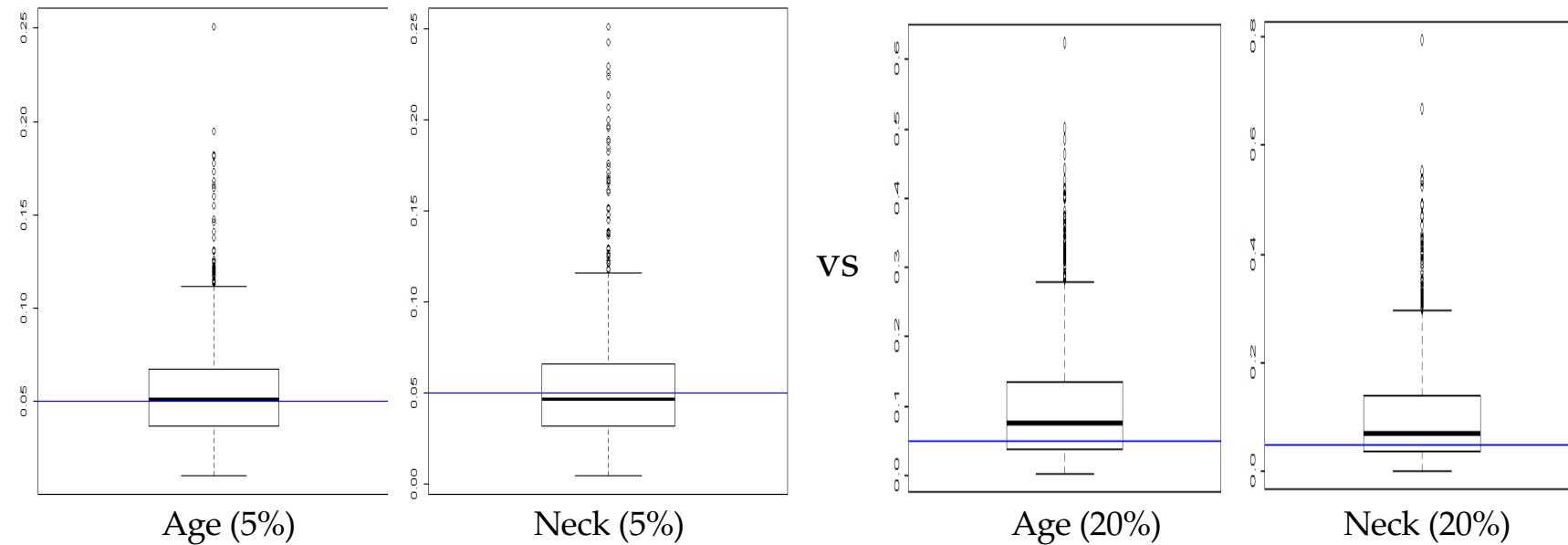
[JA]

~20% Deleted (N=50)



Conclusions seem to change ...

- Age/Neck: fail to reject the null hypothesis usually?



Still reject Forearm/Wrist most of the time

This is assuming the missing subjects' distribution does not differ from the non-missing. This would cause bias ...

Types Of Missing-ness

- Missing Completely at Random (MCAR)
- Missing at Random (MAR)
- Missing Not at Random (MNAR)

What Distinguishes Each Type of missing-ness?

- Suppose you're loitering outside of STH one day ...



Students just received their mid-semester grades

You start asking passing students their CMSC131 grades

- You don't force them to tell you or anything
- You also write down their gender and hair color

Your Sample

Hair Color	Gender	Grade
Red	M	A
Brown	F	A
Black	F	B
Black	M	A
Brown	M	
Brown	M	
Brown	F	
Black	M	B
Black	M	B
Brown	F	A
Black	F	
Brown	F	C
Red	M	
Red	F	A
Brown	M	A
Black	M	A

Summary:

- 7 students received As
- 3 students received Bs
- 1 student received a C

Nobody is failing!

- But 5 students did not reveal their grade ...

What influences a data point's Presence?

- Same dataset, but the values are replaced with a “0” if the data point is observed and “1” if it is not
- Question: for any one of these data points, what is the probability that the point is equal to “1” ...?
- What type of missing-ness do the grades exhibit?

Hair Color	Gender	Grade
0	0	0
0	0	0
0	0	0
0	0	0
0	0	1
0	0	1
0	0	1
0	0	0
0	0	0
0	0	0
0	0	1
0	0	0
0	0	1
0	0	0
0	0	0
0	0	0

MCAR: Missing Completely at Random

- If this probability is not dependent on any of the data, observed or unobserved, then the data is Missing Completely at Random (MCAR)
- Suppose that X is the observed data and Y is the unobserved data. Call our “missing matrix” R
- Then, if the data are MCAR, $P(R|X,Y) = ??????????$

•

$$P(R|X,Y) = P(R)$$

- Probability of those rows missing is independent of anything.

Totally Realistic MCAR Example



- You are running an experiment on plants grown in pots, when suddenly you have a nervous breakdown and smash some of the pots
- You will probably not choose the plants to smash in a well-defined pattern, such as height age, etc.
- Hence, the missing values generated from your act of madness will likely fall into the MCAR category

Applicability of MCAR

- A completely random mechanism for generating missing-ness in your data set just isn't very realistic
- Usually, missing data is missing for a reason:
 - Maybe older people are less likely to answer web-delivered questions on surveys
 - In longitudinal studies people may die before they have completed the entire study
 - Companies may be reluctant to reveal financial information

MAR: Missing at Random

- Missing at Random (MAR): probability of missing data is dependent on the observed data but not the unobserved data
- Suppose that X is the observed data and Y is the unobserved data. Call our “missing matrix” R
- Then, if the data are MAR, $P(R|X,Y) = ??????????$

$$\bullet \\ P(R|X,Y) = P(R|X)$$

- Not exactly random (in the vernacular sense).
- There is a probabilistic mechanism that is associated with whether the data is missing
- Mechanism takes the observed data as input

Examples?



MAR: Key Point

- We can model that latent mechanism and compensate for it
- Imputation: replacing missing data with substituted values
- Models today will assume MAR
- Example: if age is known, you can model missing-ness as a function of age
- Whether or not missing data is MAR or the next type, Missing Not at Random (MNAR), is not* testable.
- Requires you to “understand” your data

*unless you can get the missing data (e.g., post-study phone calls)

MNAR: Missing Not at Random

- MNAR: missing-ness has something to do with the missing data itself
- Examples: ??????????
- Do you binge drink? Do you have a trust fund? Do you use illegal drugs? What is your sexuality? Are you depressed?
- Said to be “non-ignorable”:
- Missing data mechanism must be considered as you deal with the missing data
- Must include model for why the data are missing, and best guesses as to what the data might be

Back to CSIC ...

- Is the missing data:
- MCAR;
- MAR; or
- MNAR?
- ????????????



Hair Color	Gender	Grade
Red	M	A
Brown	F	A
Black	F	B
Black	M	A
Brown	M	
Brown	M	
Brown	F	
Black	M	B
Black	M	B
Brown	F	A
Black	F	
Brown	F	C
Red	M	
Red	F	A
Brown	M	A
Black	M	A

Add a variable

- Bring in the GPA:
- Does this change anything?

Hair Color	GPA	Gender	Grade
Red	3.4	M	A
Brown	3.6	F	A
Black	3.7	F	B
Black	3.9	M	A
Brown	2.5	M	
Brown	3.2	M	
Brown	3.0	F	
Black	2.9	M	B
Black	3.3	M	B
Brown	4.0	F	A
Black	3.65	F	
Brown	3.4	F	C
Red	2.2	M	
Red	3.8	F	A
Brown	3.8	M	A
Black	3.67	M	A

Handling Missing Data ...



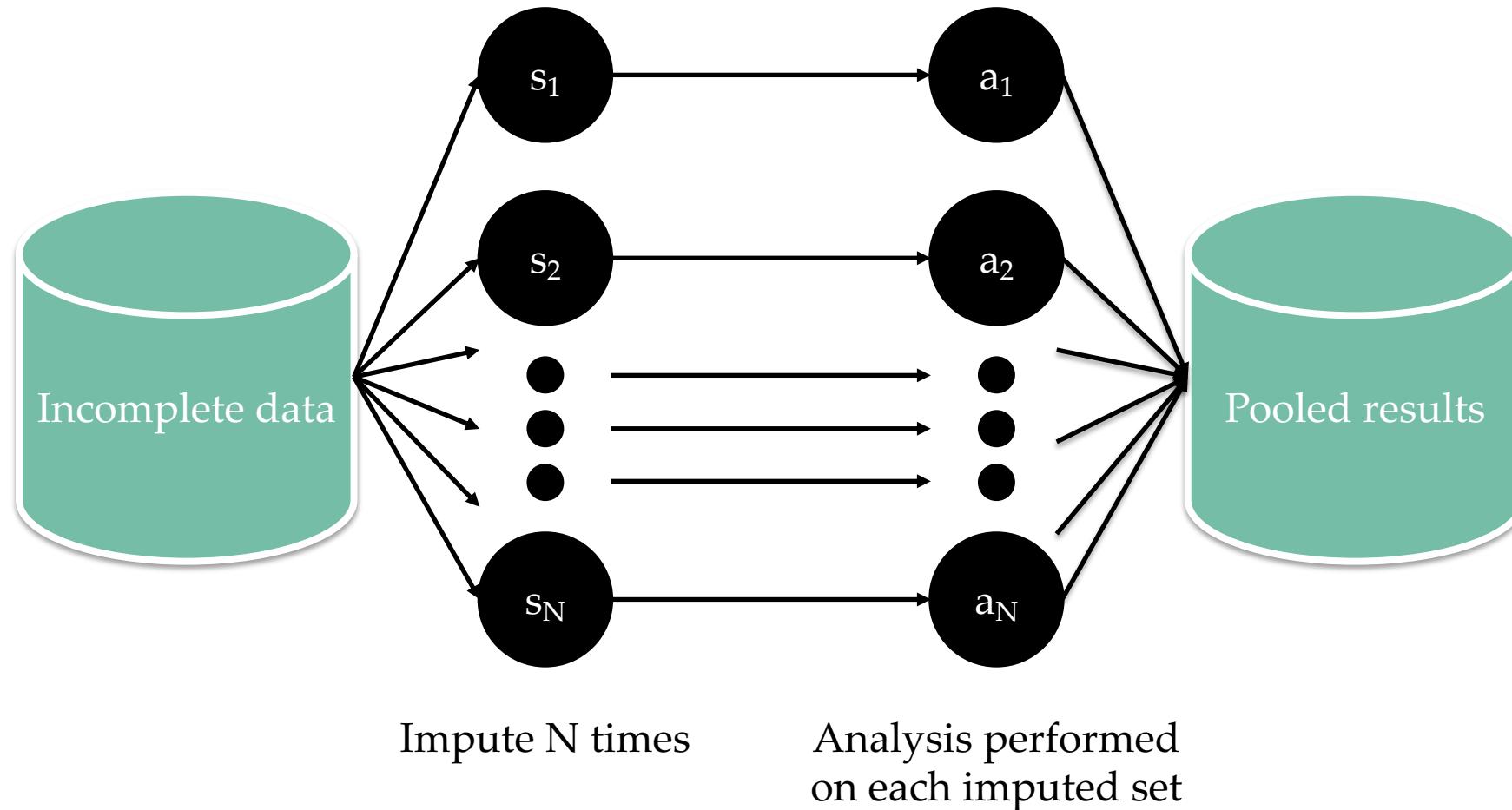
Single Imputation

- Mean imputation: imputing the average from observed cases for all missing values of a variable
- Hot-deck imputation: imputing a value from another subject, or “donor,” that is most like the subject in terms of observed variables
- Last observation carried forward (LOCF): order the dataset somehow and then fill in a missing value with its neighbor
- Cold-deck imputation: bring in other datasets
- Old and busted:
 - All fundamentally impose too much precision.
 - Have uncertainty over what unobserved values actually are
 - Developed before cheap computation

Multiple Imputation

- Developed to deal with noise during imputation
- Impute once → treats imputed value as observed
- We have uncertainty over what the observed value would have been
- Multiple imputation: generate several random values for each missing data point during imputation

Imputation Process



Tiny example

X	Y
32	2
43	?
56	6
25	?
84	5

Independent variable: X

Dependent variable: Y

We assume Y has a linear relationship with X

Let's Impute Some Data!

- Use a predictive distribution of the missing values:
- Given the observed values, make random draws of the observed values and fill them in.
- Do this N times and make N imputed datasets

X	Y
32	2
43	5.5
56	6
25	8
84	5

X	Y
32	2
43	7.2
56	6
25	1.1
84	5

For very large values of N=2 ...

Inference with Multiple Imputation

- Now that we have our imputed data sets, how do we make use of them? ????????????
- Analyze each of the separately

X	Y
32	2
43	5.5
56	6
25	8
84	5

Slope -0.8245
 Standard error 6.1845

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

X	Y
32	2
43	7.2
56	6
25	1.1
84	5

Slope 4.932
 Standard error 4.287

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

Pooling analyses

- Pooled slope estimate is the average of the N imputed estimates
- Our example, $\beta_{1p} = \frac{\beta_{11} + \beta_{12}}{2} = (4.932 - .8245) \times 0.5 = 2.0538$
- The pooled slope variance is given by
- $s^2 = \frac{\sum z_i}{m} + \left(1 + \frac{1}{m}\right) \times \frac{1}{m-1} * \sum (\beta_{1i} - \beta_{1p})^2$
- Where Z_i is the standard error of the imputed slopes
- Our example: $(4.287 + 6.1845)/2 + (3/2)*(16.569) = 30.08925$
- Standard error: take the square root, and we get 5.485

Predicting the missing data given the observed data

- Given events A, B; and $P(A) > 0 \dots$

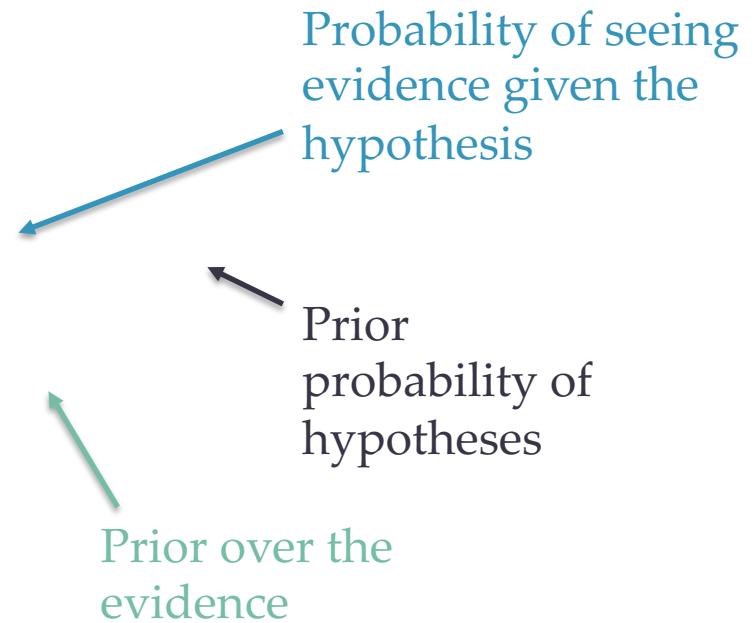
- Bayes' Theorem:

- $$P(B|A) = \frac{P(A|B)*P(B)}{P(A)}$$

- In our case:

- $$P(H|E) = \frac{P(E|H)*P(H)}{P(E)}$$

Posterior probability of the hypothesis given the evidence



Bayesian Imputation

- Establish a prior distribution:
- Some distribution of parameters of interest θ before considering the data, $P(\theta)$
- We want to estimate θ
- Given θ , can establish a distribution $P(X_{obs}/\theta)$
- Use Bayes Theorem to establish $P(\theta/X_{obs}) \dots$
- Make random draws for θ
- Use these draws to make predictions of Y_{miss}

How big should N Be?

- Number of imputations N depends on:
 - Size of dataset
 - Amount of missing data in the dataset
- Some previous research indicated that a small N is sufficient for efficiency of the estimates, based on:
 - $(1 + \frac{\lambda}{N}) - 1$
 - N is the number of imputations and λ is the fraction of missing information for the term being estimated [Schaffer 1999]
 - More recent research claims that a good N is actually higher in order to achieve higher power [Graham et al. 2007]



More advanced methods

- Interested? Further reading:
- Regression-based MI methods
- Multiple Imputation Chained Equations (MICE) or Fully Conditional Specification (FCS)
 - Readable summary from JHU School of Public Health:
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3074241/>
- Markov Chain Monte Carlo (MCMC)
 - We'll cover this a bit, but also check out CMSC422!