

Class 3: - Missing Values, Functions, Group Variables, and Intro to ML

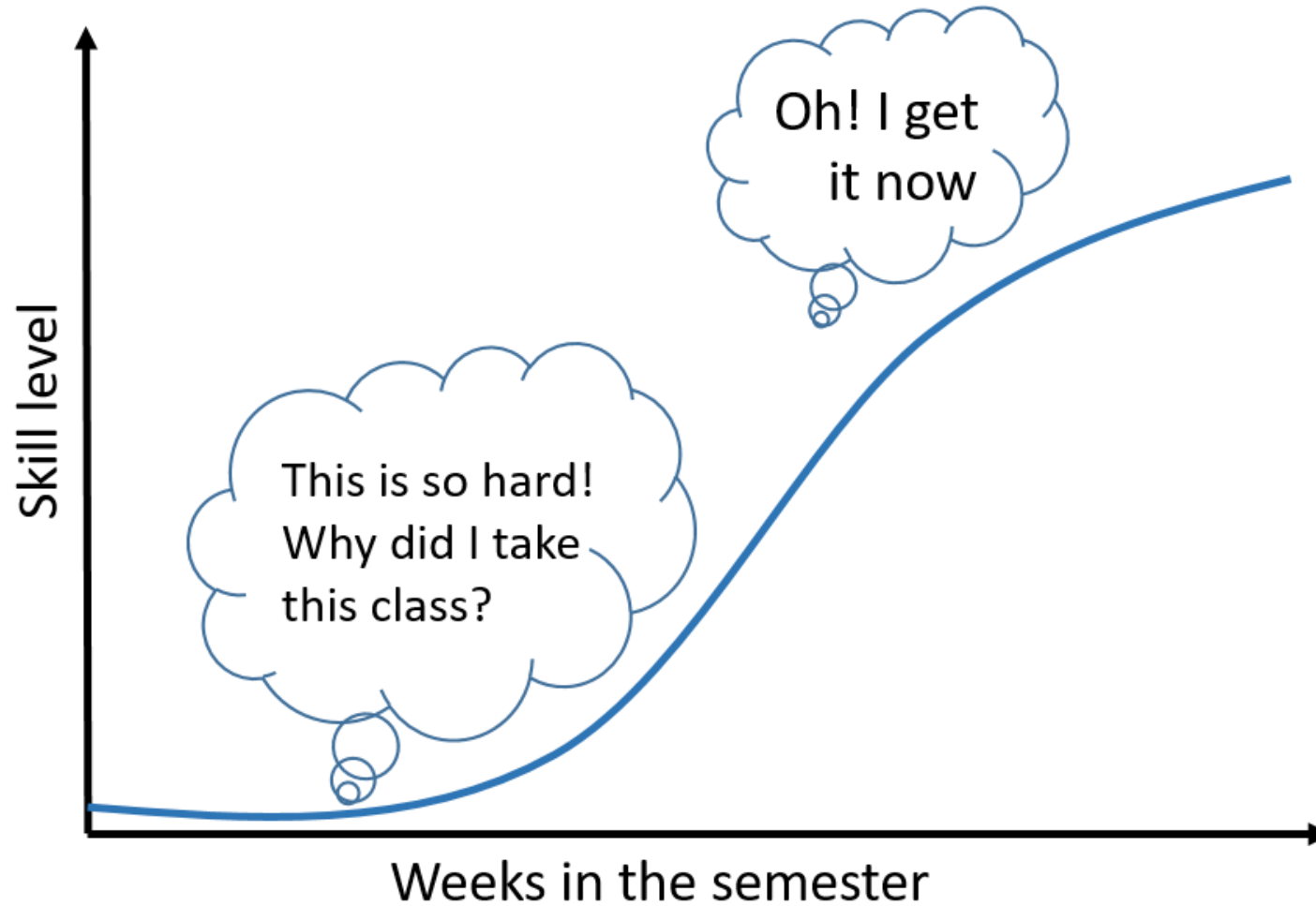
BUS 696

Prof. Jonathan Hersh

Class 3: Announcements

1. Data Analytics Week! October 5-9
 - May use, don't have to
2. Office Hours
 - TA: Wed 12-1; Thur 5-6
 - Mon 11-12noon; Wed 5-6;
3. Problem Set 1 posted – Due Sept 25
 - Must submit compiled HTML file
using Rmarkdown
4. Posted Rmarkdown Homework
Template

I apologize but we are at peak difficulty



Data Analytics Industry Week

Register on Handshake to get access to the following virtual events!

Data Analytics Accelerator Program Info Session

Monday, October 5 | 11 a.m. PST

Interested in pursuing a career in the growing field of data analytics? The Argyros School of Business is proud to present the new career skills-focused Analytics Accelerator Program. Learn more about what hard skills are needed to land a successful career in data analytics. Hear from Professor Toplansky and Dr. Hersh about how you can propel your success and prepare for 21st Century jobs that pay a premium.

Careers in Data Analytics

Tuesday, October 6 | 12 p.m. PST

Hear from the renowned authors of Build a Career in Data Science, Jacqueline Nolis and Emily Robinson about careers in data analytics.

Data Analytics Industry Panel

Thursday, October 8 | 4:30 p.m. PST

This data analytics panel will feature industry experts in analytics from entertainment, healthcare, technology, and more.

Entertainment Analytics: Turning Data Into Insights

Friday, October 9 | 12 p.m. PST

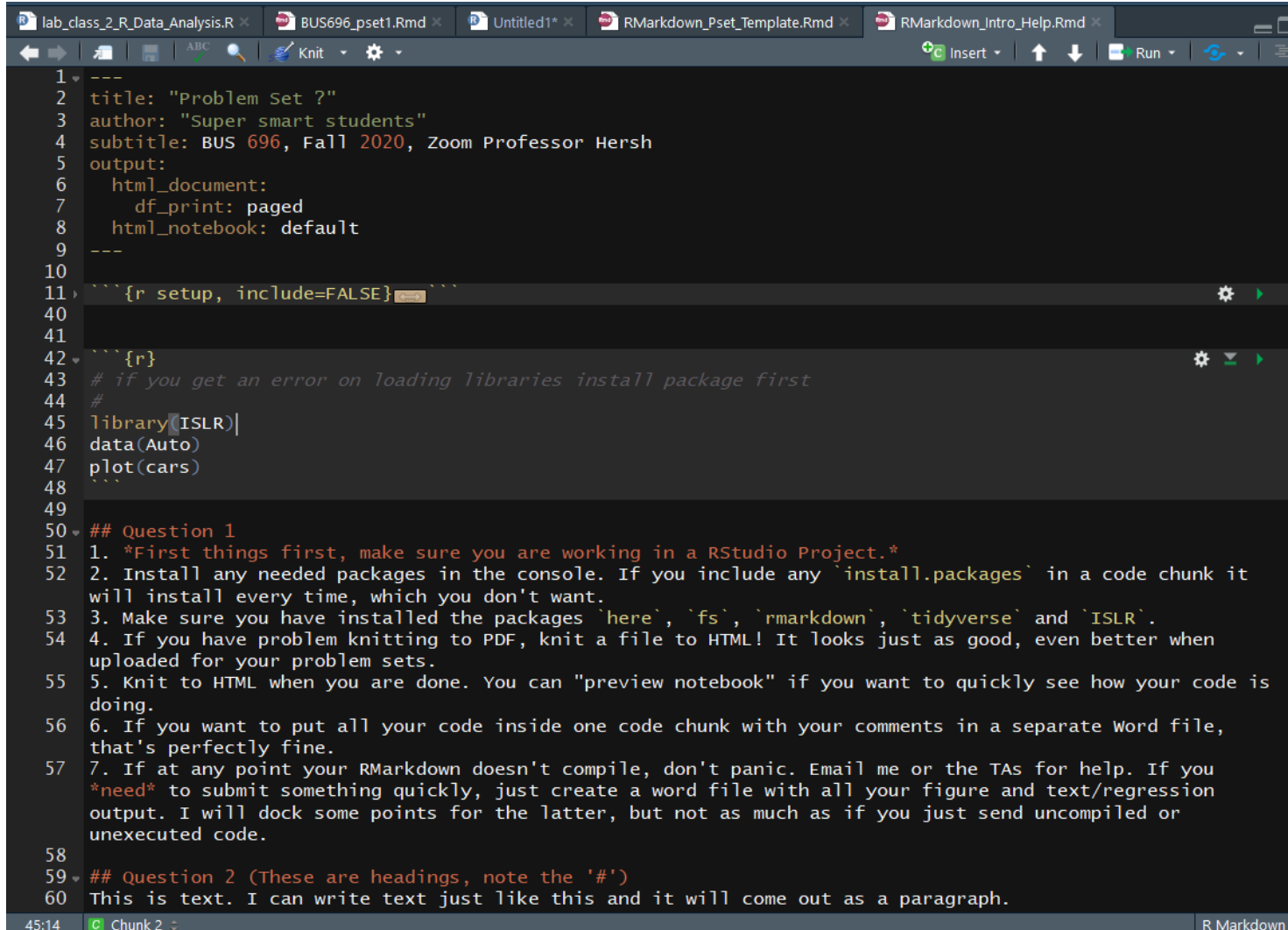
Come see a live demo and learn about turning data into actionable insights in Entertainment Analytics with Andre Vargas Head of the data department at leading entertainment and sports agency, Creative Artists Agency (CAA).

May Use Problem Set Rmarkdown Template In Assignments Module

```
lab_class_2_R_Data_Analysis.R x BUS696_pset1.Rmd x Untitled1* x RMarkdown_Pset_Template.Rmd x
1 ---
2 title: "Problem Set ?"
3 author: "Super smart students"
4 subtitle: BUS 696 Problem Set Template
5 output:
6   html_document:
7     df_print: paged
8   html_notebook: default
9 ---
10
11 {r_setup, include=FALSE}
12
13
14
15 {r_setup-2}
16
17 # load all your libraries here
18 library('tidyverse')
19 # note, do not run install.packages() inside a code chunk. install them in the console outside of a code
20 chunk.
21
22
23
24 ## Question 1
25
26 1a) Response to part a.
27
28 {r}
29
30 # code for part a
31
32
33
34 1b) Response to part b.
35
36
```

Assignments			
	Problem Set 1 (Basic R Programming, Data Manipulation, Plotting and Statistical Uncertainty)	✓	
	Sep 25 30 pts		
	Problem Set 2 (Linear Regression and Classification)	⊘	
	Oct 7 30 pts		
	RMarkdown_Pset_Template.Rmd	✓	
	RMarkdown_Intro_Help.Rmd	✓	

IF You Want to Learn More About Rmarkdown Go Through Rmarkdown Help



```
1 ---
2 title: "Problem Set ?"
3 author: "Super smart students"
4 subtitle: BUS 696, Fall 2020, Zoom Professor Hersh
5 output:
6   html_document:
7     df_print: paged
8   html_notebook: default
9 ---
10
11 {r setup, include=FALSE}
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42 {r}
43 # if you get an error on loading libraries install package first
44 #
45 library(ISLR)
46 data(Auto)
47 plot(cars)
48
49
50 ## Question 1
51 1. *First things first, make sure you are working in a RStudio Project.*
52 2. Install any needed packages in the console. If you include any `install.packages` in a code chunk it
53    will install every time, which you don't want.
54 3. Make sure you have installed the packages `here`, `fs`, `rmarkdown`, `tidyverse` and `ISLR`.
55 4. If you have problem knitting to PDF, knit a file to HTML! It looks just as good, even better when
56    uploaded for your problem sets.
57 5. Knit to HTML when you are done. You can "preview notebook" if you want to quickly see how your code is
58    doing.
59 6. If you want to put all your code inside one code chunk with your comments in a separate Word file,
60    that's perfectly fine.
61 7. If at any point your RMarkdown doesn't compile, don't panic. Email me or the TAs for help. If you
62    *need* to submit something quickly, just create a word file with all your figure and text/regression
63    output. I will dock some points for the latter, but not as much as if you just send uncompiled or
64    unexecuted code.
65
66 ## Question 2 (These are headings, note the '#')
67 This is text. I can write text just like this and it will come out as a paragraph.
```

Class 3: Outline

1. Qs from last week?

2. Basic Data Analysis

- Missing values
- Loops
- mutate to transform variables
- Remove duplicates with distinct
- Outputting “clean” data file”

3. Data Analysis by Groups

- group_by() function
- summarize() to create group variables

4. Data Analysis Lab Class 3

5. P-values

6. Introductory Machine Learning Concepts

Missing Values

```
# -----  
# MISSING VALUES are values that are unknown in your dataset  
# -----  
# R stores missing values as NAs  
is.na(NA)  
1 > NA  
1 + 1 == NA  
NA == NA  
y <- NA  
y  
x <- 1  
y == x
```

lab_class_4_R_Exploratory_Data_Analysi... x					movies x				
Filter									
	actor_1_facebook_likes	gross	genres	actor_1_name					
	11000	200074175	Action Adventure Thriller	Christoph Waltz					
	27000	448130642	Action Thriller	Tom Hardy					
	131	NA	Documentary	Doug Walker					
	640	73058679	Action Adventure Sci-Fi	Daryl Sabara					
	24000	336530303	Action Adventure Romance	J.K. Simmons					
	799	200807262	Adventure Animation Comedy Family Fantasy Musical ...	Brad Garrett					

Loops in R

```
# -----  
# LOOP through numbers using the FOR loop  
# -----  
  
# for loops are created using the syntax  
# for(i in start:end){  
# do something with i  
# }
```

```
> for(i in 1:10){  
+   print(i)  
+ }  
[1] 1  
[1] 2  
[1] 3  
[1] 4  
[1] 5  
[1] 6  
[1] 7  
[1] 8  
[1] 9  
[1] 10
```

LOOP through numbers using the FOR loop

```
# how to see how many missings you have in each column?  
# well, we want to sum through every column using a for loop  
# then print the variable name using names(movies[i])  
# then print the sum of is.na() for just that variable  
  
# for each column in the movies  
for(i in 1:ncol(movies)){  
  
  # print the following  
  print(  
  
    # first print "Variable: "  
    paste0("Variable: ",  
  
          # then print the variable name, then "NAs: "  
          names(movies)[i], " NAs: ",  
  
          # then print the sum of the number of missing values  
          # for that variable  
          sum(is.na(movies %>% select(i)))  
    )  
  )  
}
```

Functions in R

```
# -----  
# Creating functions  
# -----  
# we create a function in R by writing  
# function_name <- function(input1, input2,...){  
#   # function arguments  
# }  
  
print_names <- function(data_frame){  
  print(names(data_frame))  
}
```

```
> print_names(movies)  
[1] "color" "director_name" "num_critic_for_reviews"  
[4] "duration" "director_facebook_likes" "actor_3_facebook_likes"  
[7] "actor_2_name" "actor_1_facebook_likes" "gross"  
[10] "genres" "actor_1_name" "movie_title"  
[13] "num_voted_users" "cast_total_facebook_likes" "actor_3_name"  
[16] "facenumber_in_poster" "plot_keywords" "movie_imdb_link"  
[19] "num_user_for_reviews" "language" "country"  
[22] "content_rating" "budget" "title_year"  
[25] "actor_2_facebook_likes" "imdb_score" "aspect_ratio"  
[28] "movie_facebook_likes"
```

Build a function that prints number of missing values for each variable

```
# Let's take the code we wrote above and translate  
# it to a function called "num_missing".  
# We can then call the function and pass our movies dataframe  
# to it to export  
num_missing <- function(data_frame){  
  for(i in 1:ncol(movies)){  
    print(  
      paste0("Variable: ",  
            names(movies)[i], " NAs: ",  
            sum(is.na(movies %>% select(i)))  
    )  
  }  
}
```

```
> num_missing(movies)  
[1] "Variable: color NAs: 0"  
[1] "Variable: director_name NAs: 0"  
[1] "Variable: num_critic_for_reviews NAs: 50"  
[1] "Variable: duration NAs: 15"  
[1] "Variable: director_facebook_likes NAs: 104"  
[1] "Variable: actor_3_facebook_likes NAs: 23"  
[1] "Variable: actor_2_name NAs: 0"  
[1] "Variable: actor_1_facebook_likes NAs: 7"  
[1] "Variable: gross NAs: 884"  
[1] "Variable: genres NAs: 0"  
[1] "Variable: actor_1_name NAs: 0"  
[1] "Variable: movie_title NAs: 0"  
[1] "Variable: num_voted_users NAs: 0"  
[1] "Variable: cast_total_facebook_likes NAs: 0"  
[1] "Variable: actor_3_name NAs: 0"  
[1] "Variable: facenumber_in_poster NAs: 13"  
[1] "Variable: plot_keywords NAs: 0"  
[1] "Variable: movie_imdb_link NAs: 0"  
[1] "Variable: num_user_for_reviews NAs: 21"  
[1] "Variable: language NAs: 0"  
[1] "Variable: country NAs: 0"  
[1] "Variable: content_rating NAs: 0"  
[1] "Variable: budget NAs: 492"  
[1] "Variable: title_year NAs: 108"  
[1] "Variable: actor_2_facebook_likes NAs: 13"  
[1] "Variable: imdb_score NAs: 0"  
[1] "Variable: aspect_ratio NAs: 329"  
[1] "Variable: movie_facebook_likes NAs: 0"
```

MUTATE to Transform variables in your dataset

```
# -----  
# MUTATE to Transform variables in your dataset  
# -----  
  
# adding new variables using mutate()  
# note %<>% == DF <- DF %>%  
# let's create new variables budgetM and grossM that  
# are budget and gross in units of millions  
movies %<>% mutate(budgetM = budget/1000000,  
                  grossM = gross/1000000,  
                  profitM = grossM - budgetM)  
  
movies %>% glimpse()  
  
# so it looks like there's some outliers  
# The most expensive movie ever made was Pirates of  
# the Caribbean: On Stranger Tides  
# which cost $387.8m. Any movies with a budget higher  
# than this must be a data anomaly  
  
# Let's use the filter command to remove these  
movies_clean <- movies %>% filter(budgetM < 400)
```

Find Duplicate Rows with duplicated()

```
# -----  
# Find Duplicate Rows with duplicated()  
# and find_duplicates() (must install hablar package)  
# -----  
# number of duplicated rows  
movies %>% duplicated() %>% sum()  
  
# view duplicated rows  
# install.packages(hablar)  
movies %>% hablar::find_duplicates()
```

Output final clean version of dataset

```
# -----  
# Output final clean version of dataset  
# -----  
# remove duplicate rows, create new budget and gross variables,  
# rename director and title  
# remove budgets greater than 400M,  
# order title, year, budget, director and gross first, then store in new file  
movies_clean <-  
  movies %>%  
  distinct() %>%  
  mutate(budgetM = budget/1000000,  
         grossM = gross/1000000,  
         profitM = grossM - budgetM) %>%  
  rename(director = director_name,  
         title = movie_title,  
         year = title_year) %>%  
  relocate(title, year, country, director, budgetM, grossM, imdb_score) %>%  
  filter(budgetM < 400)  
  
movies_clean %>% glimpse()
```

- Generally we do pre-processing on our dataset starting from a raw file.
- After these transformations we save a “clean” version of the dataset that is used for analysis

Create summary statistics by GROUP using group_by()

```
# -----  
# Create summary statistics by GROUP using group_by()  
# -----  
# group summaries using summarise and group_by  
director_avg <-  
  movies_clean %>%  
    # group_by() is used to indicate the grouping variable  
    group_by(director) %>%  
  
    # summarize creates a new variable based on this group  
    # here we create averages by director using the 'mean'  
    # function |  
    summarize(gross_avg_director = mean(grossM, na.rm = TRUE))  
  
# view results  
director_avg %>% arrange(-gross_avg_director) %>% print()
```


Create averages, count and standard deviation by groups

```
# -----  
# Create grouped variables using the Summarize function  
# n() creates counts by  
# sd() creates standard deviations  
# -----  
# let's create budget by director, gross by director, profit by director,  
# number films by director  
director_df <-  
  movies_clean %>%  
  group_by(director) %>%  
  summarize(  
  
    # create average budget by director  
    budget_avg_director = mean(budgetM, na.rm = TRUE),  
    # create average gross by director  
    gross_avg_director = mean(grossM, na.rm = TRUE),  
    # create average movie profit by director  
    profit_avg_director = mean(profitM, na.rm = TRUE),  
    # create variable that lists number of films  
    # by director  
    num_films = n(),  
    # create a standard deviation of profit  
    # by director  
    profit_sd_director = sd(profitM, na.rm = TRUE)  
  
  )
```

Lab Exercises

1. Print a dataframe with the film director name, and number of films for the 10 directors with the most films in the dataset
2. What movie genres have the highest average profit? (hint, must use a new `group_by()` command)
3. Which countries have the most films in the top 5000 IMDB database? (hint, must use a new `group_by()` command)
4. How many missing values are there for the `profit_avg_director`?
5. Why do some directors have “NA” for `profit_avg_director`?

Class 3: Outline

1. Qs from last week?

2. Basic Data Analysis

- Missing values
- Loops
- mutate to transform variables
- Remove duplicates with distinct
- Outputting “clean” data file”

3. Data Analysis by Groups

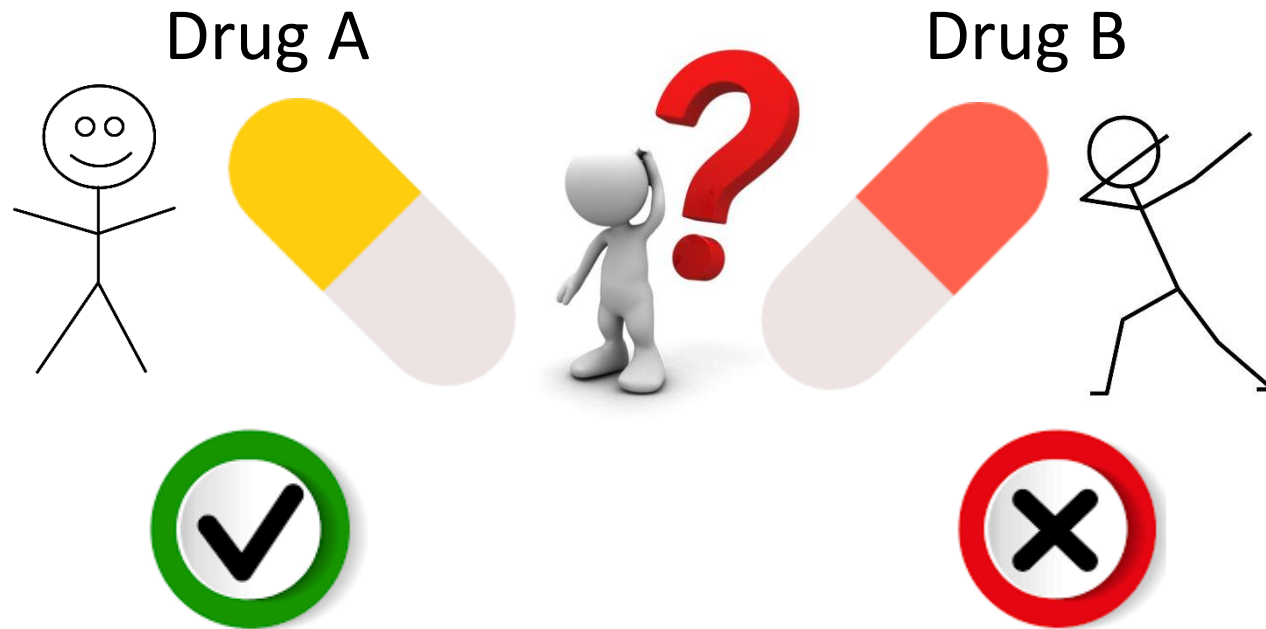
- group_by() function
- summarize() to create group variables

4. Data Analysis Lab Class 3

5. P-values

6. Introductory Machine Learning Concepts

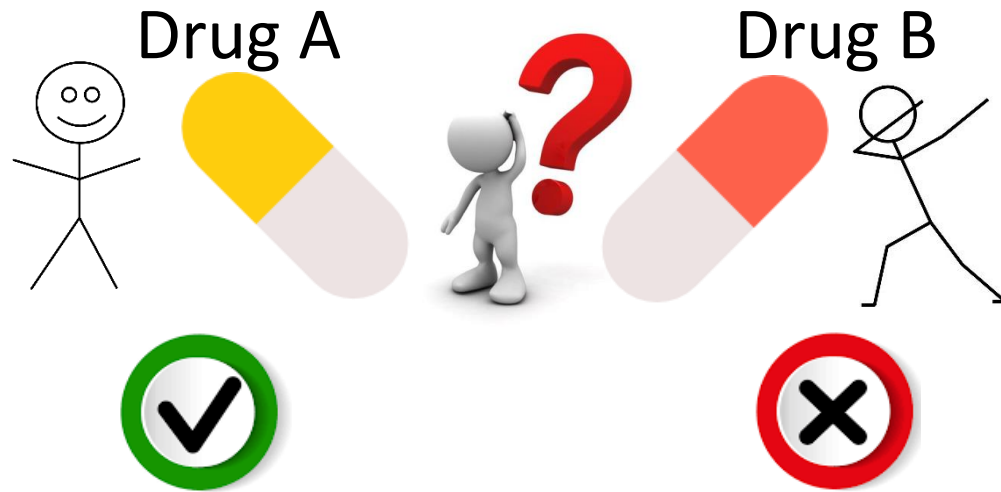
What Are P-Values?



- Suppose we want to know the effectiveness of Drug A vs Drug B
- We can give Drug A to 1 person, and give drug B to 1 person.
- Suppose person A gets better, and person B does not.
- Can we conclude drug A is better than drug B?

Adapted from the excellent StatQuest: <https://youtu.be/vemZtEM63GY>

What Are P-Values?



- **No! Perhaps...**
 - Person B didn't follow the instructions
 - Person B has pre-existing conditions
 - Person A is healthier
- Only by repeating the experiment many times can we learn whether Drug A > Drug B
- Ideally we express our confidence as a quantitative number, how likely it is that we find Drug A > Drug B only due to chance?

Repeating Experiment With More Observations

Drug A



Cured!	Not Cured!
73	125

37%
Cured!

Drug B

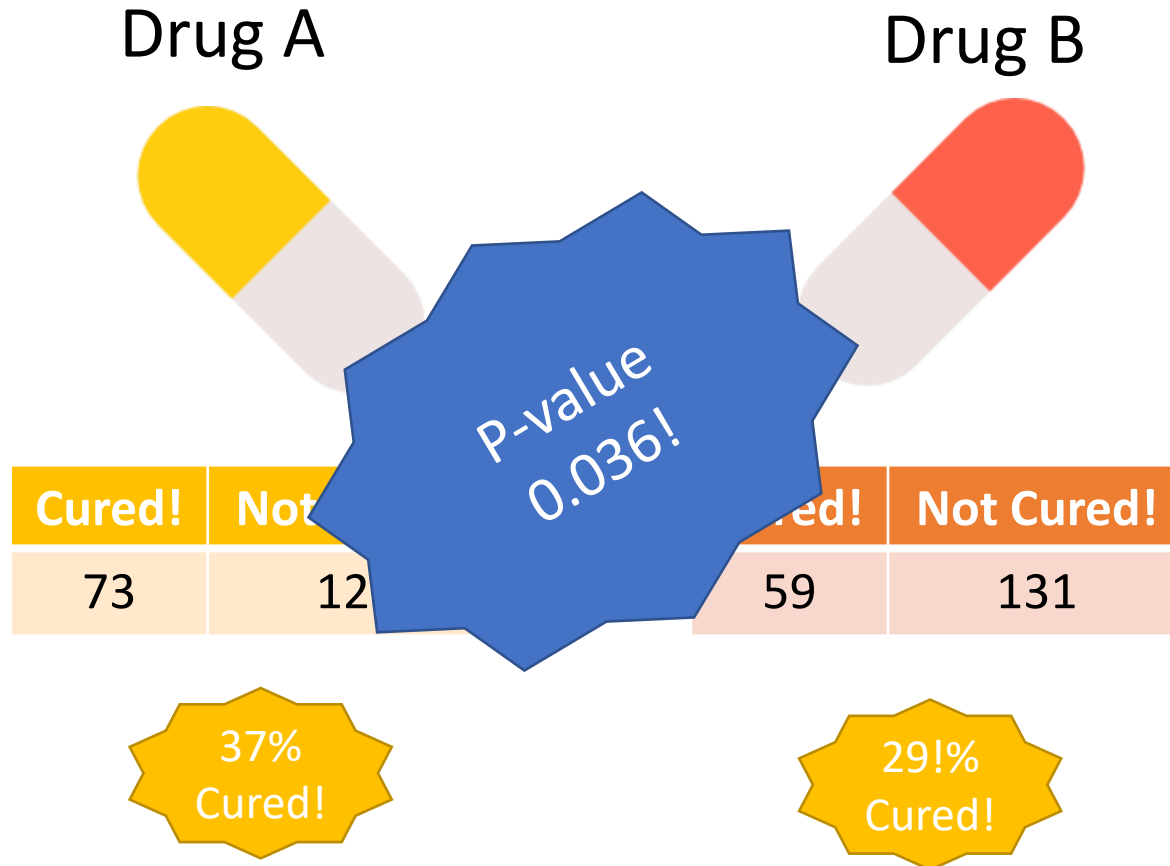


Cured!	Not Cured!
59	131

29!%
Cured!

- Suppose we repeat the experiment with ~400 patients
- We find that drug A has a cure rate of 37%, drug B a cure rate of 29%.
- But: No study is perfect. There are always random things that could influence drug A vs B
- **The p-value is a number between 0 and 1 that tells us how confident we should be between one hypothesis (H_0 or null) and another (H_1 , alternative)**

Repeating Experiment With More Observations



- P value close to 0:
 - Result not likely due to chance
- P value close to 1:
 - More likely result is due to chance
- What is “enough evidence”?
 - Alpha = critical value
 - Commonly use ~ 0.05 , 0.01 , or 0.001
 - i.e. 5%, 1% or 0.1% chance difference due to randomness and there’s no difference

What do p-values measure?

P-value tells us the likelihood that, if the null hypothesis were true, we would receive a result as extreme as the one seen

Small p-values: unlikely that we would receive a result as extreme as the one seen if the null is true

P-value does not measure

- Size of effect
- Importance of a result
- Probability the alternative hypothesis is true

If You Find This Confusing You Are In Good Company

THE AMERICAN STATISTICIAN
2016, VOL. 70, NO. 2, 129–133
<http://dx.doi.org/10.1080/00031305.2016.1154108>



EDITORIAL

The ASA's Statement on p -Values: Context, Process, and Purpose

In February 2014, George Cobb, Professor Emeritus of Mathematics and Statistics at Mount Holyoke College, posed these questions to an ASA discussion forum:

Q: Why do so many colleges and grad schools teach $p = 0.05$?

A: Because that's still what the scientific community and journal editors use.

Q: Why do so many people still use $p = 0.05$?

A: Because that's what they were taught in college or grad school.

Cobb's concern was a long-worrisome circularity in the sociology of science based on the use of bright lines such as $p < 0.05$: "We teach it because it's what we do; we do it because it's what we teach." This concern was brought to the attention of the ASA Board.

The ASA Board was also stimulated by highly visible discussions over the last few years. For example, ScienceNews (Siegfried 2010) wrote: "It's science's dirtiest secret: The 'scientific method' of testing hypotheses by statistical analysis stands on a flimsy foundation." A November 2013, article in Phys.org Science News Wire (2013) cited "numerous deep flaws" in null hypothesis significance testing. A ScienceNews article (Siegfried 2014) on February 7, 2014, said "statistical techniques for testing hypotheses ... have more flaws than Facebook's privacy policies." A week later, statistician and "Simply Statistics" blogger Jeff Leek responded. "The problem is not that people use P -values poorly," Leek wrote, "it is that the vast majority of data analysis is not performed by people properly trained to perform data analysis" (Leek 2014). That same week, statistician and science writer Regina Nuzzo published an article in *Nature* entitled "Scientific Method: Statistical Errors" (Nuzzo 2014). That article is now one of the most highly viewed *Nature* articles, as reported by altmetric.com (<http://www.altmetric.com/details/2115792#score>).

Of course, it was not simply a matter of responding to some

2014) and a statement on risk-limiting post-election audits (American Statistical Association 2010). However, these were truly policy-related statements. The VAM statement addressed a key educational policy issue, acknowledging the complexity of the issues involved, citing limitations of VAMs as effective performance models, and urging that they be developed and interpreted with the involvement of statisticians. The statement on election auditing was also in response to a major but specific policy issue (close elections in 2008), and said that statistically based election audits should become a routine part of election processes.

By contrast, the Board envisioned that the ASA statement on p -values and statistical significance would shed light on an aspect of our field that is too often misunderstood and misused in the broader research community, and, in the process, provides the community a service. The intended audience would be researchers, practitioners, and science writers who are not primarily statisticians. Thus, this statement would be quite different from anything previously attempted.

The Board tasked Wasserstein with assembling a group of experts representing a wide variety of points of view. On behalf of the Board, he reached out to more than two dozen such people, all of whom said they would be happy to be involved. Several expressed doubt about whether agreement could be reached, but those who did said, in effect, that if there was going to be a discussion, they wanted to be involved.

Over the course of many months, group members discussed what format the statement should take, tried to more concretely visualize the audience for the statement, and began to find points of agreement. That turned out to be relatively easy to do, but it was just as easy to find points of intense disagreement.

Q: Why do so many colleges and grad schools teach $p = 0.05$?

A: Because that's still what the scientific community and journal editors use.

Q: Why do so many people still use $p = 0.05$?

A: Because that's what they were taught in college or grad school.

Cobb's concern was a long-worrisome circularity in the sociology of science based on the use of bright lines such as $p < 0.05$: "We teach it because it's what we do; we do it because it's what we teach." This concern was brought to the attention of the ASA Board.

P-Values Commonly Misused or Misunderstood

Statistical tests, P values, confidence intervals, and power: a guide to misinterpretations

Sander Greenland¹ · Stephen J. Senn² · Kenneth J. Rothman³ · John B. Carlin⁴ · Charles Poole⁵ · Steven N. Goodman⁶ · Douglas G. Altman⁷

Received: 9 April 2016 / Accepted: 9 April 2016 / Published online: 21 May 2016
© The Author(s) 2016. This article is published with open access at Springerlink.com

Abstract Misinterpretation and abuse of statistical tests, confidence intervals, and statistical power have been decried for decades, yet remain rampant. A key problem is that there are no interpretations of these concepts that are at once simple, intuitive, correct, and foolproof. Instead, correct use and interpretation of these statistics requires an attention to detail which seems to tax the patience of working scientists. This high cognitive demand has led to an epidemic of shortcut definitions and interpretations that are simply wrong, sometimes disastrously so—and yet these misinterpretations dominate much of the scientific

literature. In light of this problem, we provide a discussion of basic statistics that are more accurate and critical than typically found in traditional expositions. Our goal is to provide a resource for teachers, researchers, and consumers of statistics who have limited but who wish to avoid and spot misinterpretations. We emphasize how violation of often unspoken assumptions (such as selecting analyses for presentation on the P values they produce) can lead to misleading conclusions even if the declared test hypothesis is correct. We then provide an explanatory list of 25 misinterpretations of P values, confidence intervals, and power, along with guidelines for improving statistical interpretation and reporting.

Editor's note This article has been published online as supplementary material with an article of Wasserstein RL, Lazar NA. The ASA's statement on p-values: context, process and purpose. The American Statistician 2016.

Albert Hofman, Editor-in-Chief EJE.

Common misinterpretations of single P values

1. The P value is the probability that the test hypothesis is true; for example, if a test of the null hypothesis gave $P = 0.01$, the null hypothesis has only a 1 % chance of being true; if instead it gave $P = 0.40$, the null hypothesis has a 40 % chance of being true. No! The P value *assumes* the test hypothesis is true—it is *not* a hypothesis probability and may be far from any reasonable probability for the test hypothesis. The P value simply indicates the degree to which the data conform to the pattern predicted by the test hypothesis and all the other assumptions used in the test (the underlying statistical model). Thus $P = 0.01$ would indicate that the data are not very close to what the statistical model (including the test hypothesis) predicted they should be, while $P = 0.40$ would indicate that the data are much closer to the model prediction, allowing for chance variation.

Class 3: Outline

1. Qs from last week?

2. Basic Data Analysis

- Missing values
- Loops
- mutate to transform variables
- Remove duplicates with distinct
- Outputting “clean” data file”

3. Data Analysis by Groups

- group_by() function
- summarize() to create group variables

4. Data Analysis Lab Class 3

5. P-values

6. Introductory Machine Learning Concepts

Supervised vs Unsupervised Learning

Supervised Learning:

- For every x_i we observe some y_i
- Ex: random forests to predict loan default (y_i) based on applicant characteristics (x_i)

Supervised Learning



Unsupervised Learning



Unsupervised Learning:

- We only observe x_i
- Ex: clustering loan applicants based on characteristics (x_i)

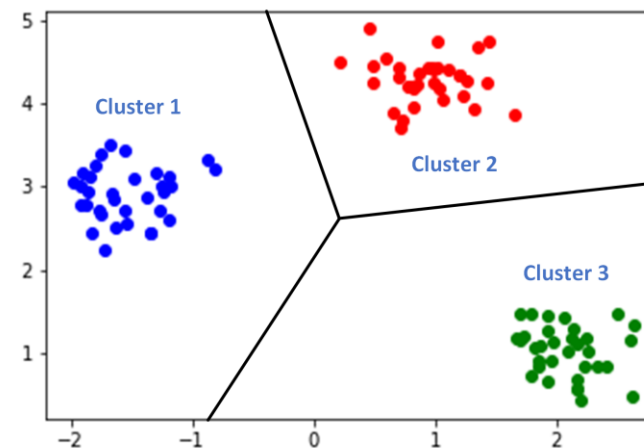
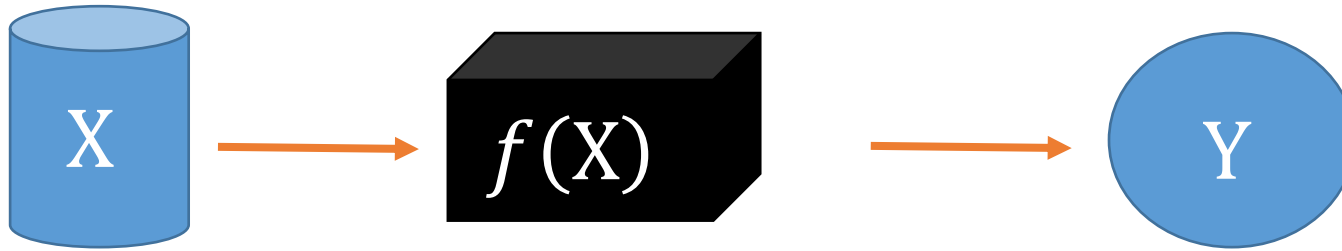


Fig.1. An Example Of Data Clustering

Supervised learning: learning $f(X)$ our predicted out given inputs

$$Y = f(X) + \epsilon$$



ϵ = “epsilon” (unexplained portion)

“Estimating” $\hat{f}(X)$

- $Y = f(X) + \epsilon$ is the true value
- We can only use data to “guess” at $f(X)$
- We call this guess $\hat{f}(X)$

How do we know when we’ve selected a “good” $\hat{f}(X)$?

- We reserve a portion of our data into a “test” set, estimate a model on the other part, and see how our model performs on this test set

Testing Training Data Subsets

Training set: (observation-wise) subset of data used to develop models



Testing/Training Split

Training set: (observation-wise) subset of data used to develop models

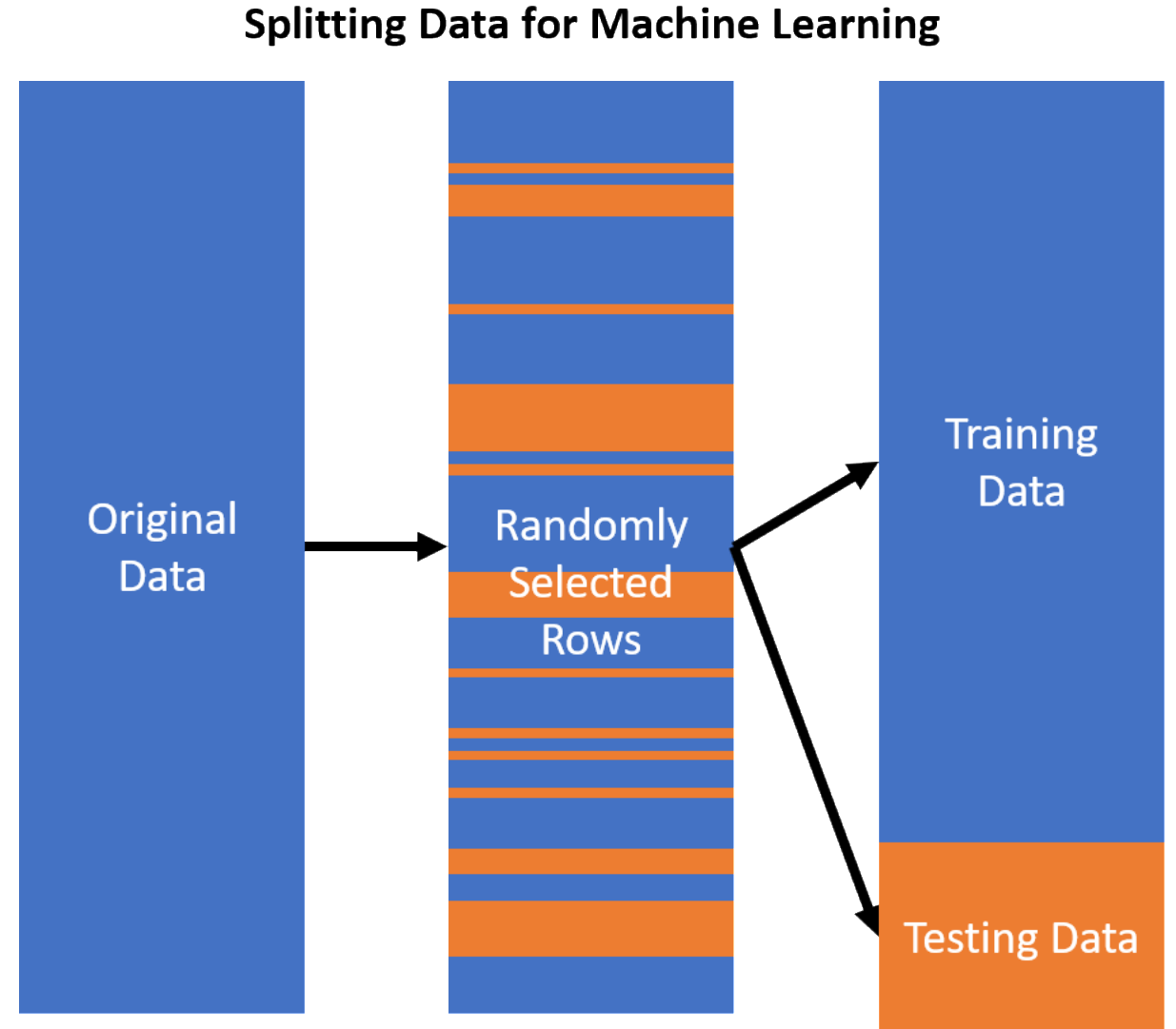
Test set: subset of data used during intermediate stages to “tune” model parameters

Rule of thumb 75% training 25% test -ish



Randomly Selecting Rows for Test or Training Sets

- Observations are randomly selected into either testing or training splits of the data



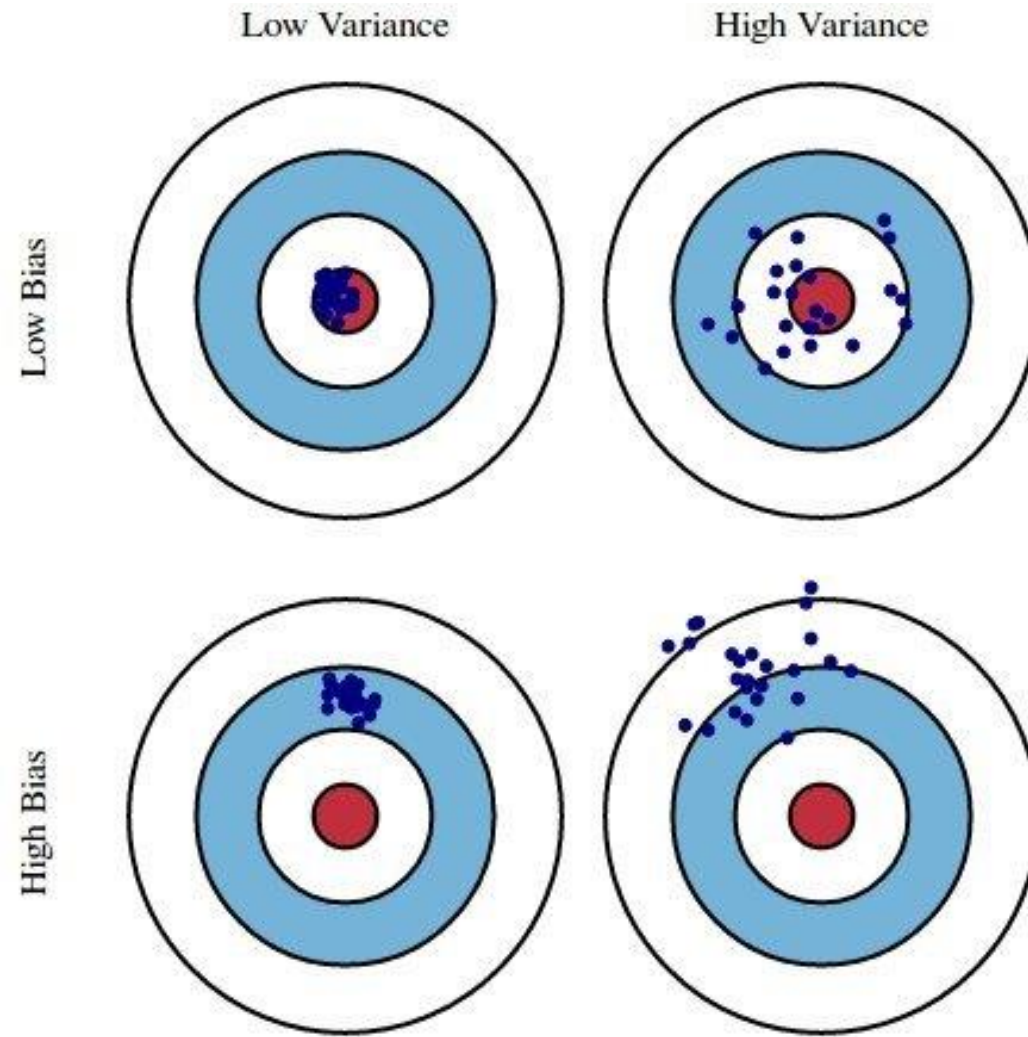
Bias and Variance

Bias: Tendency of an in-sample statistic to over or under estimate the statistic in the *population*

Variance: Tendency to noisily estimate a statistic.

E.g., sensitivity to small fluctuations in the training dataset.

Bias-Variance Tradeoff



Class 3 Summary

- **Missing values** (NAs) indicate we don't know the value of a variable for that observation
 - Will need to make assumptions on how to treat these that can influence our results!
- **Functions** create “more readable” code.
- “**Clean**” version of datasets have been processed and are ready for analysis
- Use **group_by()** and **summarize()** to create statistics by groups (averages, standard deviations)
- P-values measure the discrepancy of the fit of a model or “null hypothesis” to data.
 - P-value tells us the likelihood that, if the null

hypothesis were true, we would receive a result as extreme as the one seen

- **Supervised** models contain a y_i (target/outcome variable) for every x_i (descriptor variables)
- **Unsupervised** models contain only x_i
- **Training** data is the data we will use to estimate our model parameters
- **Testing** data is the data used to evaluate our model performance
- **Bias:** tendency of an in-sample statistic to over or underestimate the true value
- **Variance:** tendency to noisily estimate that statistic