

INTRO TO DATA SCIENCE

PT2 ORGANIZING DATA

How do we get data into R and how do we make it workable?

MIDTERM 1

R & R STUDIO

- Data science is statistics meets computer science, so much of the work we do is in R Studio w/ coding language R.
- Some basics about R...
 - we will work w/ Variables, numbers, vectors, tibbles/dataframes and more
 - we can use normal calculator functions in R too (! + -)
 - typically functions are called in this gen manner...
function(argument1 = __, arg2 = __, ...)

- In R Studio we have 4 places we work in...
 - Editor → where we write code
 - Console → see our code run
 - Environment → shows all obj. you created

	ENV
Editor	
Console	Output

R MARKDOWN

- R markdown is a file type in R studio that allows you to mix code and normal writing
- Use Knit to export to a .doc, .html, or .pdf.
- Code chunks are started w/ `{r}` and end w/ `}`
- We can include some helpful options in headers...

Stops...	Run Code	Show Code	Output	plots	mess.	warn.
eval=F	x		x	x	x	x
include=F		x	x	x	x	x
echo=F		x				
results="hide"			x			
fig.show="hide"				x		
message=F					x	
warning=F						x

• Can also name chunks...
`{r Name, ...}`
• chunks named "setup" will automatically run before any other code.

CODE USED

- `seq(start, stop, length.out =)` makes seq sequence of ints
- `getwd()` finds working directory
- Knitr::kable(data, options) prints df as tables in knitted file

LIBRARIES USED

• Knitr

CODE USED

- `filter(data, var ==)`
- `is.na()` returns TRUE if value NA
- `arrange(data, var)`
- `select(data, column)`
- `rename(data, newvar = var)`
- `mutate(data, newvar = var)`
- `transmute()` mutate but only keeps new variables
- `summarize(data, newvar = func(var))`
- `group_by()` groups data based on a var. (good w/ sum, avg, group)
- `read_csv(file path/.csv, col.names = c("var", ...))`
- `read_delim()`
- `guess_encoding()`

- `parse_*()` ex: `parse_double()`, `parse_logical()`, ...
- `write_csv(data, filename.csv)`
- `as_tibble()` turns df to tibble
- `tibble()` creates new tibble
- `print()` prints df or tibble
- `as.data.frame()` df to df
- `gather()` this is to combine key = "old become" value = "new"
- `spread(data, key = var to split, value = var of value)`
- `separate(data, col to split, into = c("new name", sep = "separator"))`
- `unite(data, new col, col to come, sep = ",")`

DATA IMPORT

- we can use the `readr` library for a convenient way to import data files
- works best w/ CSV's but can also specify the delimiter
- but all CSV's only contain character type vectors so to change the type of a column we use a `parse_*()` function
- R also expresses characters in strings w/ underlying code and if that encoding can't express certain characters we get weird strings.
- Use `guess_encoding()` to try and guess
- can also use `readr` to save to csv

TIBBLES

- Tibbles → better data frames
- main differences...
 - tibble doesn't print whole dataframes
 - can extract from a tibble by name or position

- `complete()` turns implicit missing values into explicit missing
- `fill()` fill in NAs or w/ other values.

LIBRARIES USED

• dplyr

• tibbles

• readr

• tidy

DATA TRANSFORMATION

- we will use 5 main functions...
 - `filter()` → pick observations
 - `arrange()` → reorder rows
 - `select()` → pick variables by name
 - `mutate()` → new variables made w/ func. of existing variables
 - `summarize()` → collapse many values to a single summary
- Data transformation often use logical arguments and operators
- when transforming data we often use the pipe `%>%` which allows us to put multiple steps into one.
ex: `F(g(x)) %>% x %>% g() %>% F()`
- when using `summarize` we typically also have to use the `group_by()` function, most of these funcs have helper functions that allow you to do more complicated things w/ them.

TIDY DATA

- you can often represent the same data in many different ways, but some ways are better than others.
- Tidy data (for our purposes) means 1 each variable has its own column, 2 each observation has its own row, 3 each value has its own cell.
- Some untidy examples...
 - problem: column names are actually values of a variable
solution: `gather()`
 - problem: one observation is in multiple rows
solution: `spread()`
 - problem: two pieces of data in each cell of a column
solution: `separate()`
 - problem: info from two sep. columns needs to be tag. in one column
solution: `unite()`
- Missing values are either explicitly missing (NAs) or implicitly missing (we should have an entry for this month but we don't). We have to decide how to handle both types. (get rid of? comment on?)

PT3 EXPLORING DATA

DATA VISUALIZATION

- We start by making visualizations for data exploration (later we'll make some for communication)
- we use `ggplot2` for graphing.
- the `ggplot` graphics take arguments in the "aes" category for splitting up variables.
- you can get many variables in one graph by mapping variables to the color of your graph using `facet_wrap()`
- `geom` objects are the functions we use to tell `ggplot2` what kind of graph to make.
- we can add text, labels, legends, and themes to our graphs too!
- general graphing layout...
`ggplot(data) + geom_* (aes(variables)) + theme_* (theme and text)`
- Some common graphs → line graph, scatter plot, box plot, tile plot, bar plot, histogram, ...

CODE USED

- `ifelse(logical, do.it, otherwise)` tests if something is T
- `ggplot()` start a plot
- `geom_*()` adds type of graph to your plot...
 - `point()`
 - `smooth()`
 - `bar()`
 - `boxplot()`
 - `tile()`
 - `histogram()`
- `aes()` the argument in `ggplot` where you specify variables
- 2 functions get info on a function.
- `facet_wrap(nvar)` create multiple graphs grouped by a specific variable
- `facet_grid()`
- `position = "` control data position in a grid.
- `coord_*()` controls coord graph.
- `labs()` add labels to graph and legend.
- `geom_text()` label individual points
- `theme_*()` set a theme for your graph (color, text, positions)

LIBRARIES USED

• ggplot2

EXPLORATORY DATA ANALYSIS

- EDA → Exploratory Data Analysis.
- generate questions → search for answers → refine questions → visualize, transform, and model
- Questions mainly guide EDA
 - Two types of questions → what kind of relationship exists within my var? what type of relationship occurs between my vars?
 - Relationship
 - how variables change between obsv.
 - visualize w/ histograms, bar charts
 - look for unexpected values (outliers)
 - look for subgroups
 - Comparison
 - tendency for 2 variables to vary together in a related way
 - cat, 3 cont. → boxplot
 - 2 cat. → tile plot
 - 2 cont. → scatter plots!
- Once we see patterns we get to ask more questions about what causes these patterns!

Good Luck!