

STREAMLINING WITH R

Meghan Hall
NEAIR
July 12, 2022

HOUSEKEEPING

- Intro 🖐️
- Workshop materials ⬇️
- Break 🕒
- By the end of today ✓
- Today's plan 📋

TODAY'S PLAN

1. What is R? How can it ease the burden of repeated reporting?
2. Basic functions for manipulating data
3. Using R effectively
4. More data manipulation
5. Visualizing data
6. A peek at advanced topics

WHAT IS R?

12

WHAT IS R?

12

R is an open-source (**free**!) scripting language for working with data

THE BENEFITS OF R

12

My personal Excel nightmare

The magic of R is that it's **reproducible** (by someone else or by yourself in six months)

Keeps data separate from code (data preparation steps)

You need the R language

And also the software

NAVIGATING RSTUDIO

12

The screenshot shows the RStudio IDE with the following components and annotations:

- Console:** Displays the R version (4.2.0) and startup information. A red annotation "code can go here" points to the prompt line.
- Environment:** Shows the Global Environment. A red annotation "imported data shows up here" points to the empty environment area.
- Files:** Shows the project files in the "sample" directory. A red annotation "project files are here" points to the file list.

Name	Size	Modified
..		
sample.Rproj	258 B	Jul 10, 2022, 4:45 PM
NEAIR_code.R	2.3 KB	Jul 10, 2022, 1:50 PM
courses.csv	3.8 KB	Jul 7, 2022, 11:59 AM
faculty.csv	16.6 KB	Jul 7, 2022, 11:27 AM

NAVIGATING RSTUDIO

12

The screenshot shows the RStudio interface with the following components:

- Source Editor:** Contains R code for loading and filtering data. The code is as follows:

```
1 library(tidyverse)
2
3 faculty <- read_csv("faculty.csv")
4 courses <- read_csv("courses.csv")
5
6 # filter----
7
8 # the `==` operator tests for equality
9
10 faculty %>%
11   filter(dept1 == "Sociology")
12
13 # the `|` operator signifies "or"
14
15 faculty %>%
16   filter(dept1 == "Sociology" |
17          dept1 == "Physics")
18
19 # the `%in%` operator allows for multiple options in a list
20
```
- Environment:** Shows the Global Environment with two data objects:

Object	Observations	Variables
courses	104 obs.	6 variables
faculty	392 obs.	5 variables
- Files:** Shows the project files in the Desktop/sample directory:

Name	Size	Modified
..		
sample.Rproj	258 B	Jul 10, 2022, 4:45 PM
NEAIR_code.R	373 B	Jul 10, 2022, 4:49 PM
courses.csv	3.8 KB	Jul 7, 2022, 11:59 AM
faculty.csv	16.6 KB	Jul 7, 2022, 11:27 AM
- Console:** Shows the output of the R code, including the column specification for the 'courses' data frame:

```
R 4.2.0 ~ Desktop/sample/
chr (4): year, rank, dept1, dept2
dbl (1): id

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
> courses <- read_csv("courses.csv")
Rows: 104 Columns: 6
— Column specification —
Delimiter: ","
chr (2): dept, level
dbl (4): semester, course_id, faculty_id, enrollment

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
>
```

Annotations in red text:

- "code can also go here" points to the Source Editor.
- "imported data shows up here" points to the Environment pane.
- "project files are here" points to the Files pane.

USING R

12

You use R via *packages*

...which contain *functions*

...which are just *verbs*



faculty

year	id	rank	dept1	dept2
2021-22	1005	Lecturer	Chemistry	
2021-22	1022	Professor	Physics	Engineering
2021-22	1059	Professor	Physics	
2021-22	1079	Lecturer	Music	
2021-22	1086	Assistant Professor	Music	
2021-22	1095	Adjunct Instructor	Sociology	

courses

semester	course_id	faculty_id	dept	enrollment	level
20212202	10605	1772	Physics	7	UG
20212202	10605	1772	Physics	32	GR
20212202	11426	1820	Political Science	8	UG
20212202	12048	1914	English	24	UG
20212202	13269	1095	Sociology	48	UG
20212202	13517	1086	Music	17	UG

BASIC DATA MANIPULATION

12

USEFUL OPERATORS

12

`<-` “save as” `opt + -`

`%>%` “and then” `Cmd + shift + m`

COMMON FUNCTIONS

12

`filter` keeps or discards rows (aka observations)

`select` keeps or discards columns (aka variables)

`arrange` sorts data set by certain variable(s)

`count` tallies data set by certain variable(s)

`mutate` creates new variables

`group_by/summarize` aggregates data (**pivot tables!**)

`str_*` functions work easily with text

SYNTAX OF A FUNCTION

```
function(data, argument(s))
```

is the same as

```
data %>%  
  function(argument(s))
```


FILTER

12

`filter` keeps or discards rows (aka observations)

the `==` operator tests for equality

```
1 faculty %>%  
2   filter(dept1 == "Sociology")
```

year	id	rank	dept1	dept2
2021-22	1095	Adjunct Instructor	Sociology	
2021-22	1118	Assistant Professor	Sociology	
2021-22	1161	Assistant Professor	Sociology	
2021-22	1191	Professor	Sociology	
2021-22	1216	Associate Professor	Sociology	American Studies
2021-22	1273	Assistant Professor	Sociology	

FILTER

12

the `|` operator signifies “or”

```
1 faculty %>%  
2   filter(dept1 == "Sociology" |  
3         dept1 == "Physics")
```

year	id	rank	dept1	dept2
2021-22	1022	Professor	Physics	Engineering
2021-22	1059	Professor	Physics	
2021-22	1095	Adjunct Instructor	Sociology	
2021-22	1118	Assistant Professor	Sociology	
2021-22	1161	Assistant Professor	Sociology	
2021-22	1191	Professor	Sociology	

the `%in%` operator allows for multiple options in a list

```
1 faculty %>%
2   filter(dept1 %in% c("Sociology",
3                     "Physics",
4                     "Music"))
```

year	id	rank	dept1	dept2
2021-22	1022	Professor	Physics	Engineering
2021-22	1059	Professor	Physics	
2021-22	1079	Lecturer	Music	
2021-22	1086	Assistant Professor	Music	
2021-22	1095	Adjunct Instructor	Sociology	
2021-22	1118	Assistant Professor	Sociology	

the **&** operator combines conditions

```
1 faculty %>%  
2   filter(dept1 %in% c("Sociology",  
3                     "Physics",  
4                     "Music") &  
5           rank == "Professor")
```

year	id	rank	dept1	dept2
2021-22	1022	Professor	Physics	Engineering
2021-22	1059	Professor	Physics	
2021-22	1191	Professor	Sociology	
2021-22	1201	Professor	Physics	
2021-22	1209	Professor	Music	
2021-22	1421	Professor	Physics	Engineering

SELECT

12

select keeps or discards columns (aka variables)

```
1 faculty %>%  
2   select(id, dept1, rank)
```

id	dept1	rank
1005	Chemistry	Lecturer
1022	Physics	Professor
1059	Physics	Professor
1079	Music	Lecturer
1086	Music	Assistant Professor
1095	Sociology	Adjunct Instructor

SELECT

12

can drop columns with **-column**

```
1 faculty %>%  
2   select(-dept2)
```

year	id	rank	dept1
2021-22	1005	Lecturer	Chemistry
2021-22	1022	Professor	Physics
2021-22	1059	Professor	Physics
2021-22	1079	Lecturer	Music
2021-22	1086	Assistant Professor	Music
2021-22	1095	Adjunct Instructor	Sociology

the pipe `%>%` chains multiple functions together

```
1 faculty %>%  
2   select(id, dept1, rank) %>%  
3   filter(rank == "Professor")
```

id	dept1	rank
1022	Physics	Professor
1059	Physics	Professor
1191	Sociology	Professor
1201	Physics	Professor
1209	Music	Professor
1407	English	Professor

ARRANGE

12

arrange sorts data set by certain variable(s)

use **desc()** to get descending order

```
1 courses %>%  
2   arrange(desc(enrollment))
```

semester	course_id	faculty_id	dept	enrollment
20212201	10511	1005	Chemistry	50
20212201	15934	1421	Physics	50
20192002	13850	1105	Chemistry	50
20181901	17773	1942	Music	50
20212202	13269	1095	Sociology	48
20202101	16202	1816	Political Science	48

ARRANGE

12

can sort by multiple variables

```
1 courses %>%  
2   arrange(dept, desc(enrollment))
```

semester	course_id	faculty_id	dept	enrollment
20212201	10511	1005	Chemistry	50
20192002	13850	1105	Chemistry	50
20202102	13850	1258	Chemistry	39
20202102	16606	1393	Chemistry	38
20202101	16540	1784	Chemistry	38
20181901	10511	1829	Chemistry	36

COUNT

12

count tallies data set by certain variable(s) (very useful for familiarizing yourself with data)

```
1 courses %>%  
2   count(dept)
```

dept	n
Chemistry	16
English	18
Music	17
Physics	19
Political Science	17
Sociology	17

COUNT

12

can use `sort = TRUE` to order results

```
1 courses %>%  
2   count(dept, level, sort = TRUE)
```

dept	level	n
Chemistry	UG	16
English	UG	16
Music	UG	16
Physics	UG	16
Political Science	UG	16
Sociology	UG	16
Physics	GR	3
English	GR	2
Music	GR	1
Political Science	GR	1
Sociology	GR	1

mutate creates new variables (with a single =)

```
1 faculty %>%
2   mutate(new = "hello!")
```

year	id	rank	dept1	dept2	new
2021-22	1005	Lecturer	Chemistry		hello!
2021-22	1022	Professor	Physics	Engineering	hello!
2021-22	1059	Professor	Physics		hello!
2021-22	1079	Lecturer	Music		hello!
2021-22	1086	Assistant Professor	Music		hello!
2021-22	1095	Adjunct Instructor	Sociology		hello!

much more useful with a conditional such as `ifelse()`,
which has three arguments:
condition, value if true, value if false

```
1 faculty %>%  
2   mutate(prof = ifelse(rank == "Professor",  
3                       1, 0)) %>%  
4   select(rank, prof)
```

rank	prof
Lecturer	0
Professor	1
Professor	1
Lecturer	0
Assistant Professor	0
Adjunct Instructor	0

MUTATE

12

the **!** operator means *not*
is.na() identifies null values

```
1 faculty %>%  
2   mutate(joint = ifelse(!is.na(dept2),  
3                         "joint", NA)) %>%  
4   select(dept1, dept2, joint)
```

dept1	dept2	joint
Chemistry		
Physics	Engineering	joint
Physics		
Music		
Music		
Sociology		

with multiple conditions, `case_when()` is much easier!

```
1 faculty %>%
2   mutate(division = case_when(dept1 %in% c("Sociology","Political Science")
3     "Social Sciences",
4     dept1 %in% c("Music","English") ~
5       "Humanities",
6     dept1 %in% c("Chemistry","Physics") ~
7       "Sciences")) %>%
8   select(dept1, division)
```

dept1	division
Chemistry	Sciences
Physics	Sciences
Physics	Sciences
Music	Humanities
Music	Humanities
Sociology	Social Sciences

GROUP BY / SUMMARIZE

12

`group_by/summarize` aggregates data (**pivot tables!**)

`group_by()` identifies the grouping variable(s) and
`summarize()` specifies the aggregation

```
1 courses %>%  
2   group_by(dept, semester) %>%  
3   summarize(enr = sum(enrollment))
```

dept	semester	enr
Chemistry	20181901	59
Chemistry	20181902	44
Chemistry	20192001	47
Chemistry	20192002	68
Chemistry	20202101	69
Chemistry	20202102	77

GROUP BY / SUMMARIZE

12

useful arguments within `summarize`:
`mean`, `median`, `sd`, `min`, `max`, `n`

```
1 courses %>%  
2   group_by(dept, semester) %>%  
3   summarize(enr = sum(enrollment),  
4             count = n_distinct(course_id))
```

dept	semester	enr	courses
Chemistry	20181901	59	2
Chemistry	20181902	44	2
Chemistry	20192001	47	2
Chemistry	20192002	68	2
Chemistry	20202101	69	2
Chemistry	20202102	77	2

USING R EFFECTIVELY

12

WORKING IN RSTUDIO

12

The screenshot displays the RStudio interface with the following components:

- Source Editor:** Contains R code for loading and filtering data. The code is as follows:

```
1 library(tidyverse)
2
3 faculty <- read_csv("faculty.csv")
4 courses <- read_csv("courses.csv")
5
6 # filter----
7
8 # the `==` operator tests for equality
9
10 faculty %>%
11   filter(dept1 == "Sociology")
12
13 # the `|` operator signifies "or"
14
15 faculty %>%
16   filter(dept1 == "Sociology" |
17          dept1 == "Physics")
18
19 # the `%in%` operator allows for multiple options in a list
20
```
- Environment:** Shows the loaded data frames:

Data	Observations	Variables
courses	104 obs.	6 variables
faculty	392 obs.	5 variables
- Files:** Shows the project files in the Desktop/sample directory:

Name	Size	Modified
..		
sample.Rproj	258 B	Jul 10, 2022, 4:45 PM
NEAIR_code.R	373 B	Jul 10, 2022, 4:49 PM
courses.csv	3.8 KB	Jul 7, 2022, 11:59 AM
faculty.csv	16.6 KB	Jul 7, 2022, 11:27 AM
- Console:** Shows the output of the R code, including column specifications and data types for the loaded datasets.

Annotations on the image:

- "code can also go here" points to the Source Editor.
- "imported data shows up here" points to the Environment pane.
- "project files are here" points to the Files pane.

Typing in the console

- think of it like a post-it: useful for quick notes but disposable
- actions are saved but code is not
- one chunk of code is run at a time (**Return**)

Typing in a code file

- script files have a **.R** extension
- code is saved and sections of any size can be run (**Cmd + Return**)
- do ~95% of your typing in a code file instead of the console!

WORKING WITH PACKAGES

12

packages need to be installed on each computer you use

```
1 # only need to do this once (per computer)
2 install.packages("tidyverse")
```

packages need to be loaded/attached with `library()` at the beginning of every session

```
1 # always put the necessary packages at the top of a code file
2 library(tidyverse)
```

can access help files by typing `??tidyverse` or `??mutate` in the console

ORGANIZING WITH PROJECTS

12

highly recommend using *projects* to stay organized

keeps code files and data files together, allowing for easier file path navigation and better reproducible work habits

File -> New Project

more guidance: [here](#) and [here](#)

ORGANIZING WITH PROJECTS

12

The screenshot displays the RStudio interface with the following components:

- Source Editor:** Contains R code for loading and filtering data. The code includes comments explaining the use of `filter()`, `%>%`, and `%in%` operators.
- Environment:** Shows the loaded data frames: `courses` (104 obs. of 6 variables) and `faculty` (392 obs. of 5 variables).
- Files:** Lists project files in the `sample` directory: `sample.Rproj`, `NEAIR_code.R`, `courses.csv`, and `faculty.csv`.

Annotations in red text highlight key features:

- code can also go here** points to the Source Editor.
- imported data shows up here** points to the Environment pane.
- project files are here** points to the Files pane.

```
1 library(tidyverse)
2
3 faculty <- read_csv("faculty.csv")
4 courses <- read_csv("courses.csv")
5
6 # filter----
7
8 # the `==` operator tests for equality
9
10 faculty %>%
11   filter(dept1 == "Sociology")
12
13 # the `|` operator signifies "or"
14
15 faculty %>%
16   filter(dept1 == "Sociology" |
17          dept1 == "Physics")
18
19 # the `%in%` operator allows for multiple options in a list
20
22:1 # filter
```

Environment

Object	Size
courses	104 obs. of 6 variables
faculty	392 obs. of 5 variables

Files

Name	Size	Modified
..		
sample.Rproj	258 B	Jul 10, 2022, 4:45 PM
NEAIR_code.R	373 B	Jul 10, 2022, 4:49 PM
courses.csv	3.8 KB	Jul 7, 2022, 11:59 AM
faculty.csv	16.6 KB	Jul 7, 2022, 11:27 AM

Console

```
R 4.2.0 ~ Desktop/sample/
chr (4): year, rank, dept1, dept2
dbl (1): id

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
> courses <- read_csv("courses.csv")
Rows: 104 Columns: 6
— Column specification —
Delimiter: ","
chr (2): dept, level
dbl (4): semester, course_id, faculty_id, enrollment

i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
>
```

ACCESSING WORKSHOP MATERIALS

ACCESSING DATA

12

use `read_csv()` to import a csv file

```
1 # the file path is this simple if you use projects!
2 # ?read_csv() in the console will bring up the help file with more options
3 faculty <- read_csv("faculty.csv")
```

the `readxl` package is helpful for Excel files

```
1 # needs to be loaded but not installed as it's part of the tidyverse
2 library(readxl)
3 faculty <- read_excel("faculty.xlsx", sheet = 2)
```

view the data with `View(faculty)` or by clicking on the data name in the Environment pane

MORE DATA MANIPULATION

12

STRINGR FUNCTIONS

12

functions from **stringr** (which all start with **str_**) are useful for working with text data

```
1 faculty %>%  
2   filter(str_detect(rank, "Professor"))
```

year	id	rank	dept1	dept2
2021-22	1022	Professor	Physics	Engineering
2021-22	1059	Professor	Physics	
2021-22	1086	Assistant Professor	Music	
2021-22	1118	Assistant Professor	Sociology	
2021-22	1158	Assistant Professor	Political Science	
2021-22	1161	Assistant Professor	Sociology	

STRINGR FUNCTIONS

12

cheat sheet of functions is [here](#)

```
1 courses %>%
2   mutate(year = str_c(str_sub(semester, 1, 4),
3                         "-",
4                         str_sub(semester, 5, 6))) %>%
5   select(semester, year) %>%
6   unique()
```

semester	year
20212202	2021-22
20212201	2021-22
20202102	2020-21
20202101	2020-21
20192002	2019-20
20192001	2019-20
20181902	2018-19
20181901	2018-19

PIVOTING DATA

12

existing **faculty** data has one row per faculty, some with multiple departments (sometimes known as *wide* data)

year	id	rank	dept1	dept2
2021-22	1005	Lecturer	Chemistry	
2021-22	1022	Professor	Physics	Engineering
2021-22	1059	Professor	Physics	
2021-22	1079	Lecturer	Music	
2021-22	1086	Assistant Professor	Music	
2021-22	1095	Adjunct Instructor	Sociology	

PIVOTING DATA

12

what if you instead want one row per faculty *per department*?
(sometimes known as *long data*)

year	id	rank	dept_no	dept
2021-22	1005	Lecturer	dept1	Chemistry
2021-22	1022	Professor	dept1	Physics
2021-22	1022	Professor	dept2	Engineering
2021-22	1059	Professor	dept1	Physics
2021-22	1079	Lecturer	dept1	Music
2021-22	1086	Assistant Professor	dept1	Music

PIVOTING DATA

12

the `pivot_longer` function lengthens data

```
1 faculty %>%  
2   pivot_longer(dept1:dept2,  
3               names_to = "dept_no",  
4               values_to = "dept",  
5               values_drop_na = TRUE) %>%  
6   select(-year, -rank)
```

id	dept_no	dept
1005	dept1	Chemistry
1022	dept1	Physics
1022	dept2	Engineering
1059	dept1	Physics
1079	dept1	Music
1086	dept1	Music

PIVOTING DATA

12

and `pivot_wider` does the opposite!

semester	course_id	faculty_id	dept	enrollment	level
20212202	10605	1772	Physics	7	UG
20212202	10605	1772	Physics	32	GR

```
1 courses %>%  
2   pivot_wider(names_from = "level",  
3               values_from = "enrollment")
```

semester	course_id	faculty_id	dept	UG	GR
20212202	10605	1772	Physics	7	32
20212202	11426	1820	Political Science	8	
20212202	12048	1914	English	24	
20212202	13269	1095	Sociology	48	

R has many useful functions for handling **relational data**

all you need is at least one **key variable** that connects data sets

left_join is most common, but there are **more**

JOINING DATA

12

what's the average UG enrollment per year, per faculty rank?

faculty

year	id	rank	dept1	dept2
2021-22	1005	Lecturer	Chemistry	
2021-22	1022	Professor	Physics	Engineering
2021-22	1059	Professor	Physics	
2021-22	1079	Lecturer	Music	

courses

semester	course_id	faculty_id	dept	enrollment	level
20212202	10605	1772	Physics	7	UG
20212202	10605	1772	Physics	32	GR
20212202	11426	1820	Political Science	8	UG
20212202	12048	1914	English	24	UG

`faculty$id` is the same as `courses$faculty_id`

what's the average **UG** enrollment per year, per faculty rank?

semester	course_id	faculty_id	dept	enrollment	level
20212202	10605	1772	Physics	7	UG
20212202	10605	1772	Physics	32	GR
20212202	11426	1820	Political Science	8	UG
20212202	12048	1914	English	24	UG
20212202	13269	1095	Sociology	48	UG

- filter to **UG** courses only
- create our **year** variable again
- summarize **enrollment** by **year** and **faculty_id**

use the `<-` operator to create a new data frame
`courses_UG`

```
1 courses_UG <- courses %>%  
2   filter(level == "UG") %>%  
3   mutate(year = str_c(str_sub(semester, 1, 4),  
4                        "-",  
5                        str_sub(semester, 5, 6)))
```

filter to undergraduate courses only and **mutate** a new academic year variable

```
1 courses_UG <- courses %>%  
2   filter(level == "UG") %>%  
3   mutate(year = str_c(str_sub(semester, 1, 4),  
4                       "-",  
5                       str_sub(semester, 5, 6)))
```

`group_by` year and faculty member; `summarize` enrollment

```
1 courses_UG <- courses %>%
2   filter(level == "UG") %>%
3   mutate(year = str_c(str_sub(semester, 1, 4),
4                        "-",
5                        str_sub(semester, 5, 6))) %>%
6   group_by(year, faculty_id) %>%
7   summarize(enr = sum(enrollment))
```

year	faculty_id	enr
2018-19	1059	35
2018-19	1086	14
2018-19	1102	37
2018-19	1203	25

what's the average UG enrollment per year, per faculty rank?

faculty

year	id	rank	dept1	dept2
2021-22	1005	Lecturer	Chemistry	
2021-22	1022	Professor	Physics	Engineering
2021-22	1059	Professor	Physics	
2021-22	1079	Lecturer	Music	
2021-22	1086	Assistant Professor	Music	
2021-22	1095	Adjunct Instructor	Sociology	

courses_UG

year	faculty_id	enr
2021-22	1005	50
2021-22	1086	17
2021-22	1095	48
2021-22	1128	32
2021-22	1147	32
2021-22	1191	7

JOINING DATA

12

1

2

```
1 fac_enr <- faculty %>%  
2   left_join(courses_UG, by = c("id" = "faculty_id",  
3                               "year" = "year"))
```

3

1. new data frame

2. data frame you're adding data to

3. data frame where the new data is coming from

year	id	rank	dept1	dept2	enr
2021-22	1005	Lecturer	Chemistry		50
2021-22	1022	Professor	Physics	Engineering	
2021-22	1059	Professor	Physics		
2021-22	1079	Lecturer	Music		
2021-22	1086	Assistant Professor	Music		17
2021-22	1095	Adjunct Instructor	Sociology		48

what's the average UG enrollment per year, per faculty rank?

```
1 fac_enr <- faculty %>%  
2   left_join(courses_UG, by = c("id" = "faculty_id",  
3                               "year" = "year")) %>%  
4   group_by(year, rank) %>%  
5   summarize(avg_enr = mean(enr, na.rm = TRUE))
```

year	rank	avg_enr
2021-22	Adjunct Instructor	34.66667
2021-22	Assistant Professor	23.60000
2021-22	Associate Professor	17.25000
2021-22	Lecturer	31.83333
2021-22	Professor	32.16667
2021-22	Visiting Researcher	

DATA VISUALIZATION

12

`ggplot2` is the data visualization package that is loaded with the `tidyverse`

the `g`rammar of `g`raphics maps data to the aesthetic attributes of geometric points

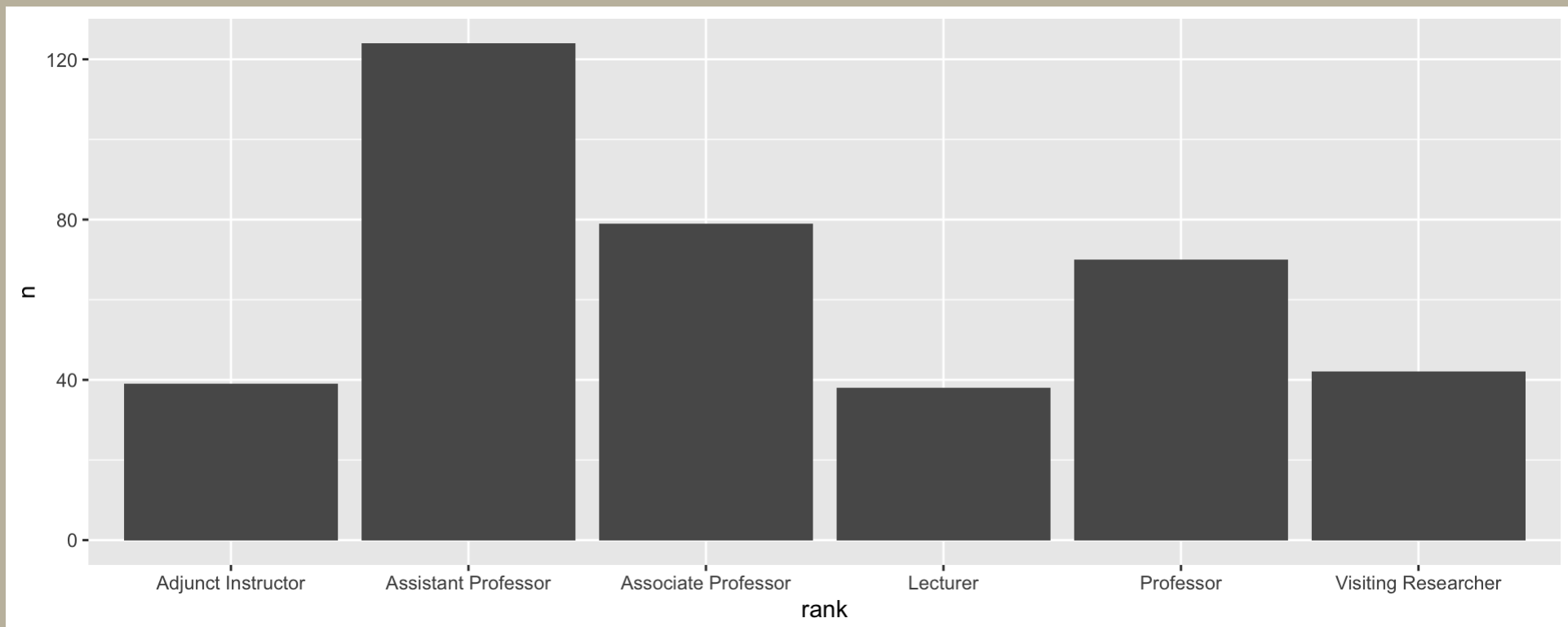
encoding data into visual cues (e.g., length, color, position, size) is how we signify changes and comparisons

BAR CHART

12

```
1 faculty %>%  
2   count(rank) %>%  
3   ggplot(aes(x = rank, y = n)) +  
4   geom_bar(stat = "identity")
```

to combine lines into one code chunk, use **+** instead of **%>%**

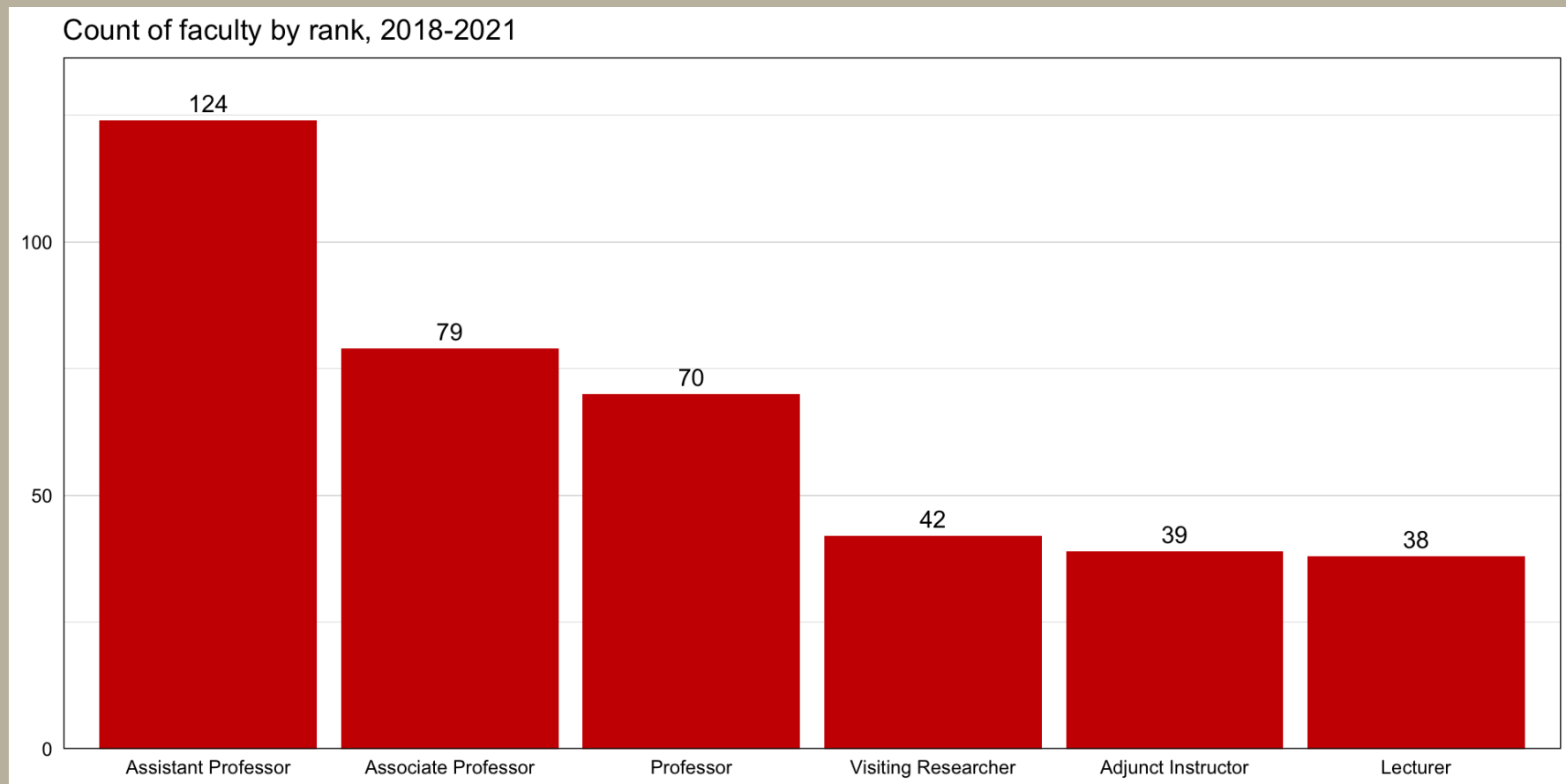


BAR CHART

12

can create a prettier plot pretty easily

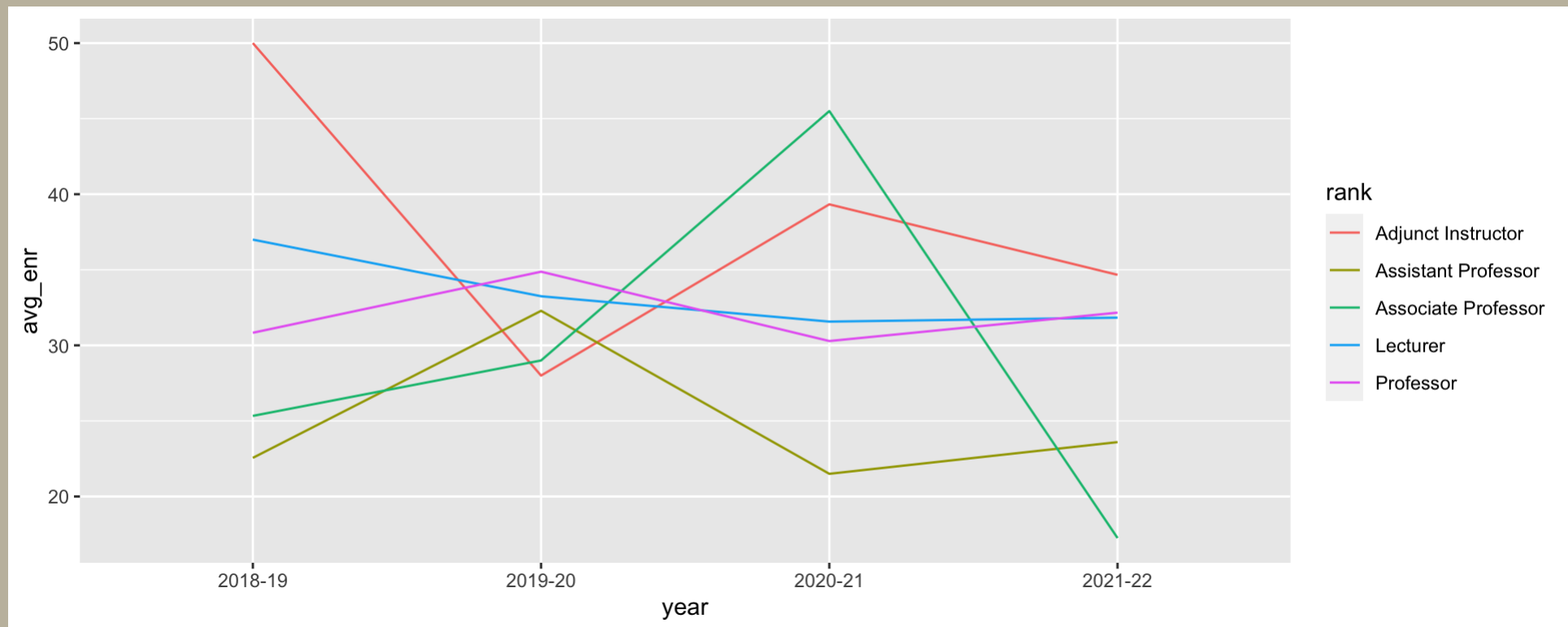
► expand for full code



LINE GRAPH

12

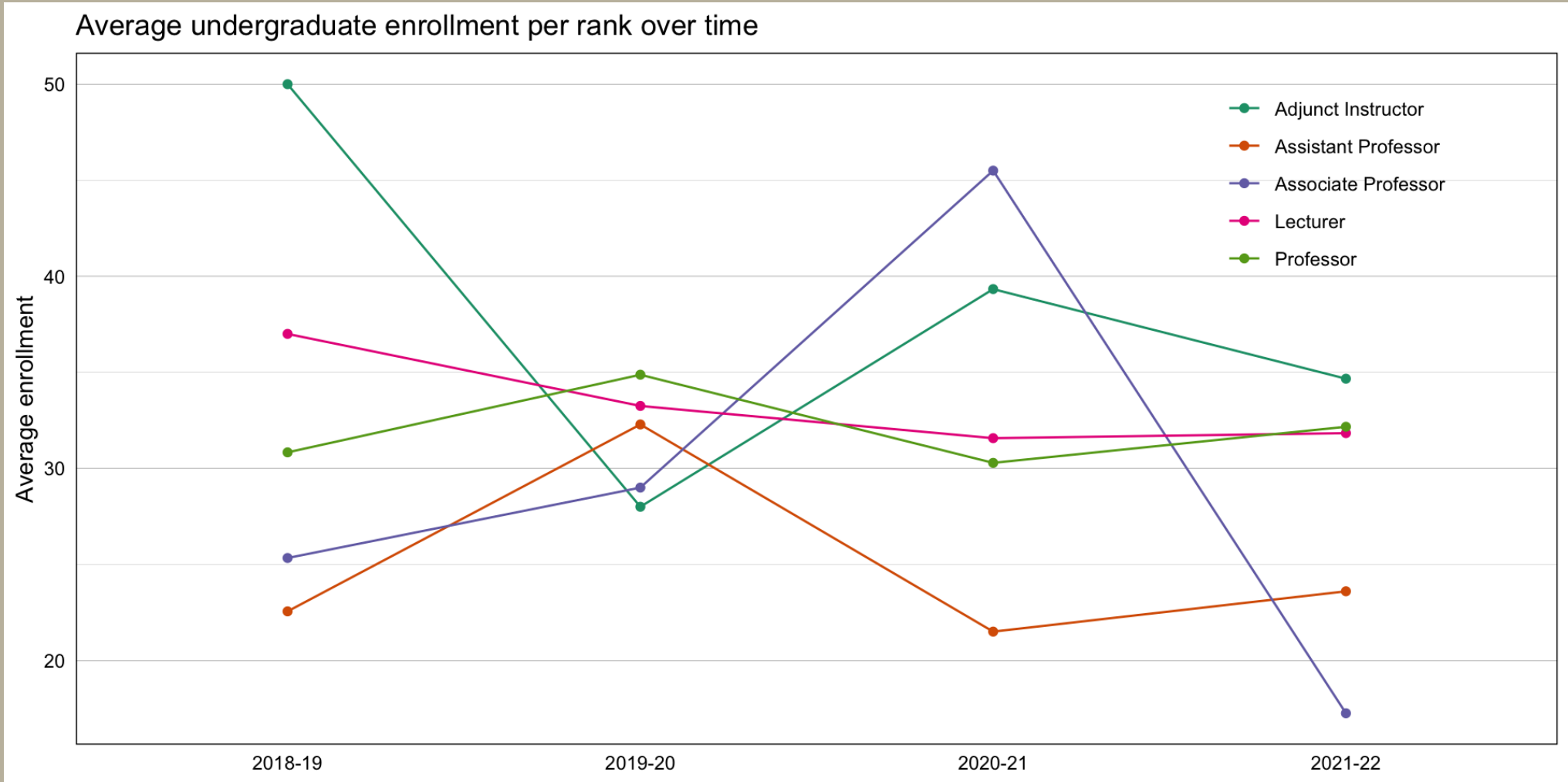
```
1 fac_enr %>%  
2   filter(!is.na(avg_enr)) %>%  
3   ggplot(aes(x = year, y = avg_enr, group = rank, color = rank)) +  
4   geom_line()
```



LINE GRAPH

12

► expand for full code



GGPLOT2 RESOURCES

12

from R for Data Science

Data Visualization: a practical introduction

creating custom themes

the ggplot2 book

the R graph gallery

PUTTING IT ALL TOGETHER

12

with what we've done so far, your **.R** file could:

- import your data files
- document all data cleaning and preparation steps and decisions
- produce a PPT-ready graphic summarizing your results

and that file would make it extremely easy for you or someone else to reproduce this analysis with new data in six months

ADVANCED TOPICS

12

using RStudio, create `. Rmd` documents that combine text, code, and graphics

many output formats: html, pdf, Word, slides

exceedingly useful for **parameterized reporting**: can create an R-based PDF report and generate it automatically for, say, each department

INTERNAL PACKAGES

12

you can also create your own packages!

your package can hold:

- common data sets that are used across projects
- custom **ggplot2** themes
- common functions and calculations (and their definitions!)

can be stored on a shared drive to facilitate collaboration

R MARKDOWN AND PACKAGE RESOURCES

12

R Markdown

the official [R Markdown website](#)

[R Markdown: The Definitive Guide](#)

internal packages

a comprehensive [theoretical explainer](#)

a [talk I gave](#) earlier this year on the topic

LEARN MORE ABOUT R

RESOURCES

R for Data Science: the ultimate guide

R for Excel users: a very useful workshop

STAT 545: an online book on reproducible data analysis in R

the RStudio Education site

the Learn tidyverse site