

**MINE ÇETINKAYA-RUNDEL**

DUKE UNIVERSITY + RSTUDIO

@minebocek

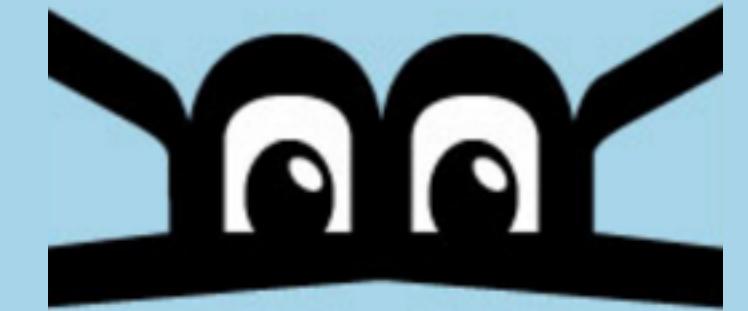
mine-cetinkaya-rundel

cetinkaya.mine@gmail.com



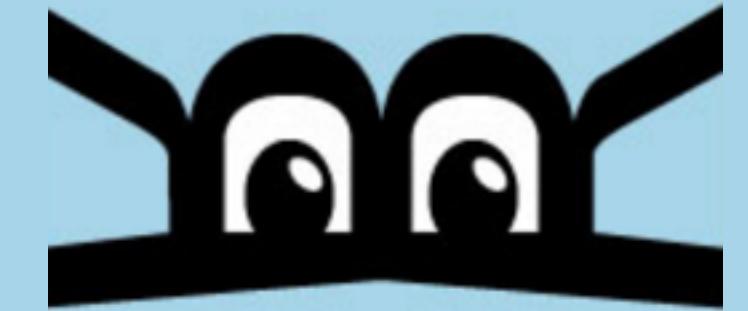
 [bit.ly/dsbox-dscwav](https://bit.ly/dsbox-dscwav)

# Three questions that keep me up at night...

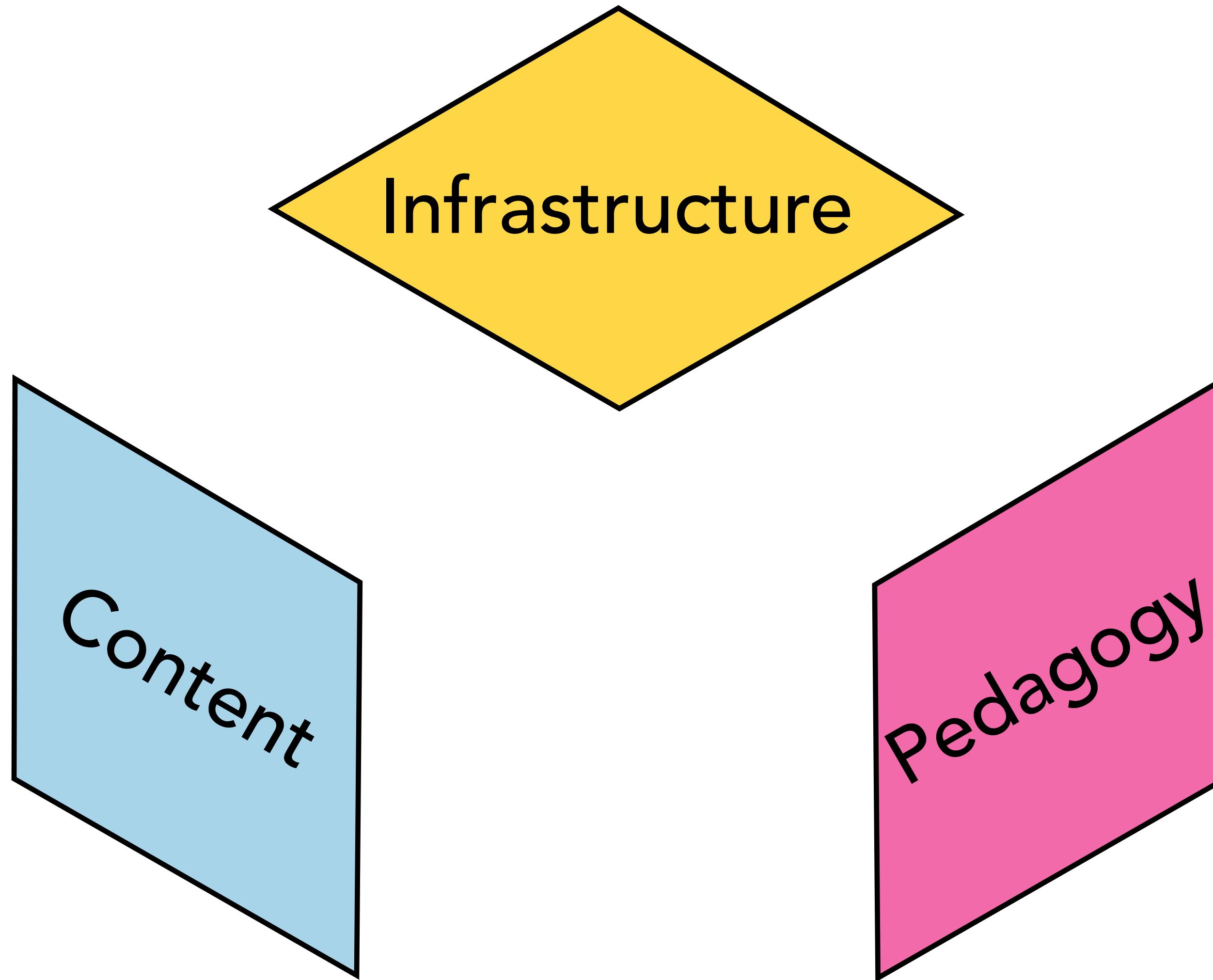


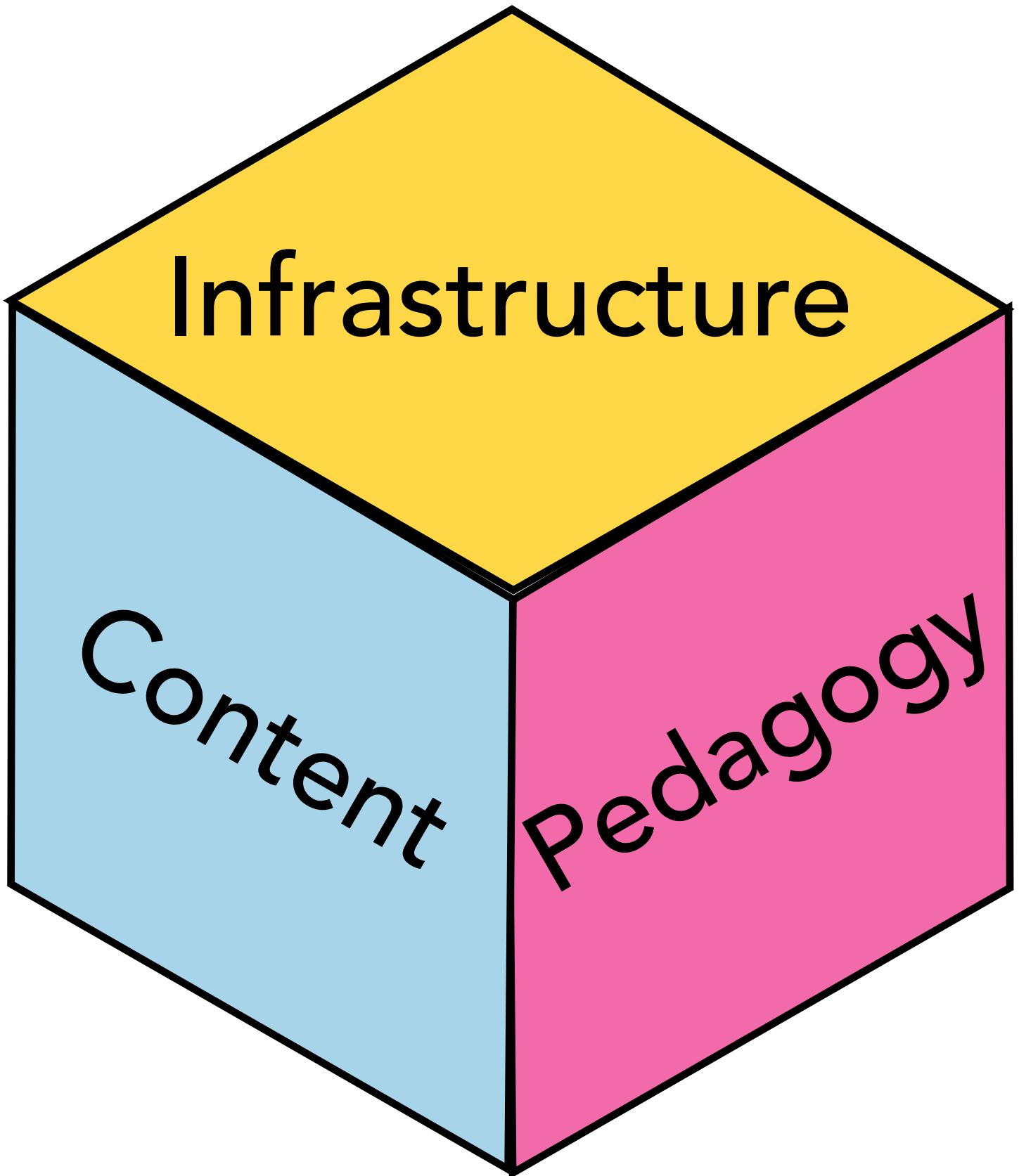
- 1 What should my students learn?
- 2 How will my students learn best?
- 3 What tools will enhance my students' learning?

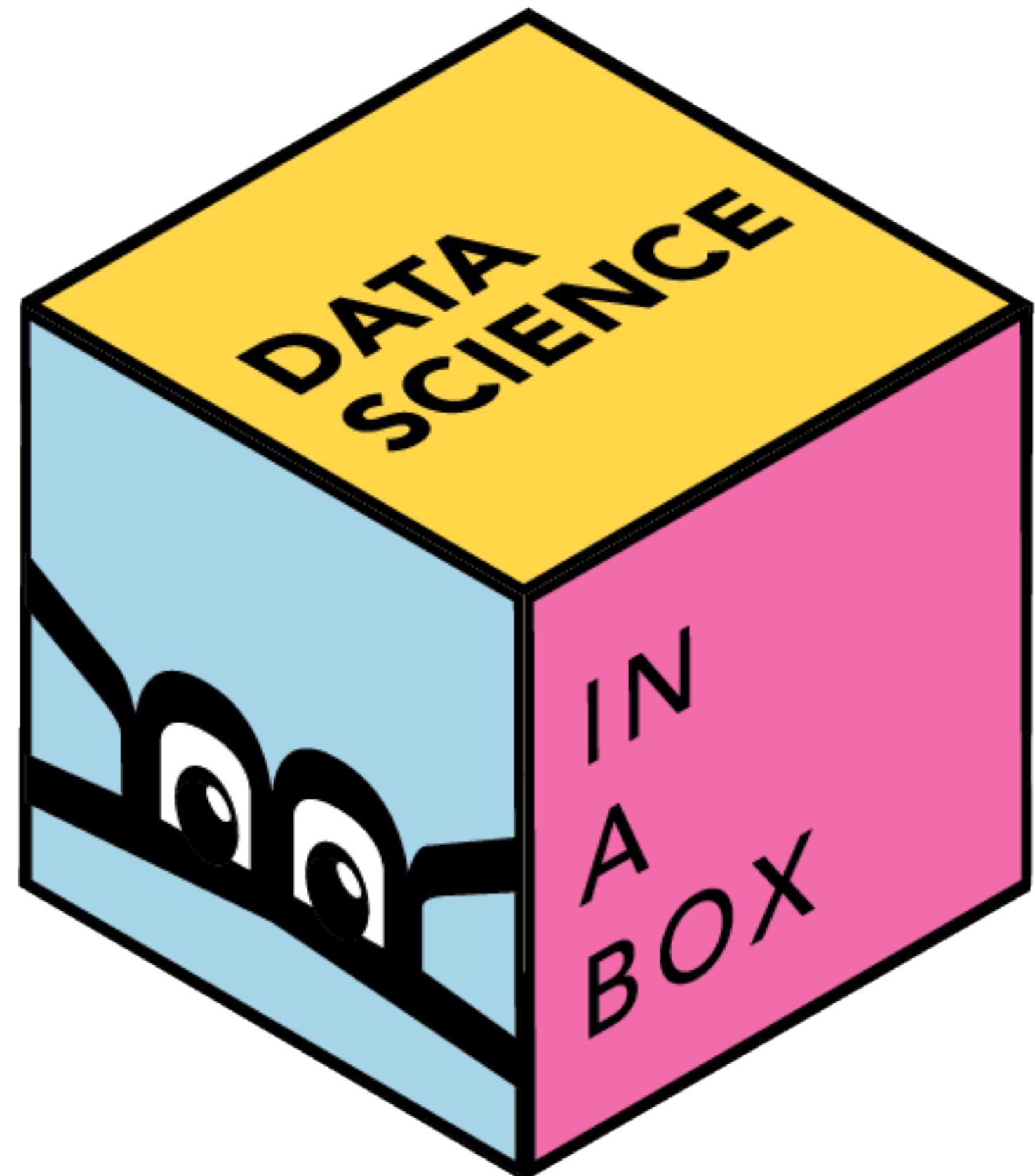
# Three questions that keep me up at night...



- |                       |   |  |
|-----------------------|---|--|
| <b>Content</b>        | 1 | What should my students learn?                 |
| <b>Pedagogy</b>       | 2 | How will my students learn best?               |
| <b>Infrastructure</b> | 3 | What tools will enhance my students' learning? |







**rstudio-education/datasiencebox.org**

## 3 Topics

The course content is organized in three units:

**Unit 1 - Hello world:** This unit is an introduction to the concepts of data science.

**Unit 2 - Exploring data:** This unit covers fundamental concepts of data exploration, including handling missing values, dealing with outliers, and understanding confounding variables, and introduces students to tidy data, data import, and data visualization.

**Unit 3 - Data science ethics:** In this unit we discuss misrepresentation of findings, particularly in data visualisations, breaches of data privacy, and algorithmic bias.

**Unit 4 - Making rigorous conclusions:** In this unit we introduce modelling and statistical inference for making data-based conclusions. We discuss building, interpreting, and selecting models, visualizing interaction effects, and prediction and model validation. Statistical inference is introduced from a simulation based

17 Sharing

**dsbox**

The goal of dsbox is to supplement the Data Science Course in a Box project. The package contains the datasets that are used in the materials in Data Science Course in a Box as well as the learnr tutorials.

## Installation

dsbox is not yet on CRAN. For now, you can install it from GitHub with

```
# install.packages("devtools")
devtools::install_github("rstudio-education/dsbox")
```

## Questions, bugs, and feature requests

You can file an issue on the GitHub repository. When filing an issue, please include a minimal reproducible example using the `reprex` package. If you haven't used `reprex` before, don't worry—seriously, `reprex` will make all of your R-question-asking endeavors easier (which is seriously cool). It's a great way to learn what it's all about. For additional `reprex` pointers, check out the [Get started with reprex](#).

If you're reporting a bug, be sure to [search issues and pull requests](#) to make sure the bug hasn't been reported and/or already fixed by someone else. You can search for issues by specifying the `repo` argument and the `is:issue` or `is:open` qualifier. You can also search for pull requests by specifying the `repo` argument and the `is:pr`, `is:closed` qualifiers as needed. For example, you'd simply remove `is:open` to search all issues in the repo, open or closed.

## Code of Conduct

Please note that the dsbox project is released with a [Contributor Code of Conduct](#). By contributing to this project, you agree to abide by its terms.

**Links**

- [Browse source code](#)
- [Report a bug](#)
- [View on GitHub](#)

**Community**

- [Code of conduct](#)
- [Citation](#)
- [Citing dsbox](#)

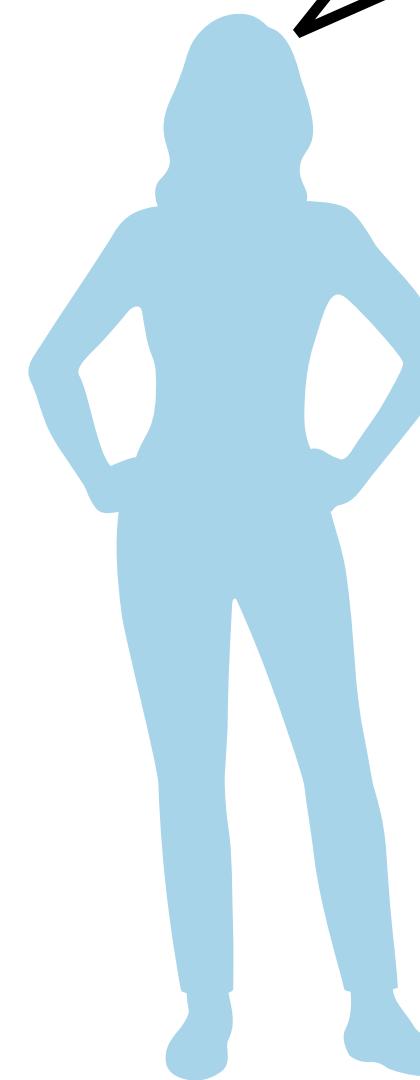
**Developers**

- Mine Çetinkaya-Rundel  
Author, maintainer
- RStudio  
Copyright holder, funder
- [More about authors...](#)

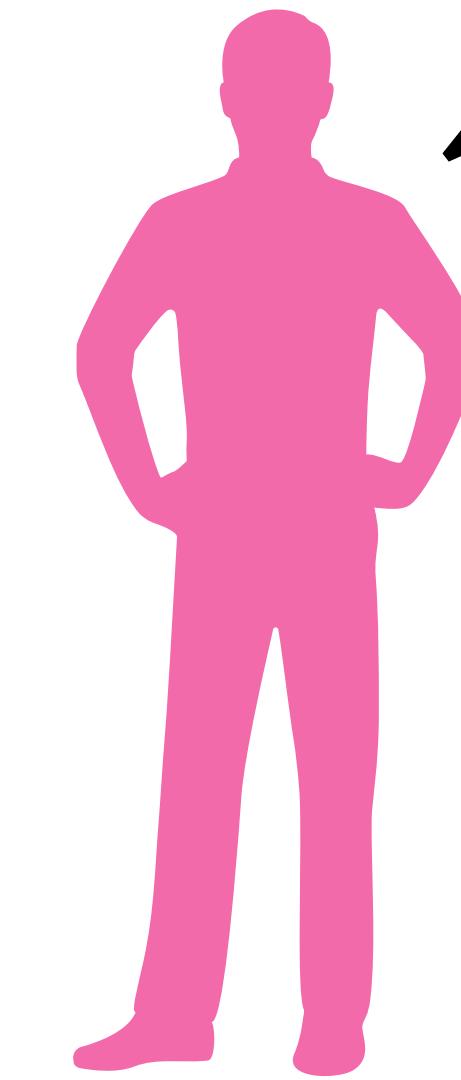
**Dev status**

- R-CMD check passing

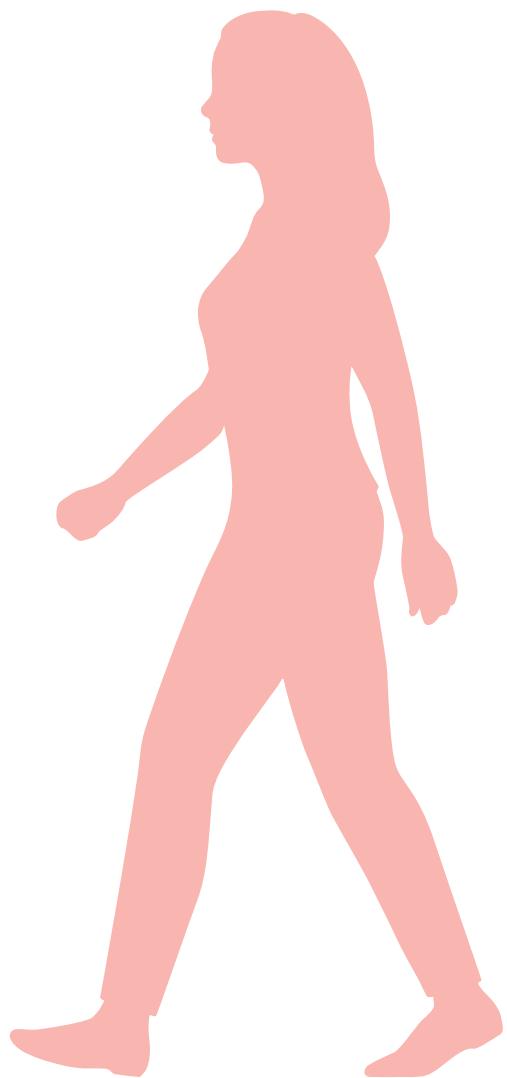
## AUDIENCE



I have been teaching with R for a while, but I want to update my teaching materials

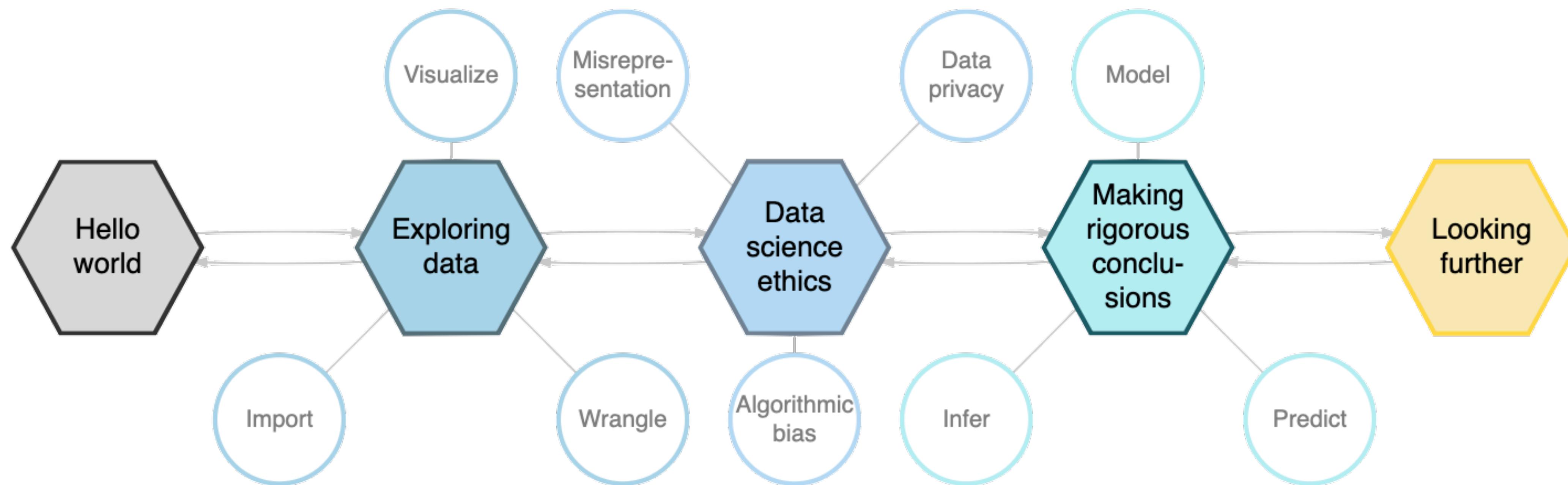


I'm new to teaching with R and need to build up my course materials



This teaching slide deck I came across on Twitter is pretty cool, but I have no idea what type of course it belongs in

# TOPICS



Fundamentals of  
data & data viz,  
confounding variables,  
Simpson's paradox  
+  
R / RStudio,  
R Markdown, simple Git

Tidy data, data frames  
vs. summary tables,  
recoding & transforming,  
web scraping & iteration  
+  
collaboration on GitHub

Building & selecting  
models,  
visualizing  
interactions,  
prediction &  
validation, inference  
via simulation

Interactive viz &  
reporting, text  
analysis,  
Bayesian inference  
+  
communication &  
dissemination

# CONTENTS



website

[datasciencebox.org](http://datasciencebox.org)



48  
slide  
decks



48  
videos



10  
application  
exercises



14  
computing  
labs



repository



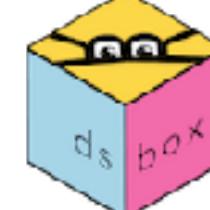
10  
homework  
assignments



2  
take-home  
exams



1  
open-ended  
project



package

dsbox



9  
interactive  
tutorials



9  
interactive  
tutorials

## DESIGN PRINCIPLES



cherish  
day one



start  
with cake



skip baby  
steps



hide the  
veggies



leverage the  
ecosystem

# DESIGN PRINCIPLES



Which kitchen would you  
rather bake a cake?



# DESIGN PRINCIPLES



Which kitchen would you  
rather bake a cake?





# Cherish day one

The screenshot shows a web browser window for "RStudio Cloud :: Data Science in a Box". The URL is <https://datasciencebox.org/infrastructure/rscloud/>. The page title is "RStudio Cloud". On the left, there's a sidebar with a "DATA SCIENCE IN A BOX" logo, a search bar, and a navigation menu with options like "Hello #dsbox", "Course content", "Infrastructure", "Pedagogy", and "Built with ❤️ and blogdown, logo by muruge". The main content area has a heading "RStudio Cloud" and a list of bullet points:

- Setting up your course in RStudio Cloud
- Projects
- Base project template
- Git integration
- Limits
- RStudio Cloud is in alpha!
- Learn more

Below the list, there's a paragraph about the RStudio IDE, followed by a section about RStudio Cloud, and two paragraphs at the bottom.

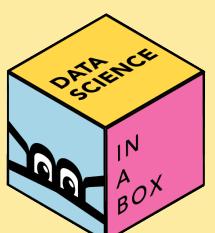
The RStudio IDE includes a viewable environment, a file browser, data viewer, and a plotting pane, which makes it less intimidating than the bare R shell. Additionally, since it is a full fledged IDE, it also features integrated help, syntax highlighting, and context-aware tab completion, which are all powerful tools that help flatten the learning curve.

RStudio Cloud is a managed cloud instance of the RStudio IDE. We recommend having students access RStudio via RStudio Cloud as opposed to using a local installation. The main reason for this choice is reducing friction at first exposure to R. Local installation can be difficult to manage, both for the student and the instructor, and can shift the focus away from data science learning at the beginning of the course. In the pedagogical decisions section we discuss in further detail the reasons for avoiding local installation at the beginning of the course and discuss when to introduce it later on in the course.

When you create an account on RStudio Cloud you get a workspace of your own, and the projects you create here are public to RStudio Cloud members. You can also add a new workspace and control its permissions, and the projects you create here can be public or private.

All student facing materials for this course have been organized in an RStudio Cloud workspace [here](#). Soon you will have the option to copy that workspace and use it to organize assignments and assessments. [Note: The workspace is currently work in progress, rest of the materials will be added soon.]

A natural way to set up a course in RStudio Cloud is using a private workspace. In this structure a classroom (a cohort of students in one semester of the course) maps to a workspace. Once a workspace is set up, instructors can invite students to the workspace via an invite link. Workspaces allow for various permission levels which can be assigned to students, teaching assistants, and instructors. Then, each assignment/project in the course maps to an RStudio Cloud project.



## Ingredients

### For the Cake:

16 ounces plain or **toasted sugar** (about 2 1/4 cups; 455g)

4 1/2 teaspoons baking powder

2 teaspoons (8g) Diamond Crystal kosher salt; for table salt, use about half as much by volume or the same weight

8 ounces unsalted butter (16 tablespoons; 225g), soft but cool, about 60°F (16°C)

3 large eggs, brought to about 65°F (18°C)

1/2 ounce vanilla extract (about 1 tablespoon; 15g)

16 ounces whole milk (about 2 cups; 455g), brought to about 65°F (18°C)

16 ounces all-purpose flour (about 3 1/2 cups, spooned; 455g)

## Directions

1. **For the Cake:** Adjust oven rack to lower-middle position and preheat to 350°F (180°C). Lightly grease three 8-inch anodized aluminum cake pans and line with parchment (explanation and tutorial [here](#)). If you don't have three pans, it's okay to bake the cakes in stages, the batter will keep at room temperature until needed.
2. In the bowl of a stand mixer fitted with the paddle attachment, combine sugar, baking powder, salt, and butter. Mix on low speed to roughly incorporate, then increase to medium and beat until fluffy and light, about 5 minutes. About halfway through, pause to scrape the bowl and beater with a flexible spatula.
3. With the mixer still running, add the eggs one at a time, letting each fully incorporate before adding the next, then dribble in the vanilla. Reduce speed to low and sprinkle in about 1/3 of the flour, then drizzle in 1/3 of the milk. Repeat with remaining flour and milk, working in thirds as before.
4. Scrape the bowl and beater with a flexible spatula, and resume mixing on medium speed for about 3 seconds to ensure everything is well combined. The batter should look creamy and thick, registering between 65 and 68°F (18 and 20°C) on a digital thermometer. (Significant
5. Fold batter once or twice from the bottom up with a flexible spatula, then divide evenly between prepared cake pans (about 20 ounces or 565g if you have a scale). Stagger pans together on the oven rack, and bake until puffed, firm, and pale gold, about 32 minutes. If your oven has very uneven heat, pause to rotate the pans after about 20 minutes. Alternatively, bake two layers at once and finish the third when they're done.
6. Cool cakes directly in their pans for 1 hour, then run a butter knife around the edges to loosen. Invert onto a wire rack, peel off the parchment, and return cakes right-side-up (covered in plastic, the cakes can be left at room temperature for a few hours). Prepare the buttercream.

# How do you prefer your cake recipes? Words only, or words & pictures?



## Ingredients

### For the Cake:

16 ounces plain or **toasted sugar**  
(about 2 1/4 cups; 455g)

4 1/2 teaspoons baking powder

2 teaspoons (8g) Diamond  
Crystal kosher salt; for table salt,  
use about half as much by  
volume or the same weight

8 ounces unsalted butter (16  
tablespoons; 225g), soft but  
cool, about 60°F (16°C)

3 large eggs, brought to about  
65°F (18°C)

1/2 ounce vanilla extract (about 1  
tablespoon; 15g)

16 ounces whole milk (about 2  
cups; 455g), brought to about  
65°F (18°C)

16 ounces all-purpose flour  
(about 3 1/2 cups, spooned;  
455g)

## Directions

1. **For the Cake:** Adjust oven rack to lower-middle position and preheat to 350°F (180°C). Lightly grease three 8-inch anodized aluminum cake pans and line with parchment (explanation and tutorial [here](#)). If you don't have three pans, it's okay to bake the cakes in stages, the batter will keep at room temperature until needed.
2. In the bowl of a stand mixer fitted with the paddle attachment, combine sugar, baking powder, salt, and butter. Mix on low speed to roughly incorporate, then increase to medium and beat until fluffy and light, about 5 minutes. About halfway through, pause to scrape the bowl and beater with a flexible spatula.
3. With the mixer still running, add the eggs one at a time, letting each fully incorporate before adding the next, then dribble in the vanilla. Reduce speed to low and sprinkle in about 1/3 of the flour, then drizzle in 1/3 of the milk. Repeat with remaining flour and milk, working in thirds as before.
4. Scrape the bowl and beater with a flexible spatula, and resume mixing on medium speed for about 3 seconds to ensure everything is well combined. The batter should look creamy and thick, registering between 65 and 68°F (18 and 20°C) on a digital thermometer. (Significant
5. Fold batter once or twice from the bottom up with a flexible spatula, then divide evenly between prepared cake pans (about 20 ounces or 565g if you have a scale). Stagger pans together on the oven rack, and bake until puffed, firm, and pale gold, about 32 minutes. If your oven has very uneven heat, pause to rotate the pans after about 20 minutes. Alternatively, bake two layers at once and finish the third when they're done.
6. Cool cakes directly in their pans for 1 hour, then run a butter knife around the edges to loosen. Invert onto a wire rack, peel off the parchment, and return cakes right-side-up (covered in plastic, the cakes can be left at room temperature for a few hours). Prepare the buttercream.

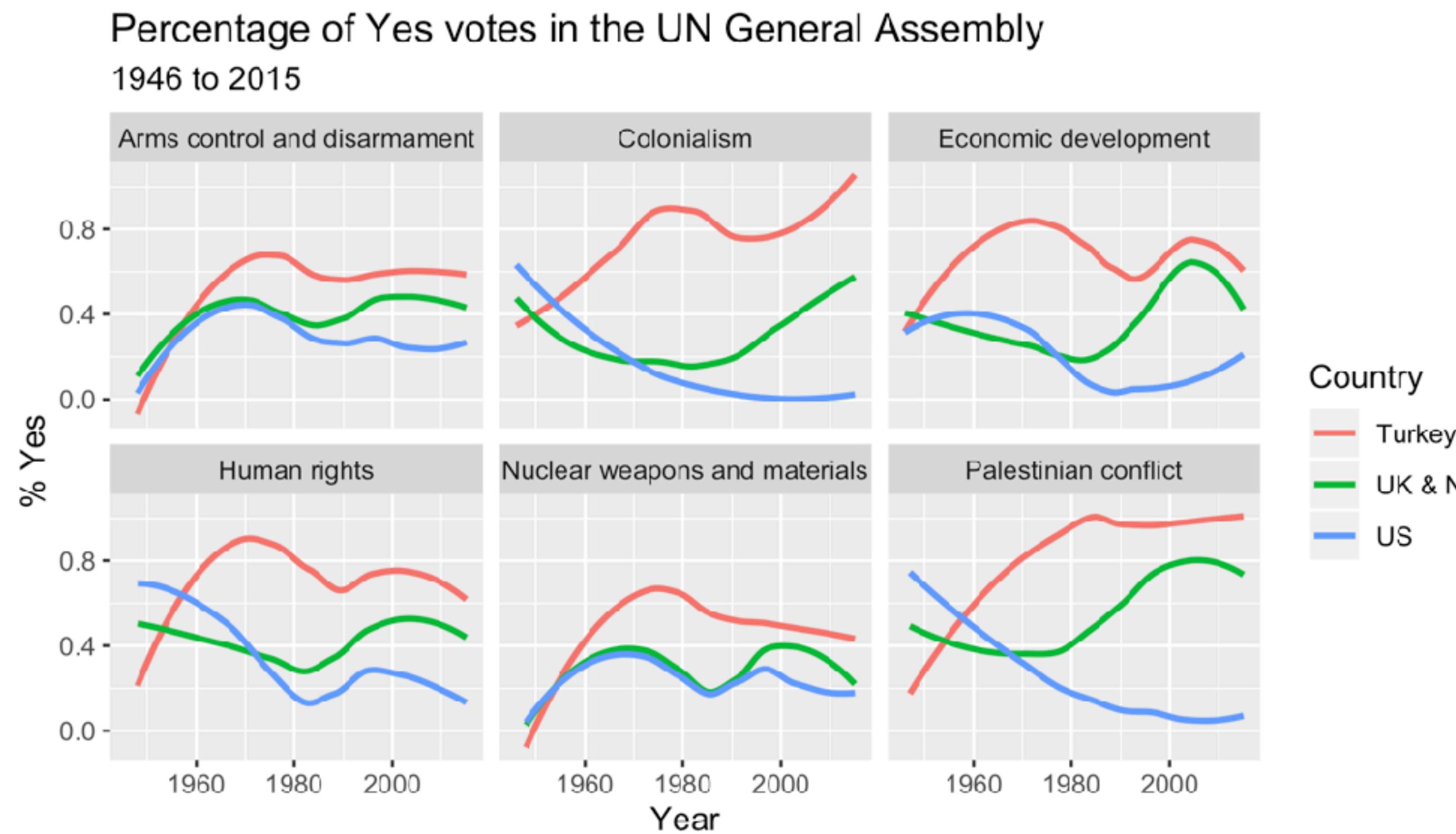
# How do you prefer your cake recipes? Words only, or words & pictures?





# Start with cake

- ▶ Open today's demo project
- ▶ Knit the document and discuss the results with your neighbor



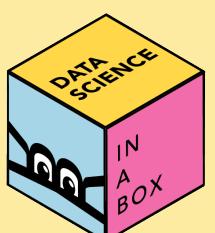
- ▶ Then, change Turkey to a different country, and plot again



# Start with cake

With great examples, comes a great amount of code...  
but let's focus on the task at hand...

- ▶ Open today's demo project
- ▶ Knit the document and discuss the results with your neighbor
- ▶ Then, change Turkey to a different country, and plot again





# Start with cake

```
un_votes %>%
  filter(country %in% c("UK & NI", "US", "Turkey")) %>%
  inner_join(un_roll_calls, by = "rcid") %>%
  inner_join(un_roll_call_issues, by = "rcid") %>%
  group_by(country, year = year(date), issue) %>%
  summarize(
    votes = n(),
    percent_yes = mean(vote == "yes")
  ) %>%
  filter(votes > 5) %>% # only use records where there are more than 5 votes
  ggplot(mapping = aes(x = year, y = percent_yes, color = country)) +
  geom_smooth(method = "loess", se = FALSE) +
  facet_wrap(~ issue) +
  labs(
    title = "Percentage of Yes votes in the UN General Assembly",
    subtitle = "1946 to 2015",
    y = "% Yes",
    x = "Year",
    color = "Country"
  )
```



# Start with cake

```
un_votes %>%  
  filter(country %in% c("UK & NI", "US", "Turkey")) %>%  
  inner_join(un_roll_calls, by = "rcid") %>%  
  inner_join(un_roll_call_issues, by = "rcid") %>%  
  group_by(country, year = year(date), issue) %>%  
  summarize(  
    votes = n(),  
    percent_yes = mean(vote == "yes")  
  ) %>%  
  filter(votes > 5) %>% # only use records where there are more than 5 votes  
  ggplot(mapping = aes(x = year, y = percent_yes, color = country)) +  
  geom_smooth(method = "loess", se = FALSE) +  
  facet_wrap(~ issue) +  
  labs(  
    title = "Percentage of Yes votes in the UN General Assembly",  
    subtitle = "1946 to 2015",  
    y = "% Yes",  
    x = "Year",  
    color = "Country"  
  )
```



# Start with cake

```
un_votes %>%
  filter(country %in% c("UK & NI", "US", "Turkey")) %>%
  inner_join(un_roll_calls, by = "rcid") %>%
  inner_join(un_roll_call_issues, by = "rcid") %>%
  group_by(country, year = year(date), issue) %>%
  summarize(
    votes = n(),
    percent_yes = mean(vote == "yes")
  ) %>%
  filter(votes > 5) %>% # only use records where there are more than 5 votes
  ggplot(mapping = aes(x = year, y = percent_yes, color = country)) +
  geom_smooth(method = "loess", se = FALSE) +
  facet_wrap(~ issue) +
  labs(
    title = "Percentage of Yes votes in the UN General Assembly",
    subtitle = "1946 to 2015",
    y = "% Yes",
    x = "Year",
    color = "Country"
  )
```



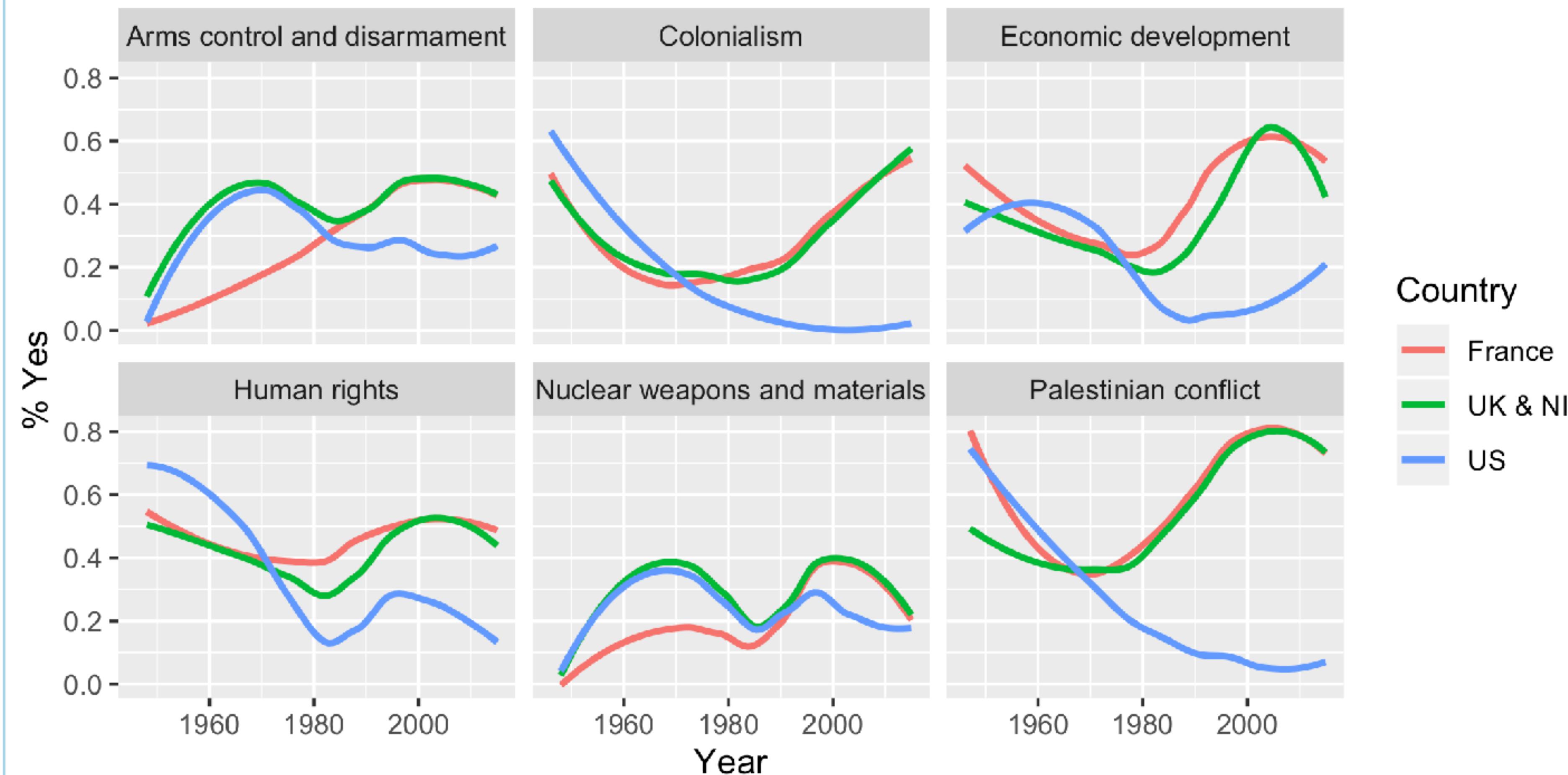
# Start with cake

```
un_votes %>%
  filter(country %in% c("UK & NI", "US", "France")) %>%
  inner_join(un_roll_calls, by = "rcid") %>%
  inner_join(un_roll_call_issues, by = "rcid") %>%
  group_by(country, year = year(date), issue) %>%
  summarize(
    votes = n(),
    percent_yes = mean(vote == "yes")
  ) %>%
  filter(votes > 5) %>% # only use records where there are more than 5 votes
  ggplot(mapping = aes(x = year, y = percent_yes, color = country)) +
  geom_smooth(method = "loess", se = FALSE) +
  facet_wrap(~ issue) +
  labs(
    title = "Percentage of Yes votes in the UN General Assembly",
    subtitle = "1946 to 2015",
    y = "% Yes",
    x = "Year",
    color = "Country"
  )
```



# Start with cake

Percentage of Yes votes in the UN General Assembly  
1946 to 2015



# DESIGN PRINCIPLES



Which motivates you more to learn how to cook: perfectly chopped onions or ratatouille?

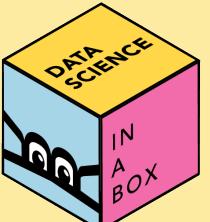


# DESIGN PRINCIPLES



Which motivates you more to learn how to cook: perfectly chopped onions or ratatouille?





# Re-insert ~~Skip baby steps~~

## Visualizing data

Data visualization with ggplot2

The data: Star Wars

Scatterplots

Setting aesthetic features

Faceting your visualizations

Data types

Univariate analysis

[Start Over](#)

## Scatterplots

How can we visualize the relationship between characters' heights and masses? Following the structure of the `ggplot` function that we laid out earlier, we pass `starwars` to the `data` argument, and map `height` and `mass` to the `x` and `y` `aes` thetics, respectively. Then, we specify on the next layer that we would like the data points to be represented by points with `geom_point`.

Fill in the blanks below to create the scatterplot.

Code

Start Over

Solution

Run Code

Submit Answer

```
1 ggplot(data = ___, mapping = aes(x = ___, y = ___)) +  
2   ___  
3   ___
```

Notice the warning that tells us that 28 of the observations have not been graphed, which means that some of the necessary information (height and mass) was missing for those rows.

Your turn!

**How would you describe the relationship between height and weight?**

- positive and nonlinear
- positive and linear
- negative and nonlinear
- negative and linear

[Submit Answer](#)

**How many outliers does the graph show?**

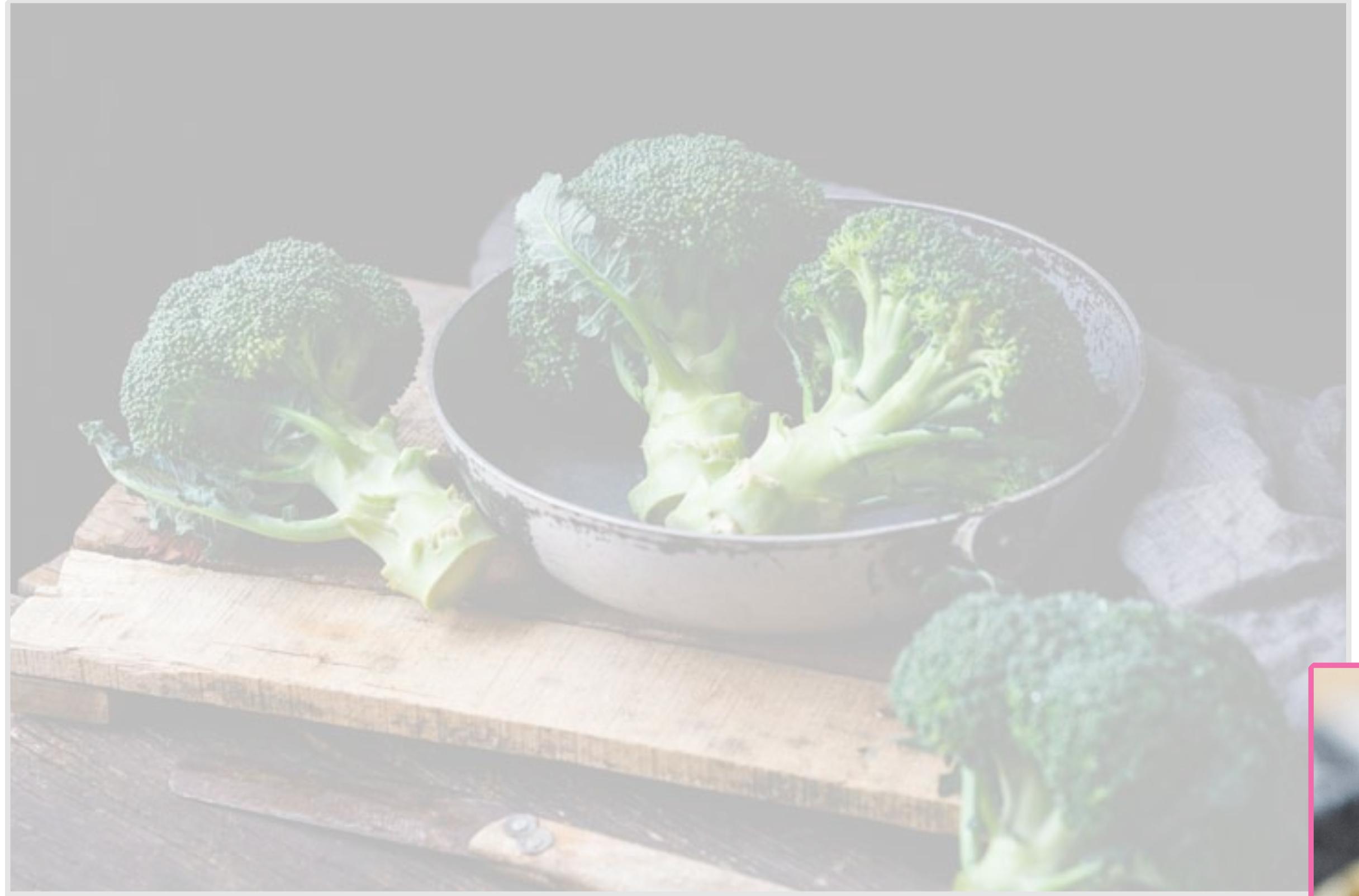
- 0
- 1
- 2

[Submit Answer](#)



Which is more likely to appeal to someone who has never tried broccoli?





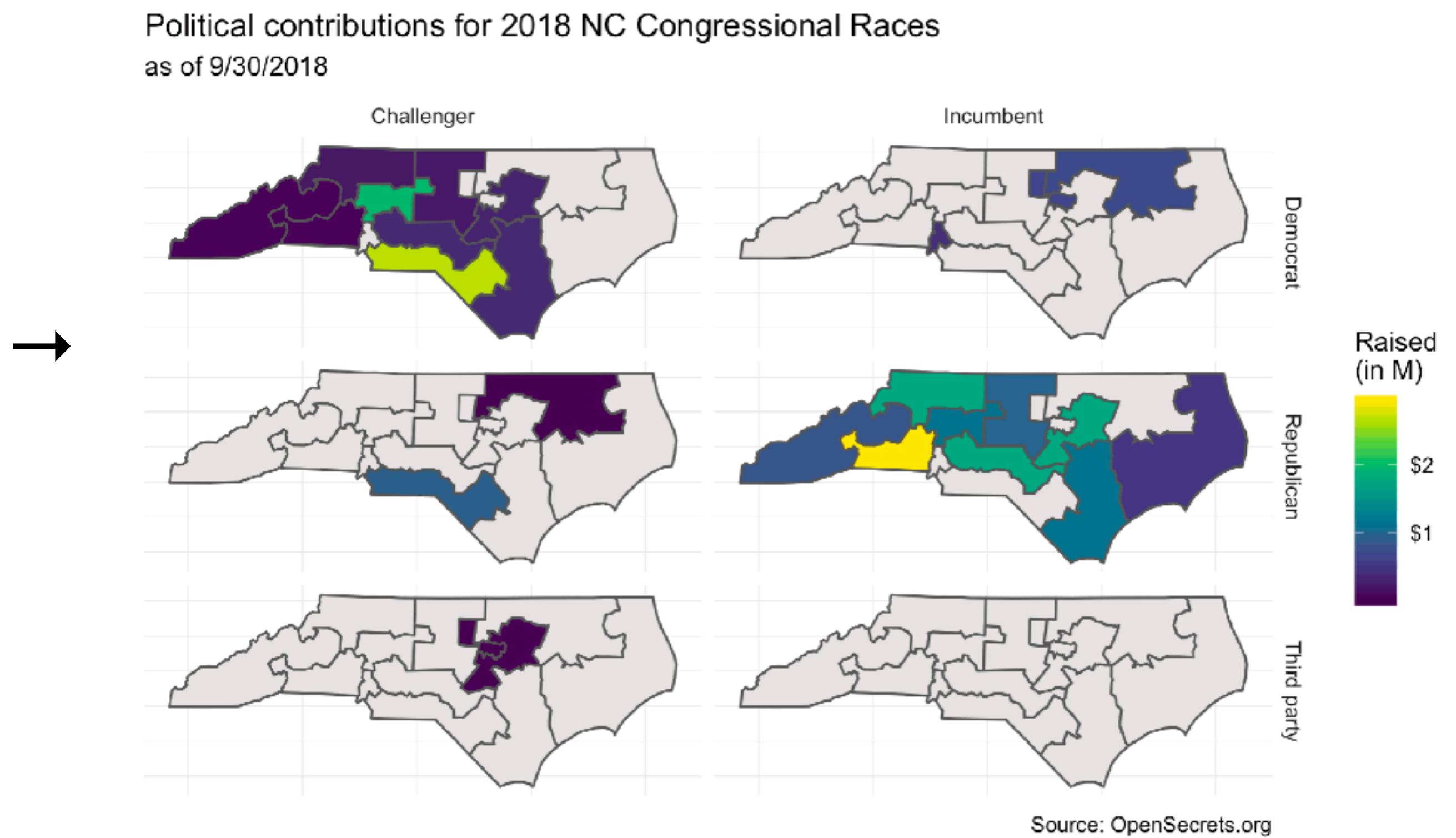
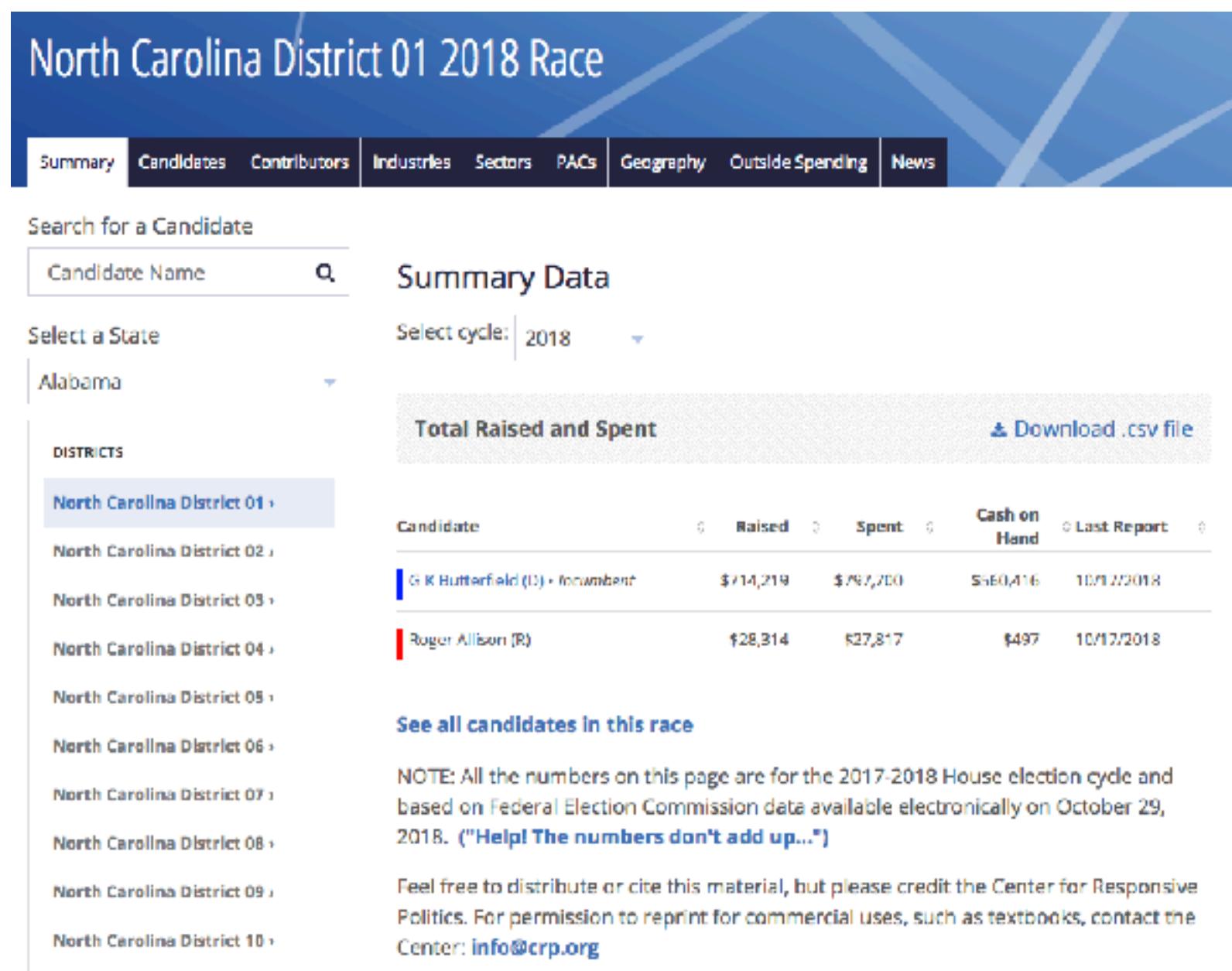
Which is more likely to appeal to someone who has never tried broccoli?





# Hide the veggies

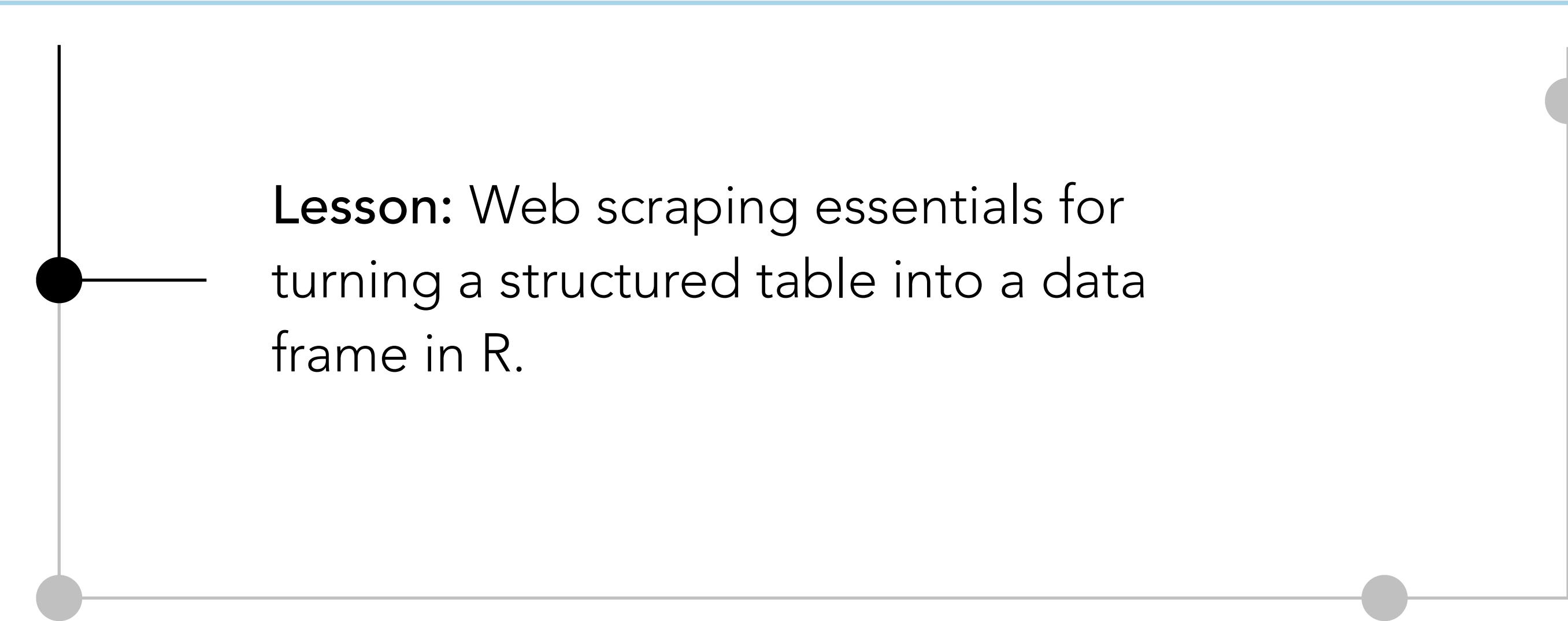
- ▶ Today we go from this to that



- ▶ And do so in a way that is easy to replicate for another state



# Hide the veggies





# Hide the veggies

Lesson: Web scraping essentials for turning a structured table into a data frame in R.

Ex 1: Scrape the table off the web and save as a data frame.

Candidate	Raised	Spent	Cash on Hand	Last Report
G K Butterfield (D) • Incumbent	\$714,219	\$797,700	\$560,416	10/17/2018
Roger Allison (R)	\$28,314	\$27,817	\$497	10/17/2018



#	candidate_info	raised	spent	cash_on_hand	last_report	race
1	G K Butterfield (D) • Incumbent	714219	797700	560416	2018-10-17	North Carolina District 01
2	Roger Allison (R)	28314	27817	497	2018-10-17	North Carolina District 01



# Hide the veggies

**Lesson:** Web scraping essentials for turning a structured table into a data frame in R.

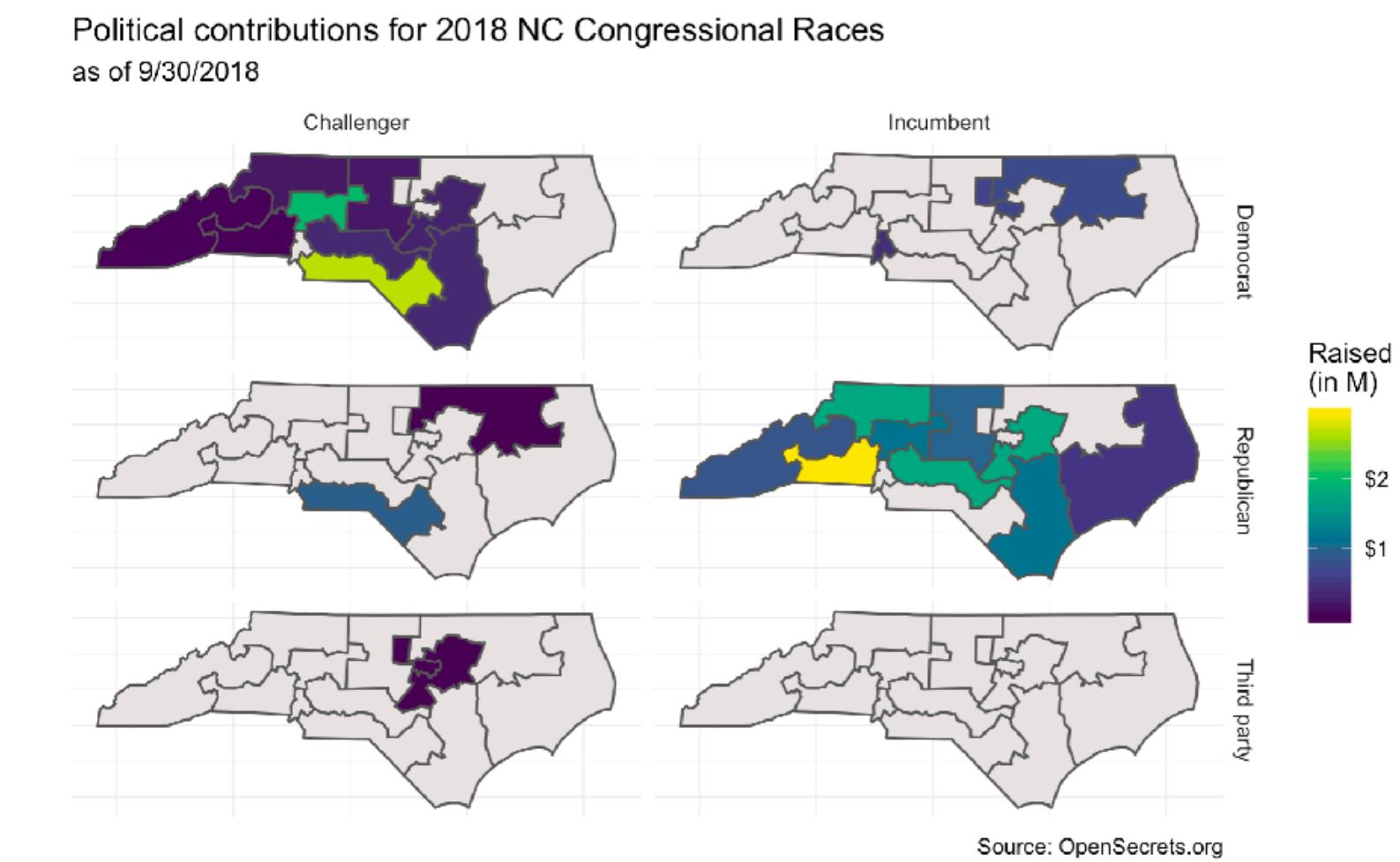
**Ex 1:** Scrape the table off the web and save as a data frame.

Candidate	Raised	Spent	Cash on Hand	Last Report
G K Butterfield (D) • Incumbent	\$714,219	\$797,700	\$560,416	10/17/2018
Roger Allison (R)	\$28,314	\$27,817	\$497	10/17/2018

↓

candidate_info	raised	spent	cash_on_hand	last_report	race
1 G K Butterfield (D) • Incumbent	714219	797700	560416	2018-10-17	North Carolina District 01
2 Roger Allison (R)	28314	27817	497	2018-10-17	North Carolina District 01

**Ex 2:** What other information do we need represented as variables to make this figure?





# Hide the veggies

**Lesson:** Web scraping essentials for turning a structured table into a data frame in R.

**Ex 1:** Scrape the table off the web and save as a data frame.

Candidate	Raised	Spent	Cash on Hand	Last Report
G K Butterfield (D) • Incumbent	\$714,219	\$797,700	\$560,416	10/17/2018
Roger Allison (R)	\$28,314	\$27,817	\$497	10/17/2018

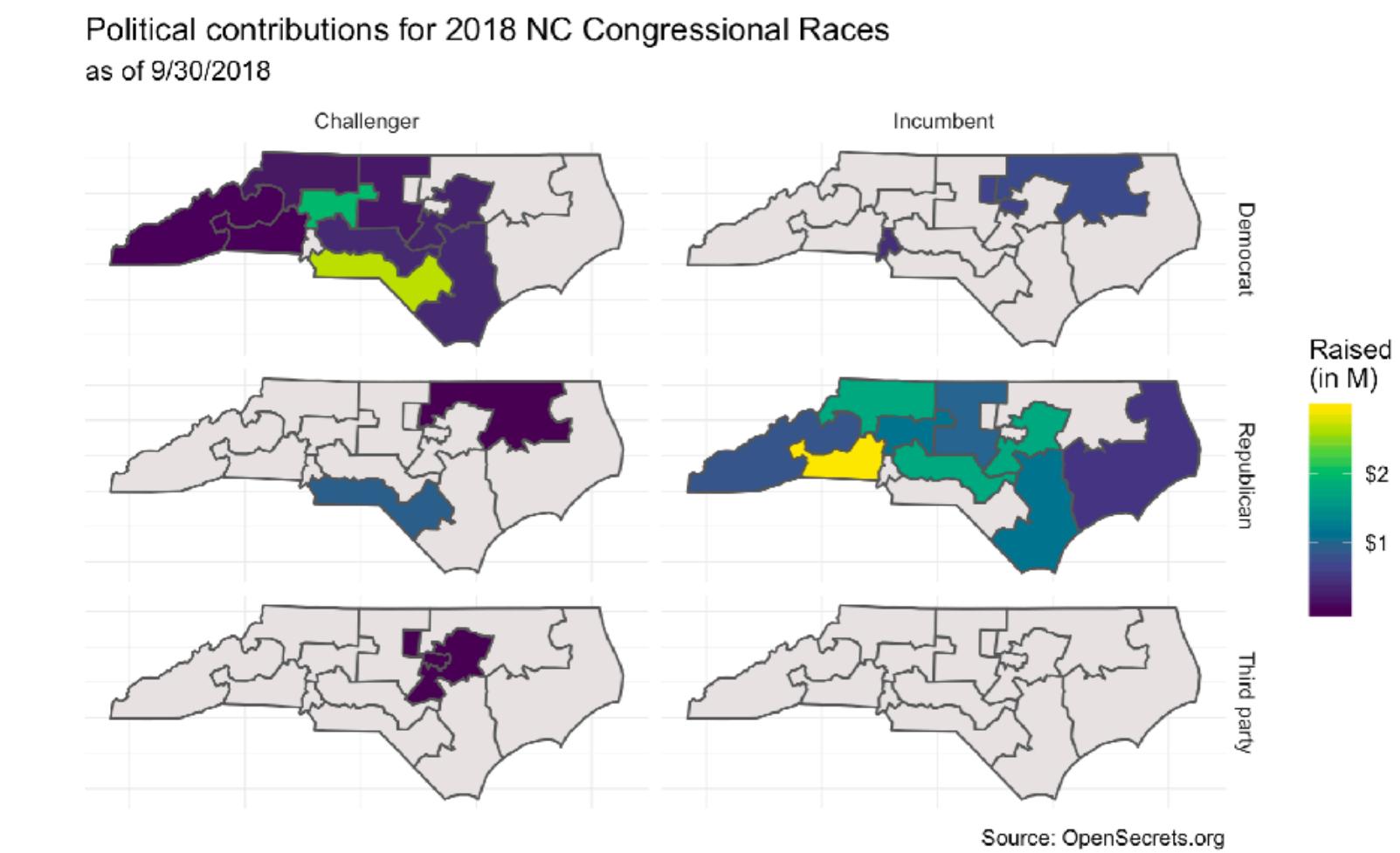
candidate_info	raised	spent	cash_on_hand	last_report	race
1 G K Butterfield (D) • Incumbent	714219	797700	560416	2018-10-17	North Carolina District 01
2 Roger Allison (R)	28314	27817	497	2018-10-17	North Carolina District 01

**Lesson:** “Just enough” regex

candidate_info
1 G K Butterfield (D) • Incumbent
2 Roger Allison (R)

candidate_name	party	status
1 G K Butterfield	Democrat	Incumbent
2 Roger Allison	Republican	Challenger

**Ex 2:** What other information do we need represented as variables to make this figure?



# DESIGN PRINCIPLES



If you are already taking a baking class, which will be easier to venture on to?



# DESIGN PRINCIPLES



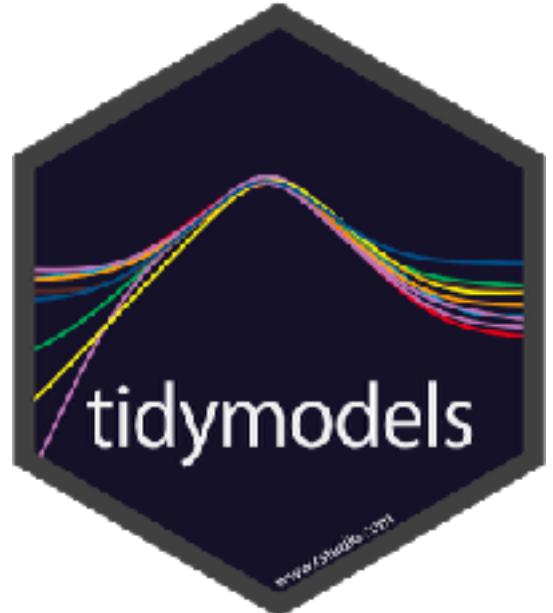
If you are already taking a baking class, which will be easier to venture on to?



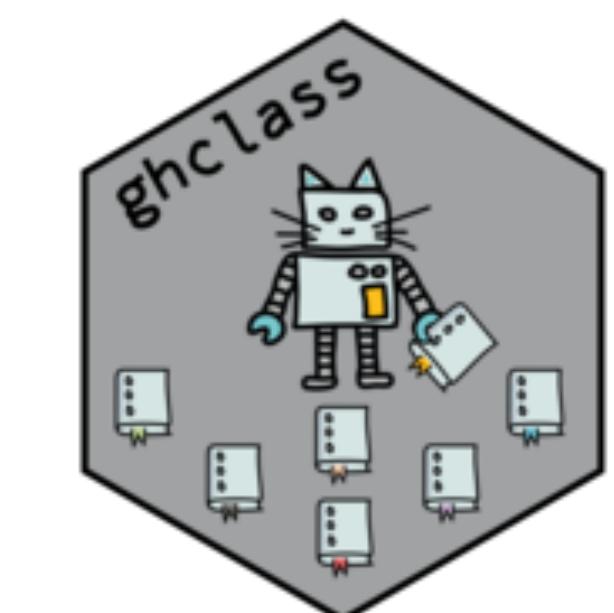


# Leverage the ecosystem

student + instructor



instructor





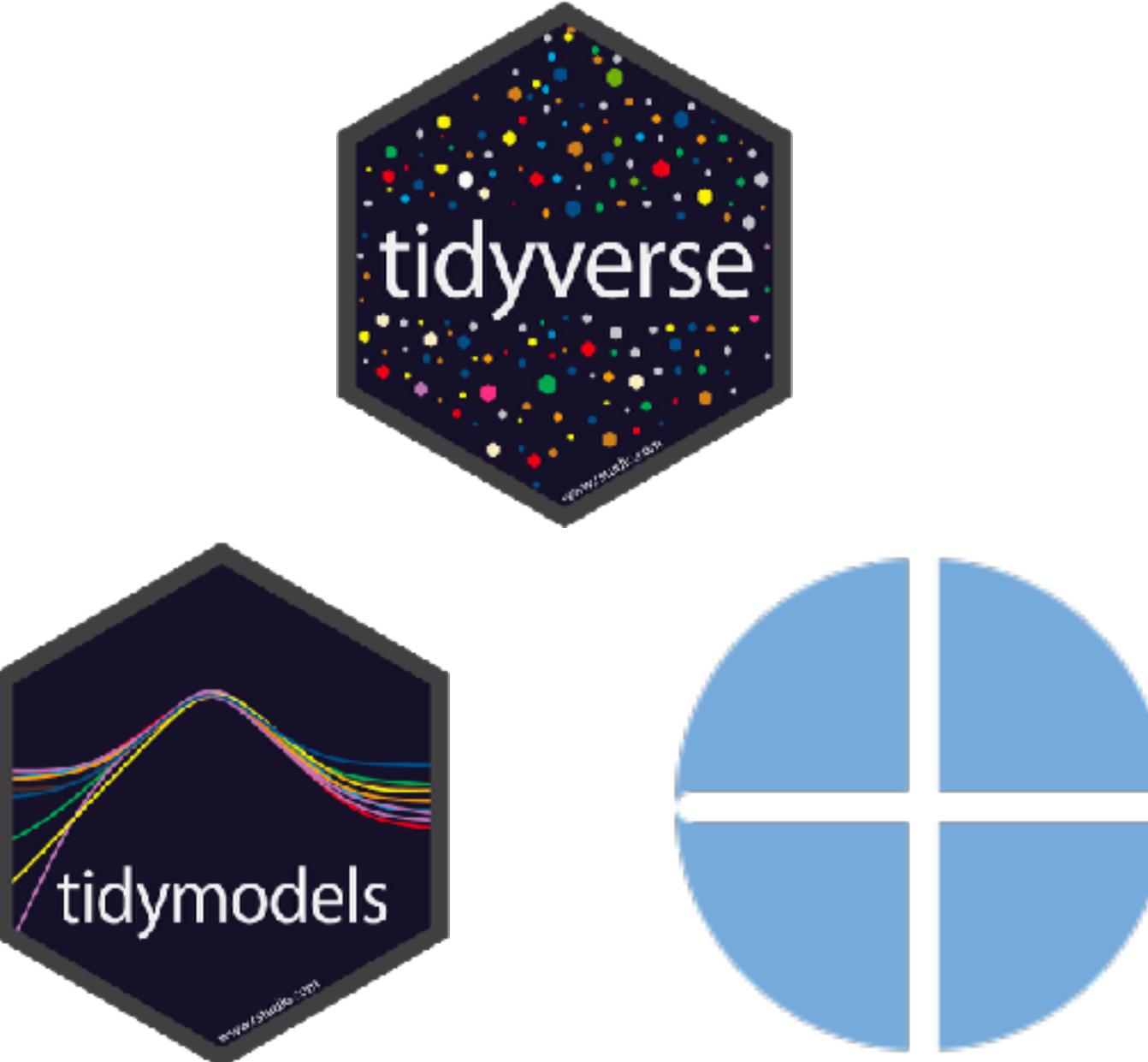
DESIGN PRINCIPLES

VERY NEAR ⚡ FUTURE

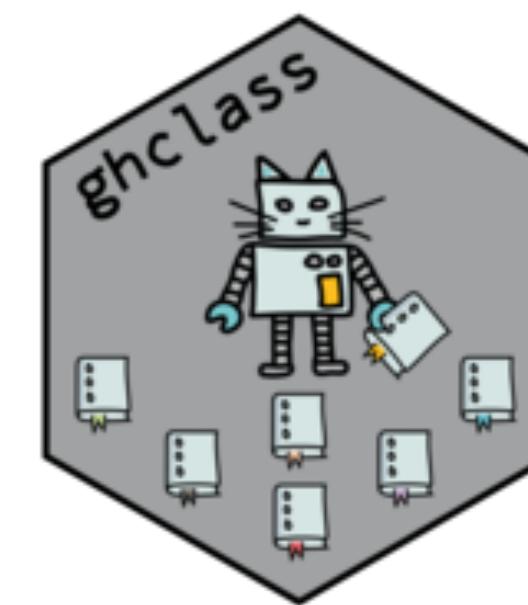


# Leverage the ecosystem

student + instructor



instructor



# USAGE



in full  
to jumpstart /  
overhaul your  
teaching

in bits & pieces  
to supplement  
your teaching



## Attribution-ShareAlike 4.0 International (CC BY-SA 4.0)

This is a human-readable summary of (and not a substitute for) the [license](#). [Disclaimer](#).

### You are free to:

**Share** — copy and redistribute the material in any medium or format

**Adapt** — remix, transform, and build upon the material for any purpose, even commercially.



The licensor cannot revoke these freedoms as long as you follow the license terms.

### Under the following terms:

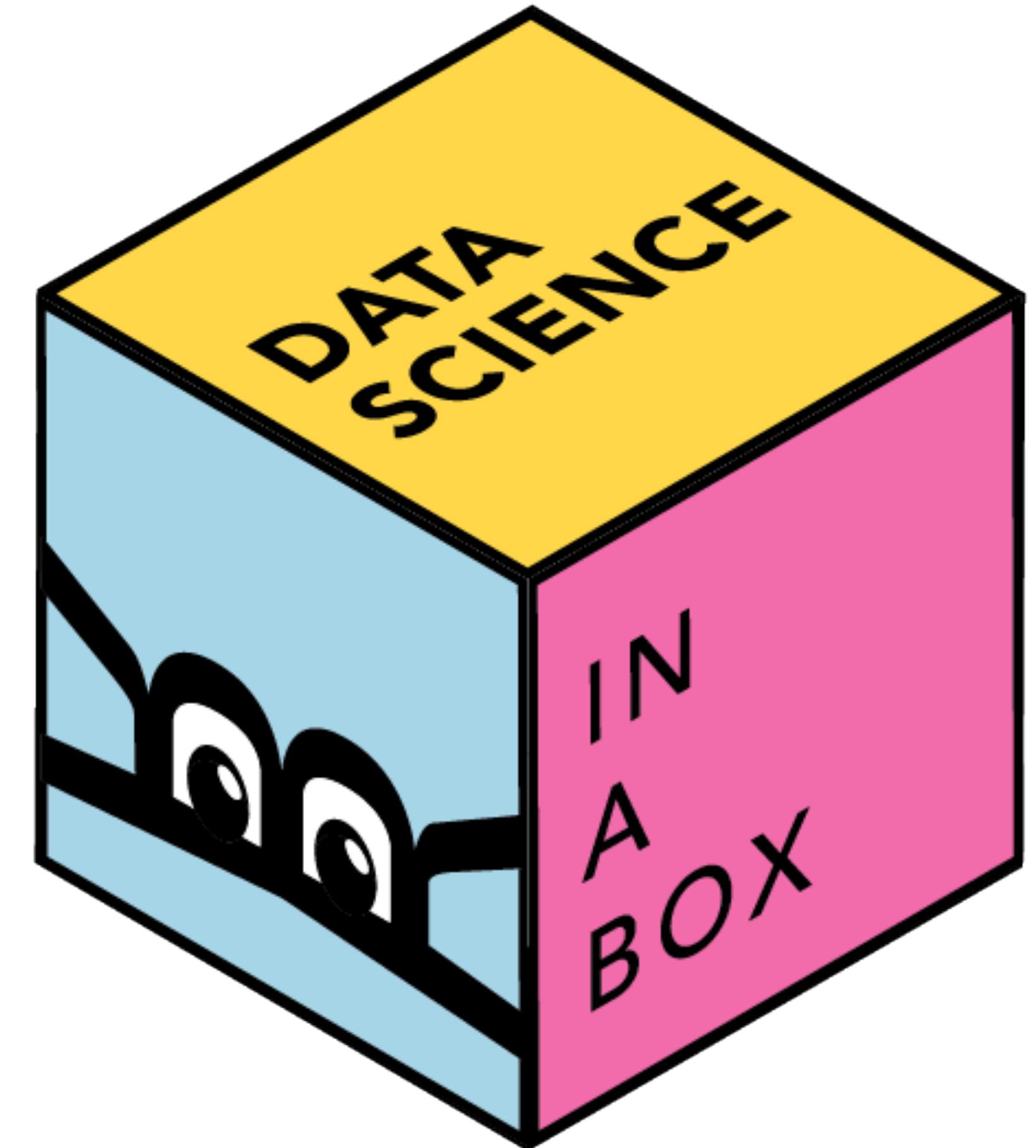


**Attribution** — You must give [appropriate credit](#), provide a link to the license, and [indicate if changes were made](#). You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.



**ShareAlike** — If you remix, transform, or build upon the material, you must distribute your contributions under the [same license](#) as the original.

**No additional restrictions** — You may not apply legal terms or [technological measures](#) that legally restrict others from doing anything the license permits.



[datasciencebox.org](http://datasciencebox.org)



[rstudio-education.github.io/dsbox](https://rstudio-education.github.io/dsbox)



[bit.ly/dsbox-dscwav](http://bit.ly/dsbox-dscwav)

**MINE ÇETINKAYA-RUNDEL**  
DUKE UNIVERSITY + RSTUDIO

@minebocek

mine-cetinkaya-rundel

cetinkaya.mine@gmail.com

