# MINE ÇETINKAYA-RUNDEL

## UNIVERSITY OF EDINBURGH + RSTUDIO

@minebocek

mine-cetinkaya-rundel

cetinkaya.mine@gmail.com

# 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…

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

26 lessons
10 application exercises
10 interactive tutorials
10 labs
10 homework assignments
2 take home exams
1 open-ended project

→

1 website
1 repo
1 package

5    design principles

If you need to bake a cake, which kitchen would you rather get started in?

DESIGN PRINCIPLES

If you need to bake a cake, which kitchen would you rather get started in?

# 🎉 Cherish day one

# When baking a cake, which do you prefer: recipe only or recipe + 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.

# 🍰 Start with cake

▸ Open today's demo project

▸ Knit the document and discuss the results with your neighbor



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

▸ 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

```r
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

```r
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

```r
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"
  )
```

DESIGN PRINCIPLES

# 🍰 Start with cake



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

Which of the two motivates you more to learn how to cook eggs?

DESIGN PRINCIPLES

Which of the two motivates you more to learn how to cook eggs?

# 👶 Skip baby steps

*But drilling through the baby steps can be useful [cite], this can happen outside of class with learnr tutorials (maybe a parson's problem here?)*

Which of the two following is more appetizing to someone who might not be a fan of broccoli?

DATA SCIENCE IN A BOX

Which of the two following is more appetizing to someone who might not be a fan of 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

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

🥦 **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?

Political contributions for 2018 NC Congressional Races as of 9/30/2018



Source: OpenSecrets.org

🥦 **Hide the veggies**

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

**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 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?

Political contributions for 2018 NC Congressional Races as of 9/30/2018



Source: OpenSecrets.org

If you already have ingredients and tools to bake a cake, which of these will be easier to also prepare?

If you already have ingredients and tools to bake a cake, which of these will be easier to also prepare?

# Leverage the ecosystem

1 Use it in full to jumpstart / overhaul your teaching

2 Use it in bits and pieces to supplement your teaching

1  Scalability

‣ More formative assessments via **learnr**

‣ Automated feedback

‣ Peer review

2  Assessment

‣ Curriculum: How are students learning?

‣ Impact: How are these resources being used?

Add link

# MINE ÇETINKAYA-RUNDEL

## UNIVERSITY OF EDINBURGH + RSTUDIO

@minebocek

mine-cetinkaya-rundel

cetinkaya.mine@gmail.com

**validated**

**Retrospective study** of 205 open ended student projects
– on **creativity**, **depth** and the complexity of **multivariate visualizations**
– compared across students who learned R using **base R** syntax vs. **tidyverse**

**validated**

**Creativity:**
1. Creation of new variable(s) based on existing variables
2. Transformation of existing variables
3. Existence of a subgroup analysis
4. Use of a subset of the dataset for all steps of the project

**validated**

**Depth:**

1.  Presence of consistent theme throughout the project
2.  Use of relevant data

**Multivariate visualizations:**

1. Presence of a visualization with 3+ variables
2. Interpretation of the multivariate visualization

validated