## Stat 331 Midterm

```
Working with Data in the Tidyverse:
Dataset %>%
filter(is.na(col))
   + This finds all the values were it is na in that specific col
   + Always check to make sure your variable names are the spelled correctly
# Filter and skim
bakeoff %>%
 filter(!is.na(us_season)) %>%
 skim()
   + Don't need to put skim(us season)
Bakers %>%
count(airedus, series) %>%
mutate(proportion bakers = n/sum(n))
   + Gives you proportion of bakers in each series for the show
Bakers %>%
count(airedus, series) %>%
count(airedus)
   + This will give you number of series that aired in the US and number of series that did not
   + If you didn't add the last count if would just go through each season and tell you which
       one was aired in the US or anywhere else.
Bakers %>%
count(result == "SB")
   + Gives counts for both not SB and SB
parse_date("17 August 2010", format = "%d %B %Y")
   + If there were dashes add them "%d - %B - %Y"
bakers by series %>%
count(baker, sort=TRUE)
   + This will get the baker that appear on the most series.
   + Def need to know sort = TRUE
ggplot(bakeoff, aes(x=episode)) +
  geom_bar() +
  facet wrap(~series)
   + Num of had a bar plot for each series and num of shows idk what it was counting exactly
desserts %>%
  arrange(desc(uk_airdate))
   + Print by descending uk airdate or i think (-) before works as well
desserts <- read_csv("desserts.csv",
            na = c("", "NA", "N/A"),
            col types = cols(technical = col number(),
             uk airdate = col date(format = "%d %B %Y"),
              result = col_factor(levels = NULL)
```

```
))
   + I guess it casts result as a factor
desserts_2 <- desserts %>%
 mutate(nut = recode(nut, "filbert" = "hazelnut",
                "no nut" = NA character,
                .default = "Other nut" ))
   + Dont forget NA_character_
# Recode channel as factor: bbc (1) or not (0)
ratings <- ratings %>%
 mutate(bbc = recode_factor(channel,
                 "Channel 4" = 0,
                 .default = 1)
# Select to look at variables to plot next
ratings %>%
 select(series, channel, bbc, viewer_growth)
# Make a filled bar chart
ggplot(ratings, aes(x = series, y = viewer growth, fill = bbc)) +
 geom_col()
# Drop 7- and 28-day episode ratings
ratings %>%
 select(-ends_with("day"))
# Move channel to front and drop 7-/28-day episode ratings
ratings %>%
 select(channel, everything(), -ends_with("day"))
# Adapt code to drop 28-day columns; keep 7-day in front
viewers_7day <- ratings %>%
  select(viewers_7day_ = ends_with("7day"),
      everything(),
      -ends_with("28day"))
tidy ratings <- ratings %>%
  # Gather and convert episode to factor every col but series
       gather(key = "episode", value = "viewers_7day", -series,
      factor_key = TRUE, na.rm = TRUE)
tidy ratings <- ratings %>%
  # Gather and convert episode to factor
```

```
gather(key = "episode", value = "viewers_7day", -series,
      factor_key = TRUE, na.rm = TRUE) %>%
      # Sort in ascending order by series and episode
  arrange(series, episode) %>%
      # Create new variable using row number()
  mutate(episode_count = row_number())
# Plot viewers by episode and series
ggplot(tidy_ratings, aes(x =episode_count,
         y = viewers_7day,
         fill = series)) +
  geom_col()
Question: Make a line plot to visualize the 7-day viewers by episode, facetted by
series. Since the x-axis is a factor, you'll use the group aesthetic to plot one line
per series
week_ratings <- ratings2 %>%
      # Select 7-day viewer ratings
  select(series, ends with("7day")) %>%
      # Gather 7-day viewers by episode
  gather(episode, viewers_7day, ends_with("7day"), na.rm = TRUE, factor_key = TRUE)
# Plot 7-day viewers by episode and series
ggplot(week_ratings, aes(x = episode,
              y = viewers_7day,
              group = series)) +
  geom line() +
  facet_wrap(~series)
Edit your code for the facetted line plot from the previous exercise to also color
the lines by series. Using guides(), remove the color guide by setting it equal
to FALSE, and add theme_minimal().
# Create week_ratings
week ratings <- ratings2 %>%
  select(series, ends_with("7day")) %>%
  gather(episode, viewers_7day, ends_with("7day"),
      na.rm = TRUE) %>%
```

# Edit your code to color by series and add a theme

mutate(episode = parse\_number(episode))

separate(episode, into = "episode", extra = "drop") %>%

```
ggplot(week_ratings, aes(x = episode, y = viewers 7day,
              group = series, color = series)) +
  geom_line() +
  facet_wrap(~series) +
  guides(color = FALSE) +
  theme_minimal()
ratings3 <- ratings2 %>%
       # Unite and change the separator
       unite(viewers_7day, viewers_millions, viewers_decimal, sep = "") %>%
       # Adapt to cast viewers as a number
       mutate(viewers_7day = parse_number(viewers_7day))
# Print to view
ratings3
Add a line after the %>% to reshape the output of count such that you have 3
columns: series, days_7, and days_28. To do this within <a href="mailto:spread()">spread()</a>, set sep =
"_" as an argument.
# Create tidy data with 7- and 28-day viewers
tidy_ratings_all <- ratings2 %>%
  gather(episode, viewers, ends_with("day"), na.rm = TRUE) %>%
  separate(episode, into = c("episode", "days")) %>%
  mutate(episode = parse_number(episode),
      days = parse_number(days))
tidy ratings all %>%
       # Count viewers by series and days
  count(series, days, wt = viewers) %>%
       # Adapt to spread counted values
  spread(days, n, sep = "_")
```

Fill in the blanks in the provided code to filter rows that contain viewer data for each series' premiere and finale episodes- don't forget to undo the grouping when you are done

Fill in the provided code to make a bar chart with one bar per series, where the height of the bar shows the bump in 7-day viewers from the premiere to the finale episodes. Do any series stand out?

```
# Calculate relative increase in viewers
bump_by_series <- first_last %>%
    spread(episode, viewers) %>%
    mutate(bump = (last - first) / first)

# Fill in to make bar chart of bumps by series
ggplot(bump_by_series, aes(x = series, y = bump)) +
    geom_col()+
    scale y continuous(labels = scales::percent) # converts to %
```

```
# Create skill variable with 3 levels
bakers skill <- bakers %>%
 mutate(skill = case when(
  star baker > technical winner ~ "super star",
  star baker < technical winner ~ "high tech",
  TRUE ~ "well rounded"
 ))
# Filter zeroes to examine skill variable
bakers skill %>%
 filter(star baker == 0 & technical winner == 0) %>%
 count(skill)
# Edit to reverse x-axis order
ggplot(bakers, aes(x = skill, fill = series winner)) +
 geom bar()
   + color= won't get it right always especially when want to split up the bars
# Add a line to extract labeled month
baker_dates_cast <- baker_dates %>%
 mutate(last date appeared us = dmy(last date appeared us),
     last_month_us = month(last_date_appeared_us, label = TRUE))
Label = TRUE allows "Jan", "Feb"
# Make bar chart by last month
ggplot(baker_dates_cast, aes(x=last_month_us)) + geom_bar()
# Add a line to create whole months on air variable
baker time <- baker time %>%
 mutate(time_on_air = interval(first_date_appeared_uk, last_date_appeared_uk),
     weeks on air = time on air / weeks(1),
     months_on_air = time_on_air %/% months(1))
```

```
# Add another mutate to replace "THIRD PLACE" with "RUNNER UP"and count
bakers <- bakers %>%
 mutate(position_reached = str_to_upper(position_reached),
     position_reached = str_replace(position_reached, "-", " "),
   position reached = str replace(position reached, "THIRD PLACE", "RUNNER UP"))
# Count rows
bakers %>%
count(position_reached)
# Summarize by country: by_country
by_country <- votes_processed %>%
 group_by(country) %>%
 summarize(total = n(),
       percent_yes = mean(vote == 1))
# Sort in ascending order of percent_yes
by_country %>%
arrange(percent_yes)
# Now sort in descending order
by_country %>%
arrange(-percent_yes)
# Load the ggplot2 package
library(ggplot2)
# Create line plot
ggplot(by_year, aes(x=year, y=percent_yes)) + geom_line()
# Change to scatter plot and add smoothing curve
ggplot(by_year, aes(year, percent_yes)) + geom_point() + geom_smooth()
```

```
# Start with by year country dataset
by year country <- votes processed %>%
 group_by(year, country) %>%
 summarize(total = n(),
       percent yes = mean(vote == 1))
# Print by_year_country
by_year_country
# Create a filtered version: UK_by_year
UK by year <- by year country %>%
filter(country =="United Kingdom")
# Line plot of percent_yes over time for UK only
ggplot(UK_by_year, aes(x=year,y=percent_yes)) + geom_line()
# Vector of four countries to examine
countries <- c("United States", "United Kingdom",
        "France", "India")
# Filter by_year_country: filtered_4_countries
filtered_4_countries <- by_year_country %>%
filter(country %in% countries)
# Line plot of % yes in four countries
ggplot(filtered_4_countries, aes(x=year,y=percent_yes,color=country)) +
geom_line()
# Line plot of % yes over time faceted by country
ggplot(filtered_6_countries, aes(x=year,y=percent_yes)) +
 geom line() +
 facet_wrap(~ country)
# if we want to scale freely
facet_wrap(~ country , scales = "free_y")
US_fit <- Im(percent_yes ~ year, data = US_by_year)
library(broom)
# Call the tidy() function on the US fit object
```

```
tidy(US fit)
# Linear regression of percent_yes by year for US
US_by_year <- by_year_country %>%
 filter(country == "United States")
US_fit <- Im(percent_yes ~ year, US_by_year)
# Fit model for the United Kingdom
UK by year <- by year country %>%
filter(country == "United Kingdom")
UK_fit <- Im(percent_yes ~ year, UK_by_year)
# Create US_tidied and UK_tidied
US_tidied <- tidy(US_fit)
UK_tidied <- tidy(UK_fit)
# Combine the two tidied models
bind_rows(US_tidied,UK_tidied)
# All countries are nested besides country
nested <- by_year_country %>%
 nest(-country)
# Print the nested data for Brazil
nested$data[[7]]
# Load tidyr and purrr
library(tidyr)
library(purrr)
# Perform a linear regression on each item in the data column
by_year_country %>%
 nest(-country) %>%
 mutate(model = map(data, \sim lm(percent yes \sim year, .)))
# Add another mutate that applies tidy() to each model
mutate(tidied = map(model,tidy))
# Add one more step that unnests the tidied column
 unnest(tidied)
```

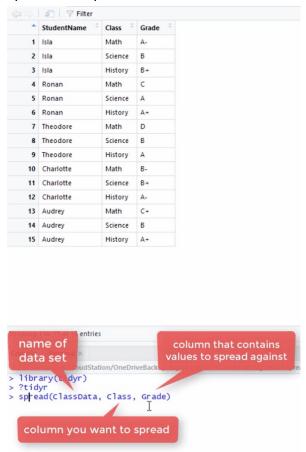
.....

```
# Filter by adjusted p-values
filtered_countries <- country_coefficients %>%
 filter(term == "year") %>%
 mutate(p.adjusted = p.adjust(p.value)) %>%
 filter(p.adjusted < .05)
glimpse(filtered_countries)
# Sort for the countries increasing most quickly
filtered countries %>%
arrange(estimate)
# Sort for the countries decreasing most guickly
filtered countries %>%
arrange(-estimate)
# Load dplyr package
library(dplyr)
# Print the votes_processed dataset
votes processed
# Print the descriptions dataset
descriptions
# Join them together based on the "rcid" and "session" columns
votes_joined <- votes_processed %>%
inner_join(descriptions, by = c("rcid", "session"))
-----
# Filter, then summarize by year: US co by year
US_co_by_year <- votes_joined %>%
 group_by(year) %>%
 filter(co == 1, country == "United States") %>%
 summarize(percent yes = mean(vote == 1))
# Graph the % of "yes" votes over time
ggplot(US_co_by_year, aes(x=year,y=percent_yes)) +
 geom line()
library(tidyr)
# Gather the six me/nu/di/hr/co/ec columns
votes joined %>%
gather(topic,has_topic,me:ec)
# Perform gather again, then filter
votes gathered <- votes joined %>%
gather(topic,has topic,me:ec) %>% filter(has topic == 1)
```

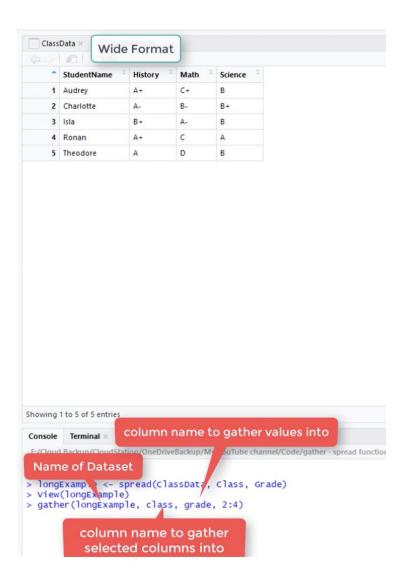
```
# Create vanuatu_by_country_year_topic
vanuatu_by_country_year_topic <- by_country_year_topic %>%
filter(country == "Vanuatu")
```

# Plot of percentage "yes" over time, faceted by topic
ggplot(vanuatu\_by\_country\_year\_topic, aes(year, percent\_yes)) +
 geom\_line() +
 facet\_wrap(~ topic)

## Spread Example



## Gather Example



^	StudentName	class	grade
1	Audrey	History	A+
2	Charlotte	History	A-
3	Isla	History	B+
4	Ronan	History	A+
5	Theodore	History	Α
6	Audrey	Math	C+
7	Charlotte	Math	B-
8	Isla	Math	A-
9	Ronan	Math	С
10	Theodore	Math	D
11	Audrey	Science	В
12	Charlotte	Science	B+
13	Isla	Science	В
14	Ronan	Science	A
15	Theodore	Science	В