

Datawrangling

Samit Kaffle

2025-03-20

load required library

```
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 4.4.2
```

```
## Warning: package 'ggplot2' was built under R version 4.4.2
```

```
## Warning: package 'readr' was built under R version 4.4.2
```

```
## Warning: package 'forcats' was built under R version 4.4.2
```

```
## Warning: package 'lubridate' was built under R version 4.4.2
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
```

```
## v dplyr      1.1.4      v readr      2.1.5
```

```
## v forcats    1.0.0      v stringr    1.5.1
```

```
## v ggplot2    3.5.1      v tibble     3.2.1
```

```
## v lubridate  1.9.3      v tidyr      1.3.1
```

```
## v purrr      1.0.2
```

```
## -- Conflicts ----- tidyverse_conflicts() --
```

```
## x dplyr::filter() masks stats::filter()
```

```
## x dplyr::lag()     masks stats::lag()
```

```
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```
library(dplyr)
```

1. 3 pts. Download two .csv files from Canvas called DiversityData.csv and Metadata.csv, and read them into R using relative file paths.

```
diversitydata=read.csv("DiversityData.csv")
metadata=read.csv("Metadata.csv",na.strings="na")
```

2. 4 pts. Join the two dataframes together by the common column 'Code'. Name the resulting dataframe alpha.

```
alpha = metadata %>%
  full_join(diversitydata, by = "Code")
```

3. 4 pts. Calculate Pielou's evenness index: Pielou's evenness is an ecological parameter calculated by the Shannon diversity index (column Shannon) divided by the log of the richness column.
 - a. Using mutate, create a new column to calculate Pielou's evenness index.
 - b. Name the resulting dataframe alpha_even.

```
alpha_even <- alpha %>%
  mutate(Pielou_evenness = shannon / log(richness))
```

4. Pts. Using tidyverse language of functions and the pipe, use the summarise function and tell me the mean and standard error evenness grouped by crop over time.
 - a. Start with the alpha_even dataframe
 - b. Group the data: group the data by Crop and Time_Point.
 - c. Summarize the data: Calculate the mean, count, standard deviation, and standard error for the even variable within each group.
 - d. Name the resulting dataframe alpha_average

```
alpha_average <- alpha_even %>%
  group_by(Crop, Time_Point) %>%
  summarise(
    mean_evenness = mean(Pielou_evenness, na.rm = TRUE),
    count = n(),
    sd_evenness = sd(Pielou_evenness, na.rm = TRUE),
    se_evenness = sd_evenness / sqrt(count)
  )
```

'summarise()' has grouped output by 'Crop'. You can override using the
'.groups' argument.

5. 4. Pts. Calculate the difference between the soybean column, the soil column, and the difference between the cotton column and the soil column
 - a. Start with the alpha_average dataframe
 - b. Select relevant columns: select the columns Time_Point, Crop, and mean.even.
 - c. Reshape the data: Use the pivot_wider function to transform the data from long to wide format, creating new columns for each Crop with values from mean.even.
 - d. Calculate differences: Create new columns named diff.cotton.even and diff.soybean.even by calculating the difference between Soil and Cotton, and Soil and Soybean, respectively.
 - e. Name the resulting dataframe alpha_average2

```
alpha_average2 <- alpha_average %>%
  select(Time_Point, Crop, mean_evenness) %>%
  pivot_wider(names_from = Crop, values_from = mean_evenness) %>%
  mutate(
    diff_cotton_even = Soil-Cotton,
    diff_soybean_even = Soil-Soybean
  )
```

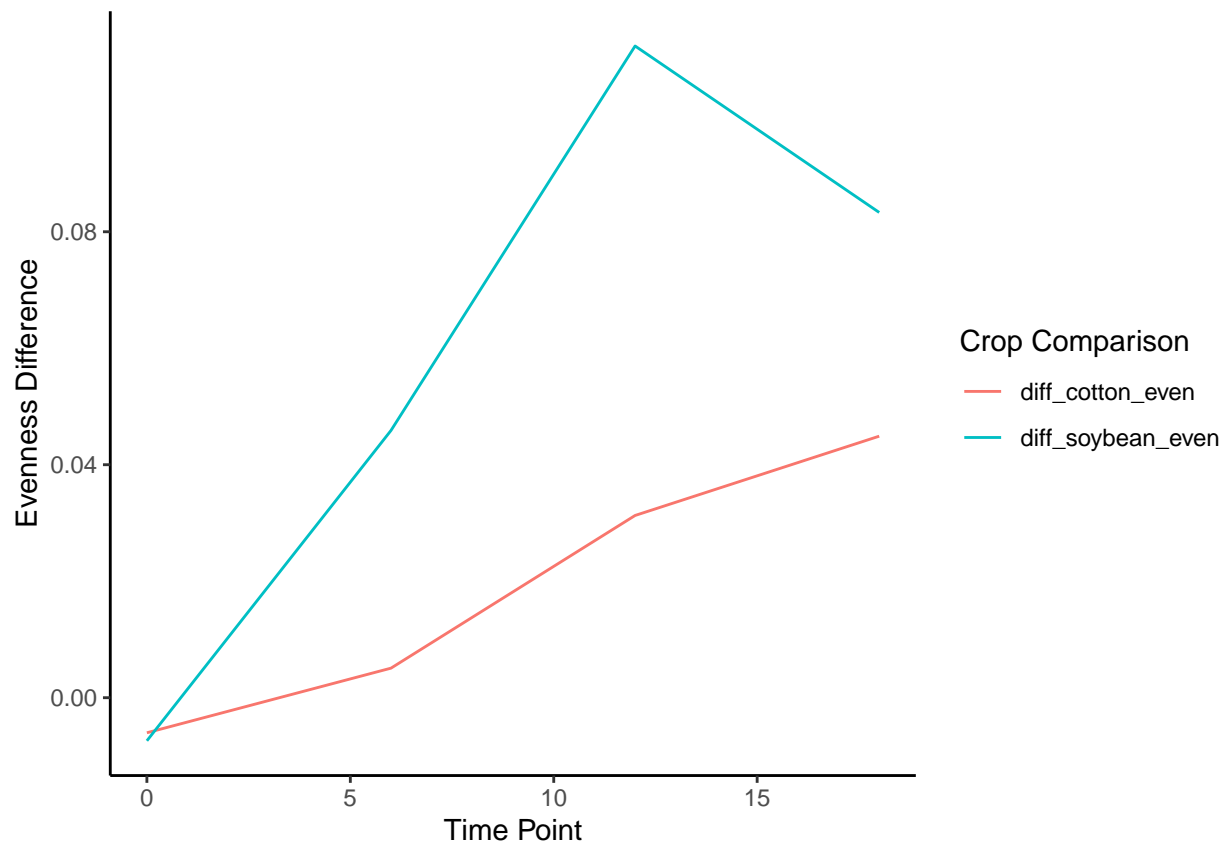
6. 4 pts. Connecting it to plots

- Start with the `alpha_average2` dataframe
- Select relevant columns: select the columns `Time_Point`, `diff.cotton.even`, and `diff.soybean.even`.
- Reshape the data: Use the `pivot_longer` function to transform the data from wide to long format, creating a new column named `diff` that contains the values from `diff.cotton.even` and `diff.soybean.even`.
- This might be challenging, so I'll give you a break. The code is below.

```
pivot_longer(c(diff.cotton.even, diff.soybean.even), names_to = "diff")
```

- Create the plot: Use `ggplot` and `geom_line()` with '`Time_Point`' on the x-axis, the column '`values`' on the y-axis, and different colors for each '`diff`' category. The column named '`values`' come from the `pivot_longer`. The resulting plot should look like the one to the right.

```
alpha_long <- alpha_average2 %>%  
  select(Time_Point, diff_cotton_even, diff_soybean_even) %>%  
  pivot_longer(cols = c(diff_cotton_even, diff_soybean_even), names_to = "diff", values_to = "values")  
  
# Create the plot  
ggplot(alpha_long, aes(x = Time_Point, y = values, color = diff)) +  
  geom_line() +  
  labs(x = "Time Point", y = "Evenness Difference", color = "Crop Comparison") +  
  theme_classic()
```



7. 2 pts. Commit and push a `gfm .md` file to GitHub inside a directory called Coding Challenge 5. Provide me a link to your github written as a clickable link in your `.pdf` or `.docx`

Challenge 5 GitHub Link