

Data wrangling notes

Mamata K C

2025-03-20

Contents

Data wrangling	1
Selecting certain columns	2
Subsetting or filtering data	2
Creating a new column	3
The pipe	4
Summarize data	4
Group by function	5
Connecting to plotting	5
Joining	6
Pivoting	7

Data wrangling

- It is essentially idea of manipulating data for handling the large amount of messy data
- tidyverse package can be installed for data wrangling. It consists of 8 different packages.

```
microbiome.data <- read.csv("Bull_richness.csv") #loading data
str(microbiome.data)
```

```
## 'data.frame':    287 obs. of  16 variables:
## $ SampleID      : chr  "Corn2017LeafObjective2Collection1T1R1CAH2" "Corn2017LeafObjective2Collecti
## $ Crop          : chr  "Corn" "Corn" "Corn" "Corn" ...
## $ Objective     : chr  "Objective 2" "Objective 2" "Objective 2" "Objective 2" ...
## $ Collection    : int   1 1 1 1 1 1 1 1 1 1 ...
## $ Compartment   : chr  "Leaf" "Leaf" "Leaf" "Leaf" ...
## $ DateSampled   : chr  "6/26/17" "6/26/17" "6/26/17" "6/26/17" ...
## $ GrowthStage   : chr  "V6" "V6" "V6" "V6" ...
## $ Treatment     : chr  "Conv." "Conv." "Conv." "Conv." ...
## $ Rep           : chr  "R1" "R1" "R1" "R1" ...
## $ Sample        : chr  "A" "B" "C" "A" ...
## $ Fungicide     : chr  "C" "C" "C" "F" ...
## $ Target_organism: chr  "Fungi" "Fungi" "Fungi" "Fungi" ...
```

```
## $ Location      : chr "Kellogg Biological Station" "Kellogg Biological Station" "Kellogg Biological Station" ...
## $ Experiment    : chr "LTER" "LTER" "LTER" "LTER" ...
## $ Year          : int 2017 2017 2017 2017 2017 2017 2017 2017 2017 2017 ...
## $ richness      : int 9 6 5 7 4 2 3 8 4 4 ...
```

Selecting certain columns

```
library(tidyverse)
```

```
## -- 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
```

```
microbiome.data1 <- select(microbiome.data, SampleID, Crop, Compartment:Fungicide, richness) #using select
```

Subsetting or filtering data

- selecting certain rows

```
head(filter(microbiome.data1, Treatment == "Conv. ")) #selecting the rows only with conventional treatment
```

```
##                               SampleID Crop Compartment DateSampled
## 1 Corn2017LeafObjective2Collection1T1R1CAH2 Corn      Leaf      6/26/17
## 2 Corn2017LeafObjective2Collection1T1R1CBA3 Corn      Leaf      6/26/17
## 3 Corn2017LeafObjective2Collection1T1R1CCB3 Corn      Leaf      6/26/17
## 4 Corn2017LeafObjective2Collection1T1R1FAC3 Corn      Leaf      6/26/17
## 5 Corn2017LeafObjective2Collection1T1R1FBD3 Corn      Leaf      6/26/17
## 6 Corn2017LeafObjective2Collection1T1R1FCE3 Corn      Leaf      6/26/17
##   GrowthStage Treatment Rep Sample Fungicide richness
## 1          V6      Conv.  R1      A          C          9
## 2          V6      Conv.  R1      B          C          6
## 3          V6      Conv.  R1      C          C          5
## 4          V6      Conv.  R1      A          F          7
## 5          V6      Conv.  R1      B          F          4
## 6          V6      Conv.  R1      C          F          2
```

```
head(filter(microbiome.data1, Treatment == "Conv." & Fungicide == "C")) #selecting the rows only with conventional treatment and fungicide C
```

```
##                               SampleID Crop Compartment DateSampled
## 1 Corn2017LeafObjective2Collection1T1R1CAH2 Corn      Leaf      6/26/17
## 2 Corn2017LeafObjective2Collection1T1R1CBA3 Corn      Leaf      6/26/17
## 3 Corn2017LeafObjective2Collection1T1R1CCB3 Corn      Leaf      6/26/17
```

```
## 4 Corn2017LeafObjective2Collection1T1R2CAF3 Corn Leaf 6/26/17
## 5 Corn2017LeafObjective2Collection1T1R2CBG3 Corn Leaf 6/26/17
## 6 Corn2017LeafObjective2Collection1T1R2CCH3 Corn Leaf 6/26/17
## GrowthStage Treatment Rep Sample Fungicide richness
## 1 V6 Conv. R1 A C 9
## 2 V6 Conv. R1 B C 6
## 3 V6 Conv. R1 C C 5
## 4 V6 Conv. R2 A C 3
## 5 V6 Conv. R2 B C 8
## 6 V6 Conv. R2 C C 4
```

```
head(filter(microbiome.data1, Sample == "A" | Sample == "B")) #selectiong rows with sample A or sample B
```

```
## SampleID Crop Compartment DateSampled
## 1 Corn2017LeafObjective2Collection1T1R1CAH2 Corn Leaf 6/26/17
## 2 Corn2017LeafObjective2Collection1T1R1CBA3 Corn Leaf 6/26/17
## 3 Corn2017LeafObjective2Collection1T1R1FAC3 Corn Leaf 6/26/17
## 4 Corn2017LeafObjective2Collection1T1R1FBD3 Corn Leaf 6/26/17
## 5 Corn2017LeafObjective2Collection1T1R2CAF3 Corn Leaf 6/26/17
## 6 Corn2017LeafObjective2Collection1T1R2CBG3 Corn Leaf 6/26/17
## GrowthStage Treatment Rep Sample Fungicide richness
## 1 V6 Conv. R1 A C 9
## 2 V6 Conv. R1 B C 6
## 3 V6 Conv. R1 A F 7
## 4 V6 Conv. R1 B F 4
## 5 V6 Conv. R2 A C 3
## 6 V6 Conv. R2 B C 8
```

Creating a new column

```
microbiome.data1$logRich <- log(microbiome.data1$richness) #previous method

#using mutate function to create new columns
head(mutate(microbiome.data1, logRich = log(richness))) #creating new column called logRich
```

```
## SampleID Crop Compartment DateSampled
## 1 Corn2017LeafObjective2Collection1T1R1CAH2 Corn Leaf 6/26/17
## 2 Corn2017LeafObjective2Collection1T1R1CBA3 Corn Leaf 6/26/17
## 3 Corn2017LeafObjective2Collection1T1R1CCB3 Corn Leaf 6/26/17
## 4 Corn2017LeafObjective2Collection1T1R1FAC3 Corn Leaf 6/26/17
## 5 Corn2017LeafObjective2Collection1T1R1FBD3 Corn Leaf 6/26/17
## 6 Corn2017LeafObjective2Collection1T1R1FCE3 Corn Leaf 6/26/17
## GrowthStage Treatment Rep Sample Fungicide richness logRich
## 1 V6 Conv. R1 A C 9 2.1972246
## 2 V6 Conv. R1 B C 6 1.7917595
## 3 V6 Conv. R1 C C 5 1.6094379
## 4 V6 Conv. R1 A F 7 1.9459101
## 5 V6 Conv. R1 B F 4 1.3862944
## 6 V6 Conv. R1 C F 2 0.6931472
```

```
head(mutate(microbiome.data1, Crop_Treatment = paste(Crop, Treatment))) #creating new column combining
```

```
##                               SampleID Crop Compartment DateSampled
## 1 Corn2017LeafObjective2Collection1T1R1CAH2 Corn      Leaf      6/26/17
## 2 Corn2017LeafObjective2Collection1T1R1CBA3 Corn      Leaf      6/26/17
## 3 Corn2017LeafObjective2Collection1T1R1CCB3 Corn      Leaf      6/26/17
## 4 Corn2017LeafObjective2Collection1T1R1FAC3 Corn      Leaf      6/26/17
## 5 Corn2017LeafObjective2Collection1T1R1FBD3 Corn      Leaf      6/26/17
## 6 Corn2017LeafObjective2Collection1T1R1FCE3 Corn      Leaf      6/26/17
##   GrowthStage Treatment Rep Sample Fungicide richness    logRich Crop_Treatment
## 1          V6      Conv. R1      A          C          9 2.1972246      Corn Conv.
## 2          V6      Conv. R1      B          C          6 1.7917595      Corn Conv.
## 3          V6      Conv. R1      C          C          5 1.6094379      Corn Conv.
## 4          V6      Conv. R1      A          F          7 1.9459101      Corn Conv.
## 5          V6      Conv. R1      B          F          4 1.3862944      Corn Conv.
## 6          V6      Conv. R1      C          F          2 0.6931472      Corn Conv.
```

The pipe

- allows us to combine the output from one function into the input of another function

```
microbiome.data %>%
  select(SampleID, Crop, Compartment:Fungicide, richness) %>% #selecting columns
  filter(Treatment == "Conv.") %>% #including only conventional treatment
  mutate(logRich = log(richness)) %>% #creating new column logRich
  head()
```

```
##                               SampleID Crop Compartment DateSampled
## 1 Corn2017LeafObjective2Collection1T1R1CAH2 Corn      Leaf      6/26/17
## 2 Corn2017LeafObjective2Collection1T1R1CBA3 Corn      Leaf      6/26/17
## 3 Corn2017LeafObjective2Collection1T1R1CCB3 Corn      Leaf      6/26/17
## 4 Corn2017LeafObjective2Collection1T1R1FAC3 Corn      Leaf      6/26/17
## 5 Corn2017LeafObjective2Collection1T1R1FBD3 Corn      Leaf      6/26/17
## 6 Corn2017LeafObjective2Collection1T1R1FCE3 Corn      Leaf      6/26/17
##   GrowthStage Treatment Rep Sample Fungicide richness    logRich
## 1          V6      Conv. R1      A          C          9 2.1972246
## 2          V6      Conv. R1      B          C          6 1.7917595
## 3          V6      Conv. R1      C          C          5 1.6094379
## 4          V6      Conv. R1      A          F          7 1.9459101
## 5          V6      Conv. R1      B          F          4 1.3862944
## 6          V6      Conv. R1      C          F          2 0.6931472
```

Summarize data

```
microbiome.data %>%
  select(SampleID, Crop, Compartment:Fungicide, richness) %>% #selecting columns
  filter(Treatment == "Conv.") %>% #including only conventional treatment
  mutate(logRich = log(richness)) %>% #creating new column logRich
  summarise(Mean.rich=mean(logRich)) #calculating overall mean
```

```
## Mean.rich
## 1 2.304395
```

```
# connect multiple summary statistics
```

```
microbiome.data %>%
  select(SampleID, Crop, Compartment:Fungicide, richness) %>% # selecting columns
  filter(Treatment == "Conv.") %>% # subsetting to only include the conventional treatment
  mutate(logRich = log(richness)) %>% # creating a new column of the log richness
  summarise(Mean.rich = mean(logRich), # calculating the mean richness, stdeviation, and standard error
            n = n(),
            sd.dev = sd(logRich)) %>%
  mutate(std.err = sd.dev/sqrt(n))
```

```
## Mean.rich n sd.dev std.err
## 1 2.304395 144 0.7024667 0.0585389
```

Group by function

```
microbiome.data %>%
  select(SampleID, Crop, Compartment:Fungicide, richness) %>% # selecting columns
  group_by(Treatment, Fungicide) %>% # grouping by treatment and fungicide to later calculate summary s
  mutate(logRich = log(richness)) %>% # creating a new column of the log richness
  summarise(Mean.rich = mean(logRich), # calculating the mean richness, stdeviation, and standard error
            n = n(),
            sd.dev = sd(logRich)) %>%
  mutate(std.err = sd.dev/sqrt(n))
```

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

```
## # A tibble: 4 x 6
## # Groups: Treatment [2]
## Treatment Fungicide Mean.rich n sd.dev std.err
## <chr> <chr> <dbl> <int> <dbl> <dbl>
## 1 Conv. C 2.53 72 0.635 0.0748
## 2 Conv. F 2.07 72 0.696 0.0820
## 3 No-till C 2.63 72 0.513 0.0604
## 4 No-till F 2.36 71 0.680 0.0807
```

Connecting to plotting

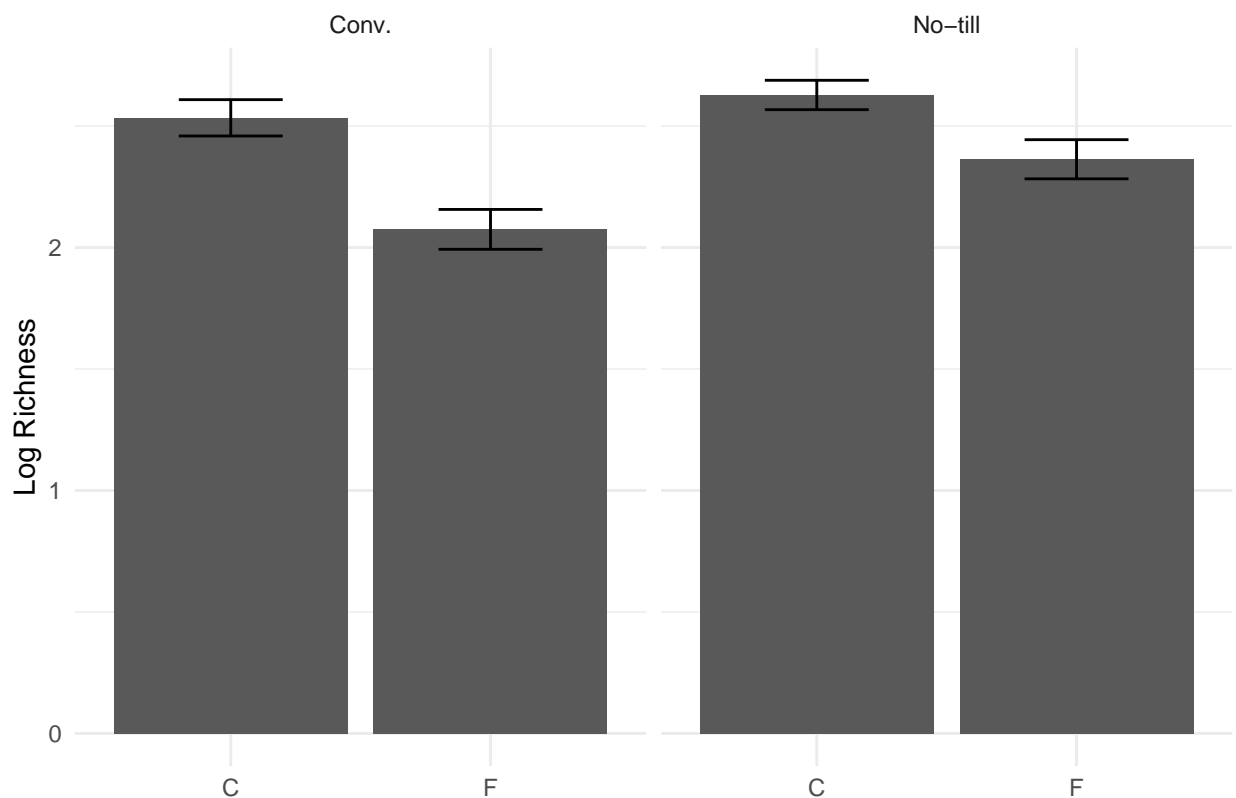
```
microbiome.data %>%
  select(SampleID, Crop, Compartment:Fungicide, richness) %>% # selecting columns
  group_by(Treatment, Fungicide) %>% # grouping by treatment and fungicide to later calculate summary s
  mutate(logRich = log(richness)) %>% # creating a new column of the log richness
  summarise(Mean.rich = mean(logRich), # calculating the mean richness, stdeviation, and standard error
            n = n(),
```

```

sd.dev = sd(logRich)) %>%
mutate(std.err = sd.dev/sqrt(n)) %>%
ggplot(aes(x = Fungicide, y = Mean.rich)) + # adding in a ggplot
geom_bar(stat="identity") +
geom_errorbar(aes(x=Fungicide, ymin=Mean.rich-std.err, ymax=Mean.rich+std.err), width=0.4) +
theme_minimal() +
xlab("") +
ylab("Log Richness") +
facet_wrap(~Treatment)

```

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



Joining

- allows us to combine multiple datasets based on common set of variables.
- it is important because sometimes we need to split our data to run different functions and again need to combine them to form a metadata.

```

# selecting just the richness and sample ID
richness <- microbiome.data %>%
select(SampleID, richness)

```

```
# selecting columns that don't include the richness
metadata <- microbiome.data %>%
  select(SampleID, Fungicide, Crop, Compartment, GrowthStage, Treatment, Rep, Sample)

head(metadata)
```

```
##                               SampleID Fungicide Crop Compartment
## 1 Corn2017LeafObjective2Collection1T1R1CAH2      C Corn      Leaf
## 2 Corn2017LeafObjective2Collection1T1R1CBA3      C Corn      Leaf
## 3 Corn2017LeafObjective2Collection1T1R1CCB3      C Corn      Leaf
## 4 Corn2017LeafObjective2Collection1T1R1FAC3      F Corn      Leaf
## 5 Corn2017LeafObjective2Collection1T1R1FBD3      F Corn      Leaf
## 6 Corn2017LeafObjective2Collection1T1R1FCE3      F Corn      Leaf
##   GrowthStage Treatment Rep Sample
## 1          V6      Conv. R1      A
## 2          V6      Conv. R1      B
## 3          V6      Conv. R1      C
## 4          V6      Conv. R1      A
## 5          V6      Conv. R1      B
## 6          V6      Conv. R1      C
```

```
head(left_join(metadata, richness, by = "SampleID")) # adding the richness data to the metadata based on SampleID
```

```
##                               SampleID Fungicide Crop Compartment
## 1 Corn2017LeafObjective2Collection1T1R1CAH2      C Corn      Leaf
## 2 Corn2017LeafObjective2Collection1T1R1CBA3      C Corn      Leaf
## 3 Corn2017LeafObjective2Collection1T1R1CCB3      C Corn      Leaf
## 4 Corn2017LeafObjective2Collection1T1R1FAC3      F Corn      Leaf
## 5 Corn2017LeafObjective2Collection1T1R1FBD3      F Corn      Leaf
## 6 Corn2017LeafObjective2Collection1T1R1FCE3      F Corn      Leaf
##   GrowthStage Treatment Rep Sample richness
## 1          V6      Conv. R1      A          9
## 2          V6      Conv. R1      B          6
## 3          V6      Conv. R1      C          5
## 4          V6      Conv. R1      A          7
## 5          V6      Conv. R1      B          4
## 6          V6      Conv. R1      C          2
```

Pivoting

```
microbiome.data %>%
  select(SampleID, Crop, Compartment:Fungicide, richness) %>% # selecting columns filter(Class == "Species")
  group_by(Treatment, Fungicide) %>% # grouping by treatment and fungicide to later calculate summary statistics
  summarise(Mean = mean(richness)) %>% # calculates the mean per Treatment and Fungicide
  pivot_wider(names_from = Fungicide, values_from = Mean) %>% # pivot to wide format
  mutate(diff.fungicide = C - F) # calculate the difference between the means.
```

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

```
## # A tibble: 2 x 4
## # Groups:   Treatment [2]
##   Treatment      C      F diff.fungicide
##   <chr>      <dbl> <dbl>      <dbl>
## 1 Conv.        14.6  9.75         4.89
## 2 No-till       15.4 13.1         2.32
```

```
microbiome.data %>%
  select(SampleID, Crop, Compartment:Fungicide, richness) %>% # selecting columns filter(Class == "Soil")
  group_by(Treatment, Fungicide) %>% # grouping by treatment and fungicide to later calculate summary statistics
  summarise(Mean = mean(richness)) %>% # calculates the mean per Treatment and Fungicide
  pivot_wider(names_from = Fungicide, values_from = Mean) %>% # pivot to wide format
  mutate(diff.fungicide = C - F) %>% # calculate the difference between the means.
  ggplot(aes(x = Treatment, y = diff.fungicide)) + # Plot it
  geom_col() +
  theme_minimal() +
  xlab("") +
  ylab("Difference in average species richness")
```

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

