

Assignment 3: Data Exploration

Danielle Butler

Fall 2024

OVERVIEW

This exercise accompanies the lessons in Environmental Data Analytics on Data Exploration.

Directions

1. Rename this file `<FirstLast>_A03_DataExploration.Rmd` (replacing `<FirstLast>` with your first and last name).
2. Change “Student Name” on line 3 (above) with your name.
3. Work through the steps, **creating code and output** that fulfill each instruction.
4. Assign a useful **name to each code chunk** and include ample **comments** with your code.
5. Be sure to **answer the questions** in this assignment document.
6. When you have completed the assignment, **Knit** the text and code into a single PDF file.
7. After Knitting, submit the completed exercise (PDF file) to the dropbox in Canvas.

TIP: If your code extends past the page when knit, tidy your code by manually inserting line breaks.

TIP: If your code fails to knit, check that no `install.packages()` or `View()` commands exist in your code.

Set up your R session

1. Load necessary packages (tidyverse, lubridate, here), check your current working directory and upload two datasets: the ECOTOX neonicotinoid dataset (`ECOTOX_Neonicotinoids_Insects_raw.csv`) and the Niwot Ridge NEON dataset for litter and woody debris (`NEON_NIWO_Litter_massdata_2018-08_raw.csv`). Name these datasets “Neonics” and “Litter”, respectively. Be sure to include the sub-command to read strings in as factors.

```
#Import packages
library(tidyverse); library(lubridate); library(here)
```

```
## -- 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
## here() starts at /home/guest/EDA_Spring2025
```

```
#Check workspace
here()
```

```
## [1] "/home/guest/EDA_Spring2025"
```

```
#Import data
Neonics <- read.csv(file = here('Data',
                                'Raw', 'ECOTOX_Neonicotinoids_Insects_raw.csv'),
                    stringsAsFactors = TRUE)

Litter <- read.csv(file = here('Data',
                                'Raw', 'NEON_NIWO_Litter_massdata_2018-08_raw.csv'),
                    stringsAsFactors = TRUE)
```

Learn about your system

2. The neonicotinoid dataset was collected from the Environmental Protection Agency's ECOTOX Knowledgebase, a database for ecotoxicology research. Neonicotinoids are a class of insecticides used widely in agriculture. The dataset that has been pulled includes all studies published on insects. Why might we be interested in the ecotoxicology of neonicotinoids on insects? Feel free to do a brief internet search if you feel you need more background information.

Answer: Neonicotinoids are insecticides that are highly toxic to insects, especially pollinators, and can have devastating ecological impacts. They are the most widely used class of insecticides in the world.

3. The Niwot Ridge litter and woody debris dataset was collected from the National Ecological Observatory Network, which collectively includes 81 aquatic and terrestrial sites across 20 ecoclimatic domains. 32 of these sites sample forest litter and woody debris, and we will focus on the Niwot Ridge long-term ecological research (LTER) station in Colorado. Why might we be interested in studying litter and woody debris that falls to the ground in forests? Feel free to do a brief internet search if you feel you need more background information.

Answer: Studying litter and woody debris on the forest floor is important because it plays a crucial role in nutrient cycling, provides habitat for various organisms, influences soil structure, and contributes significantly to the overall health and biodiversity of the forest ecosystem; essentially acting as a key component in the forest's carbon cycle and regulating water flow within the environment.

4. How is litter and woody debris sampled as part of the NEON network? Read the NEON_Litterfall_UserGuide.pdf document to learn more. List three pieces of salient information about the sampling methods here:

Answer: 1. Litter and fine woody debris are collected from elevated and ground traps, respectively. 2. Mass data is measured to the accuracy of 0.01grams for the following functional groups: leaves, needles, twigs/branches, woody material, seed, flowers and other non-woody reproductive structures, other and mixed. 3. Ground traps are sampled once per year.

Obtain basic summaries of your data (Neonics)

5. What are the dimensions of the dataset?

```
print(dim(Neonics))
```

```
## [1] 4623 30
```

```
print(dim(Litter))
```

```
## [1] 188 19
```

6. Using the `summary` function on the “Effect” column, determine the most common effects that are studied. Why might these effects specifically be of interest? [Tip: The `sort()` command is useful for listing the values in order of magnitude...]

```
sort(summary(Neonics$Effect))
```

```
##      Hormone(s)      Histology      Physiology      Cell(s)
##           1           5           7           9
##      Biochemistry      Accumulation      Intoxication      Immunological
##           11           12           12           16
##      Morphology      Growth      Enzyme(s)      Genetics
##           22           38           62           82
##      Avoidance      Development      Reproduction      Feeding behavior
##           102          136          197          255
##      Behavior      Mortality      Population
##           360          1493          1803
```

Answer: The most common effects that are studied are “Population” - This would be most interesting to study to see how the insecticides affect the population of the pollinators.

7. Using the `summary` function, determine the six most commonly studied species in the dataset (common name). What do these species have in common, and why might they be of interest over other insects? Feel free to do a brief internet search for more information if needed. [TIP: Explore the help on the `summary()` function, in particular the `maxsum` argument...]

```
sort(summary(Neonics$Species.Common.Name, maxsum = 6))
```

```
##      Bumble Bee      Carniolan Honey Bee      Buff Tailed Bumblebee
##           140           152           183
##      Parasitic Wasp      Honey Bee      (Other)
##           285           667           3196
```

Answer: Bumble Bee, Honey Bee, Parasitic Wasp and Carniolan Honey Bee Buff Tailed Bumblebee and (Other). Bees are extremely important because they pollinate plants, which helps produce food and crops and maintain the health of the ecosystem. According to the UN, 1/3 of the world’s food production depends on bees.

8. Concentrations are always a numeric value. What is the class of `Conc.1..Author.` column in the dataset, and why is it not numeric? [Tip: Viewing the dataframe may be helpful...]

```
class(Neonics$Conc.1..Author.)
```

```
## [1] "factor"
```

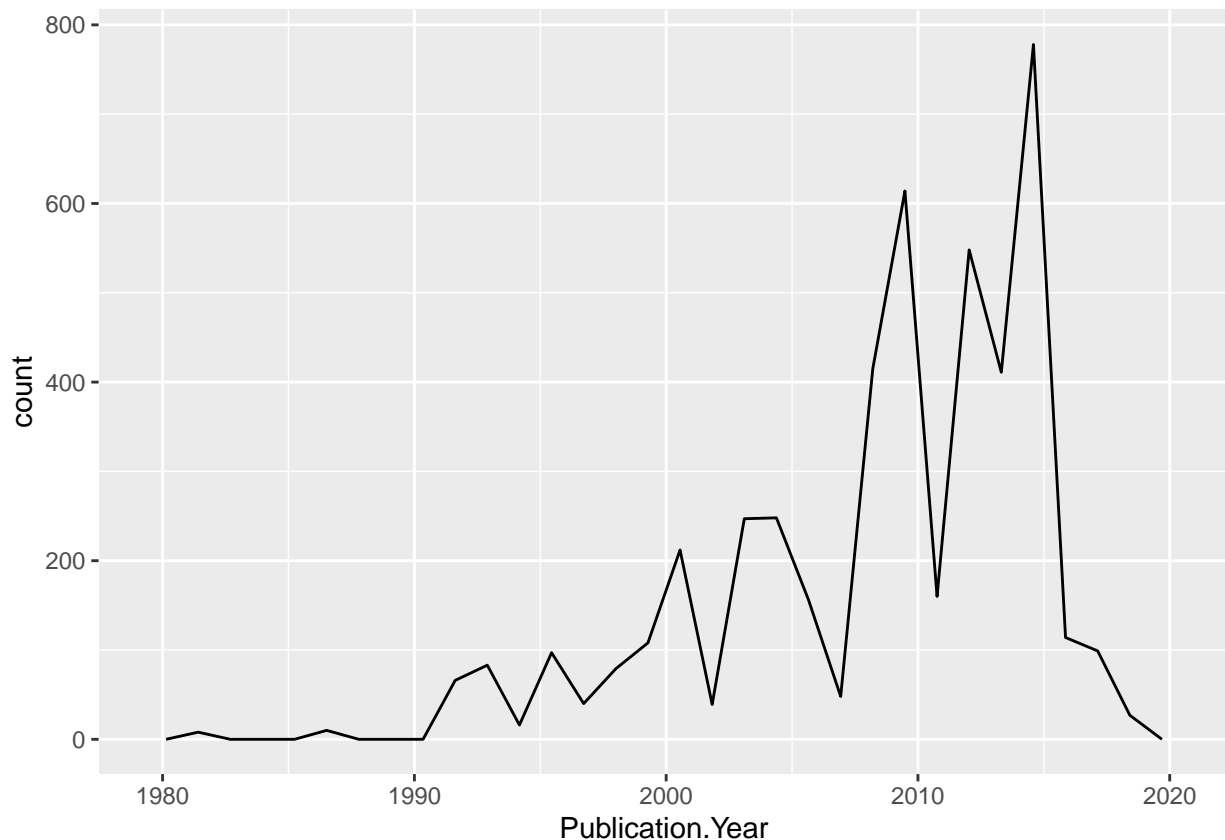
Answer: The class is “factor” - the column data is not simply numeric - there are slashes and other characters. If we want to be able to characterize by them, this class makes sense.

Explore your data graphically (Neonics)

9. Using `geom_freqpoly`, generate a plot of the number of studies conducted by publication year.

```
ggplot(Neonics,aes(x=Publication.Year)) +  
  geom_freqpoly()
```

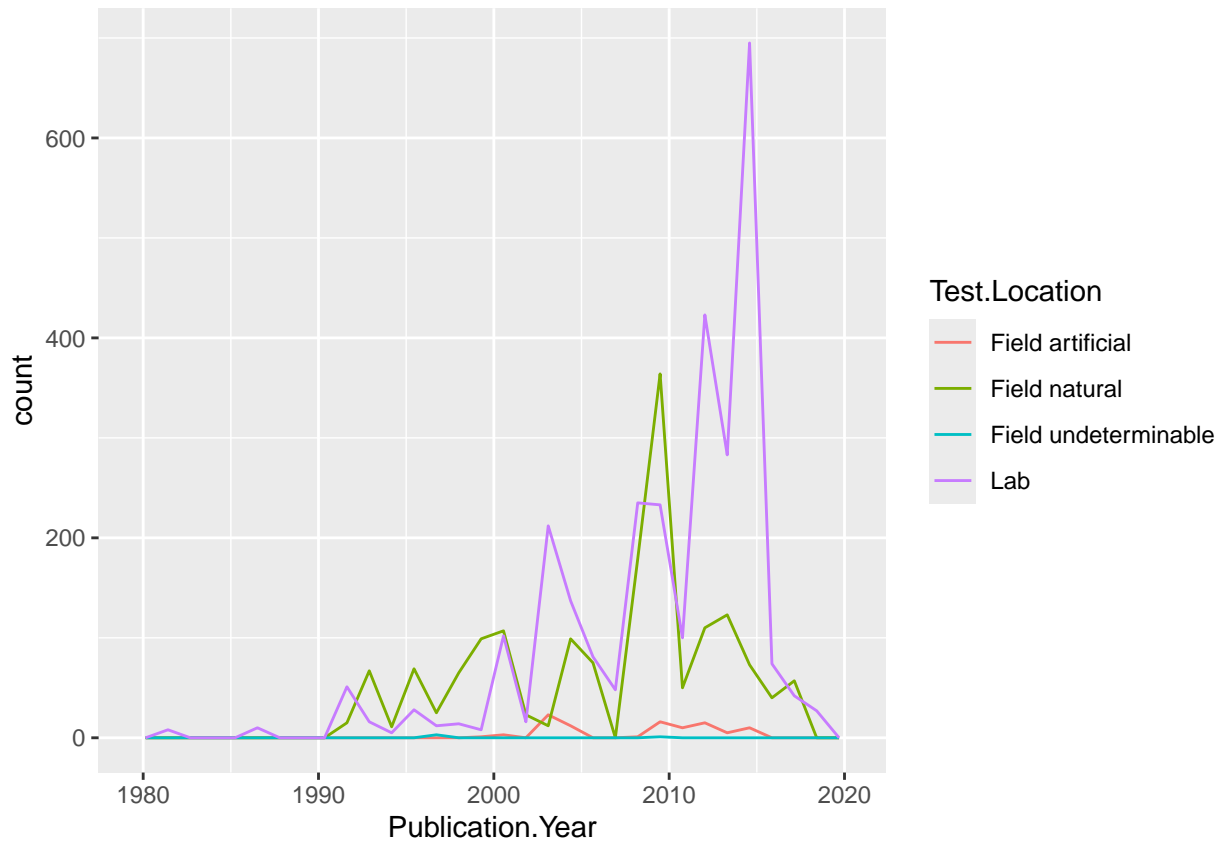
```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



10. Reproduce the same graph but now add a color aesthetic so that different Test.Location are displayed as different colors.

```
ggplot(Neonics,aes(x=Publication.Year,color=Test.Location)) +  
  geom_freqpoly()
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



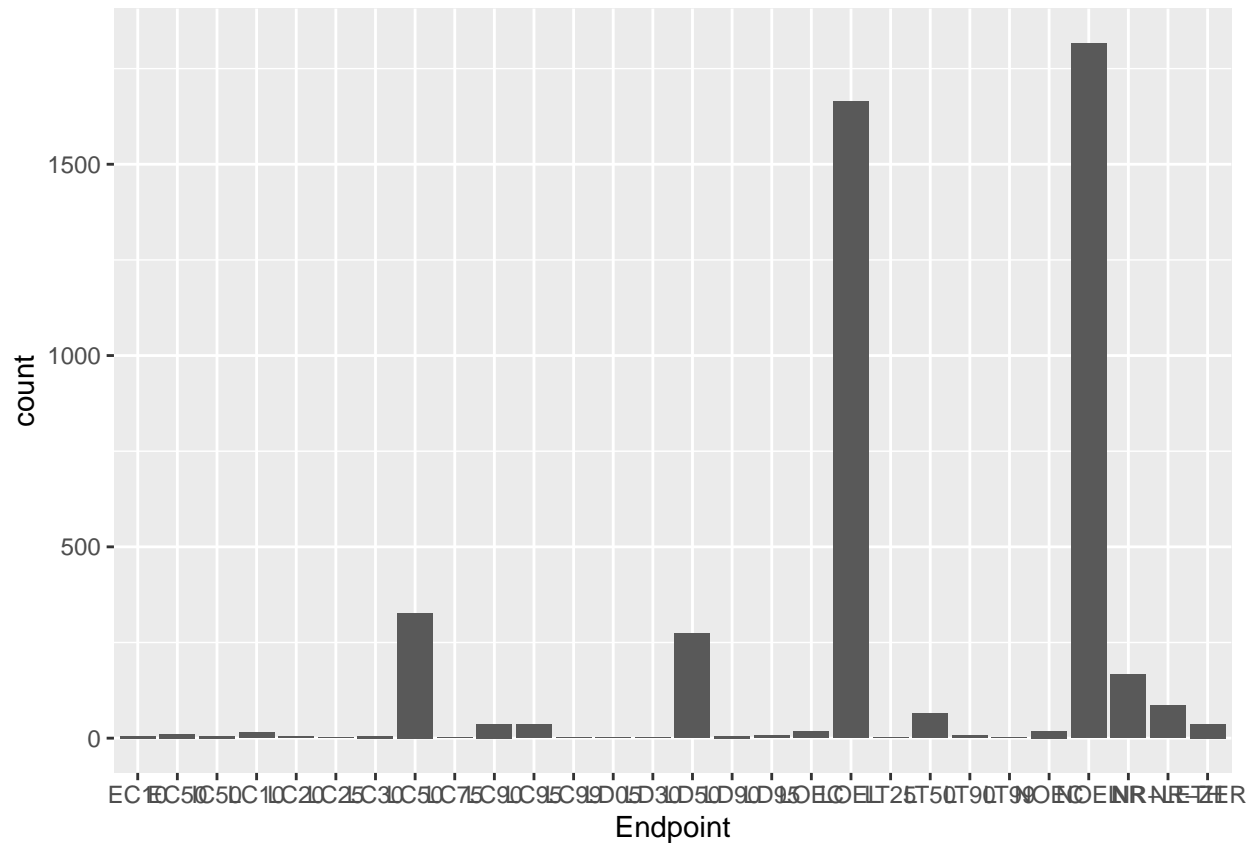
Interpret this graph. What are the most common test locations, and do they differ over time?

Answer: The more common test locations are the lab and field natural. Lab definitely takes over from 2000 onward. Field natural is the closest but in 2015, Lab is out doing Field natural by at least 3x.

11. Create a bar graph of Endpoint counts. What are the two most common end points, and how are they defined? Consult the ECOTOX_CodeAppendix for more information.

[**TIP:** Add `theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1))` to the end of your plot command to rotate and align the X-axis labels...]

```
ggplot(Neonics, aes(x=Endpoint,)) +  
  geom_bar()
```



```
theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1))
```

```
## List of 1
## $ axis.text.x:List of 11
## ..$ family      : NULL
## ..$ face         : NULL
## ..$ colour       : NULL
## ..$ size         : NULL
## ..$ hjust        : num 1
## ..$ vjust        : num 0.5
## ..$ angle        : num 90
## ..$ lineheight   : NULL
## ..$ margin       : NULL
## ..$ debug        : NULL
## ..$ inherit.blank: logi FALSE
## ..- attr(*, "class")= chr [1:2] "element_text" "element"
## - attr(*, "class")= chr [1:2] "theme" "gg"
## - attr(*, "complete")= logi FALSE
## - attr(*, "validate")= logi TRUE
```

Answer: The 2 most common endpoints are LOEL and NOEL - LOEL is Lowest-observable-effect-level and NOEL is No-observable-effect-level.

Explore your data (Litter)

12. Determine the class of collectDate. Is it a date? If not, change to a date and confirm the new class of the variable. Using the `unique` function, determine which dates litter was sampled in August 2018.

```
class(Litter$collectDate)
```

```
## [1] "factor"
```

```
Litter$collectDate = ymd(Litter$collectDate)
class(Litter$collectDate)
```

```
## [1] "Date"
```

```
unique(Litter$collectDate)
```

```
## [1] "2018-08-02" "2018-08-30"
```

Answer: The class of collectDate is “factor”, updated to date and confirmed with `class()` function. Collection dates were 8/2/2018 and 8/30/2018.

13. Using the `unique` function, determine how many different plots were sampled at Niwot Ridge. How is the information obtained from `unique` different from that obtained from `summary`?

```
unique(Litter$plotID)
```

```
## [1] NIWO_061 NIWO_064 NIWO_067 NIWO_040 NIWO_041 NIWO_063 NIWO_047 NIWO_051
## [9] NIWO_058 NIWO_046 NIWO_062 NIWO_057
## 12 Levels: NIWO_040 NIWO_041 NIWO_046 NIWO_047 NIWO_051 NIWO_057 ... NIWO_067
```

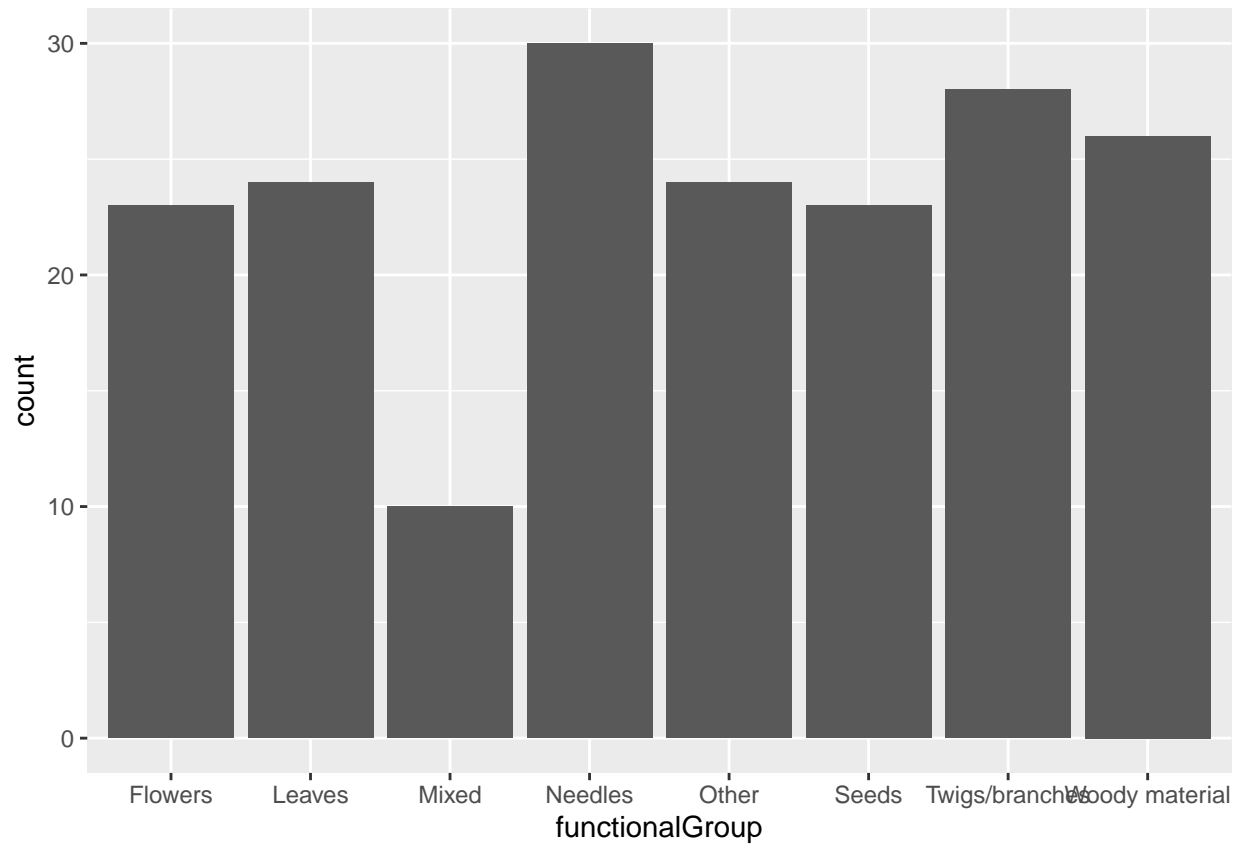
```
summary(Litter$plotID)
```

```
## NIWO_040 NIWO_041 NIWO_046 NIWO_047 NIWO_051 NIWO_057 NIWO_058 NIWO_061
##      20      19      18      15      14       8      16      17
## NIWO_062 NIWO_063 NIWO_064 NIWO_067
##      14      14      16      17
```

Answer: Using the `unique` function, 12 different plots were sampled. `Unique` provides the unique values in a field. `Summary` provides the unique values and the count numbers of them.

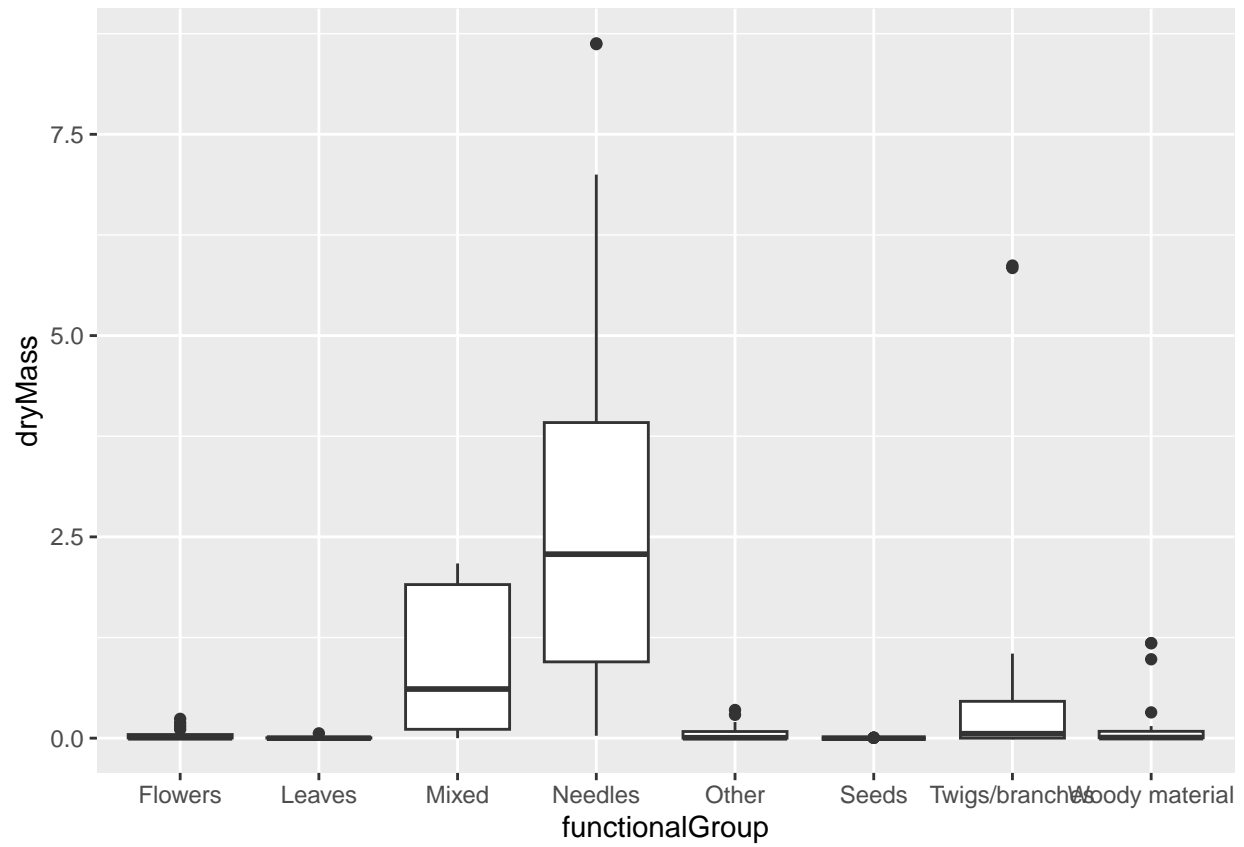
14. Create a bar graph of functionalGroup counts. This shows you what type of litter is collected at the Niwot Ridge sites. Notice that litter types are fairly equally distributed across the Niwot Ridge sites.

```
ggplot(Litter, aes(x=functionalGroup,)) +
  geom_bar()
```

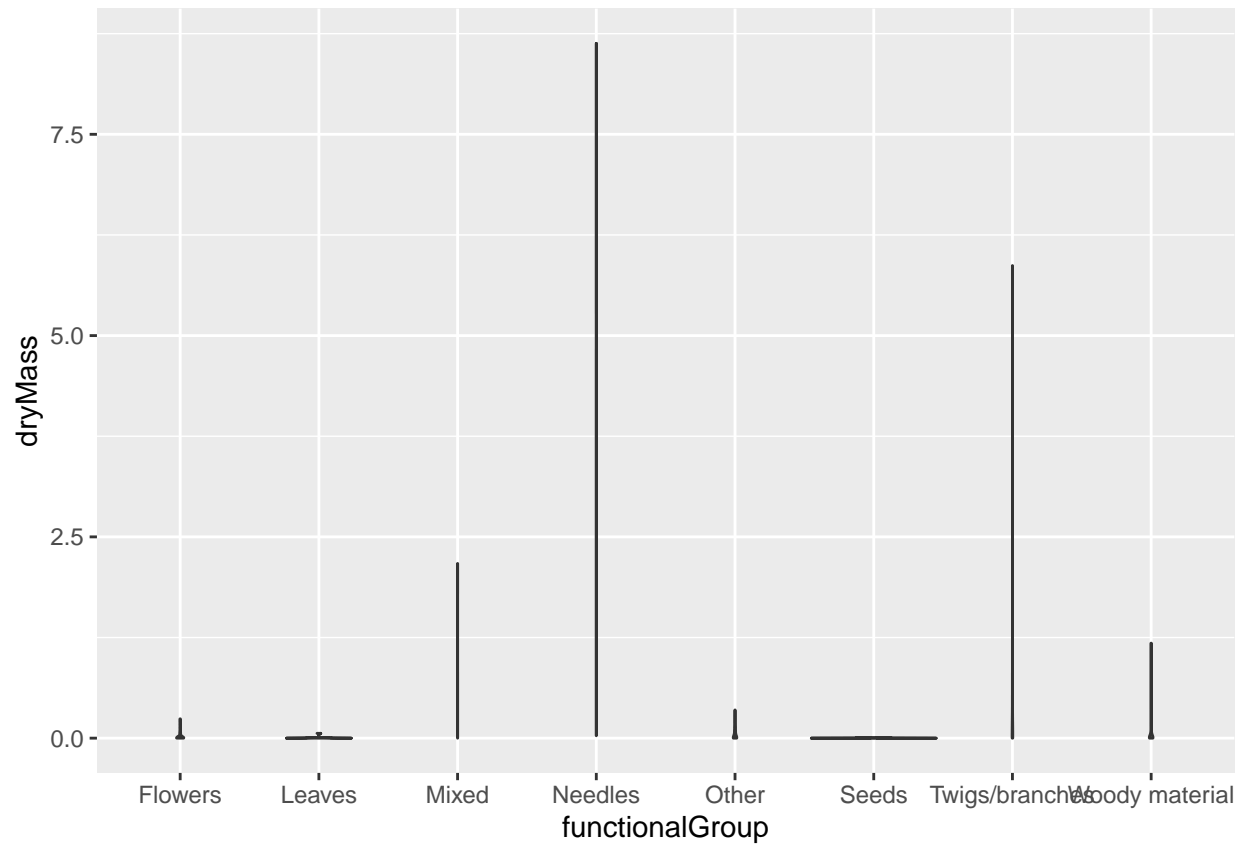


15. Using `geom_boxplot` and `geom_violin`, create a boxplot and a violin plot of `dryMass` by `functionalGroup`.

```
ggplot(Litter, aes(x = functionalGroup, y = dryMass)) +  
  geom_boxplot()
```

```
ggplot(Litter, aes(x = functionalGroup, y = dryMass)) +  
  geom_violin()
```



Why is the boxplot a more effective visualization option than the violin plot in this case?

Answer: There is no need for density distribution in this dataset. Boxplot is fine for this more simple data set.

What type(s) of litter tend to have the highest biomass at these sites?

Answer: The highest biomass seems to be mainly Needles and Mixed.