

Assignment 3: Data Exploration

Gaby Antonova, Section #3

OVERVIEW

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

Directions

1. Change “Student Name, Section #” on line 3 (above) with your name and section number.
2. Work through the steps, **creating code and output** that fulfill each instruction.
3. Be sure to **answer the questions** in this assignment document.
4. When you have completed the assignment, **Knit** the text and code into a single PDF file.
5. After Knitting, submit the completed exercise (PDF file) to the dropbox in Sakai. Add your last name into the file name (e.g., “FirstLast_A03_DataExploration.Rmd”) prior to submission.

The completed exercise is due on <>.

Set up your R session

1. Check your working directory, load necessary packages (tidyverse), 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 add the stringsAsFactors = TRUE parameter to the function when reading in the CSV files.**

```
getwd()

## [1] "C:/Users/gabri/OneDrive/Desktop/MPP Coursework/Spring 2022/ENVIRO 872/Environmental_Data_Analyt.

library(tidyverse) #load packages

Neonics <- read.csv("../Data/Processed/Lab 3/ECOTOX_Neonicotinoids_Insects_raw.csv", stringsAsFactors=TRUE)
Litter <- read.csv("../Data/Processed/Lab 3/NEON_NIWO_Litter_massdata_2018-08_raw.csv", stringsAsFactors=TRUE)

view(Neonics)

view(Litter)
```

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: We might be interested in how neonicotinoids impact insects, specifically pollinators, since they are more toxic to invertebrates. Neonicotinoids are absorbed by plants and can be found in pollen and nectar, which can harm pollinators like bees. World food crop production depends on pollinators so studying how an insecticide might harm them is crucial to avoid potential food shortages.

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: Woody debris is an important part of forest ecosystems and provide ecological services like carbon sequestration, replenishing soil nutrients, and reducing erosion. Depending on the type of woody debris, it can also be a source of increased wildlife risk which would be another reason to study it.

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: *Litter is defined as material dropped from the forest canopy with a butt end diameter of less than 2 cm and a length of less than 50 cm. Fine wood is defined as material that has a butt end diameter less than 2 cm but a length greater than 50 cm andd lastly.* Litter and fine woody debris are collected from elevated and ground traps, respectively at sites that contain woody vegetation taller than 2 meters. *Ground traps are sampled once per year and elevated traps are sampled one time every two weeks in deciduous forest (during senescence) and a time every 1-2 months at evergreen sites.

Obtain basic summaries of your data (Neonics)

5. What are the dimensions of the dataset? The Neonics dataset has 4623 rows and 30 columns.

```
dim(Neonics)
```

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

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?

```
summary(Neonics$Effect, maxsum=6)
```

```
##      Population      Mortality      Behavior Feeding behavior
##           1803           1493           360           255
##      Reproduction      (Other)
##           197           515
```

Answer: The most common effects that were studied were Population and Mortality. These are likely frequently studied because they can indicate whether the neonicotinoids are lethal to insects of interest and more broadly whether they are leading to population changes, like population decline.

- 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.

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

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

Answer: The six most commonly studies species of insect in the dataset were Honey Bees, Parasitic Wasp, Buff Tailed Bumblebees, Carniolan Honey Bees, and Bumble Bees. The top species are all members of the Hymenoptera order of insects. Almost all of the top species are bees, which are key pollinators and the only other species, the parasitic wasp can inadvertently perform pollinating services, which again are critical for world crop production.

- Concentrations are always a numeric value. What is the class of `Conc.1..Author.` in the dataset, and why is it not numeric?

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

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

Answer: The class is factor. This is because the variable contains non-numeric variables like "NR."

Explore your data graphically (Neonics)

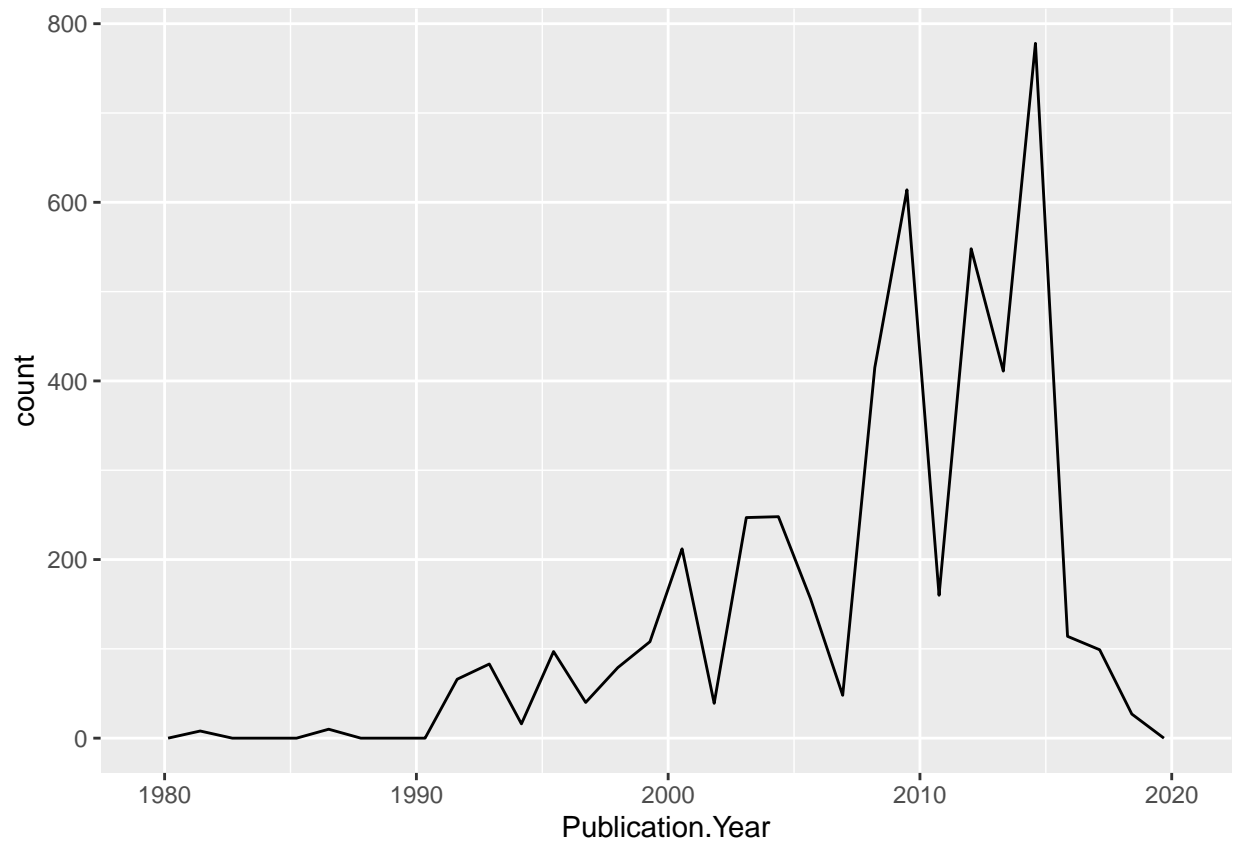
- Using `geom_freqpoly`, generate a plot of the number of studies conducted by publication year.

```
class(Neonics$Publication.Year)
```

```
## [1] "integer"
```

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

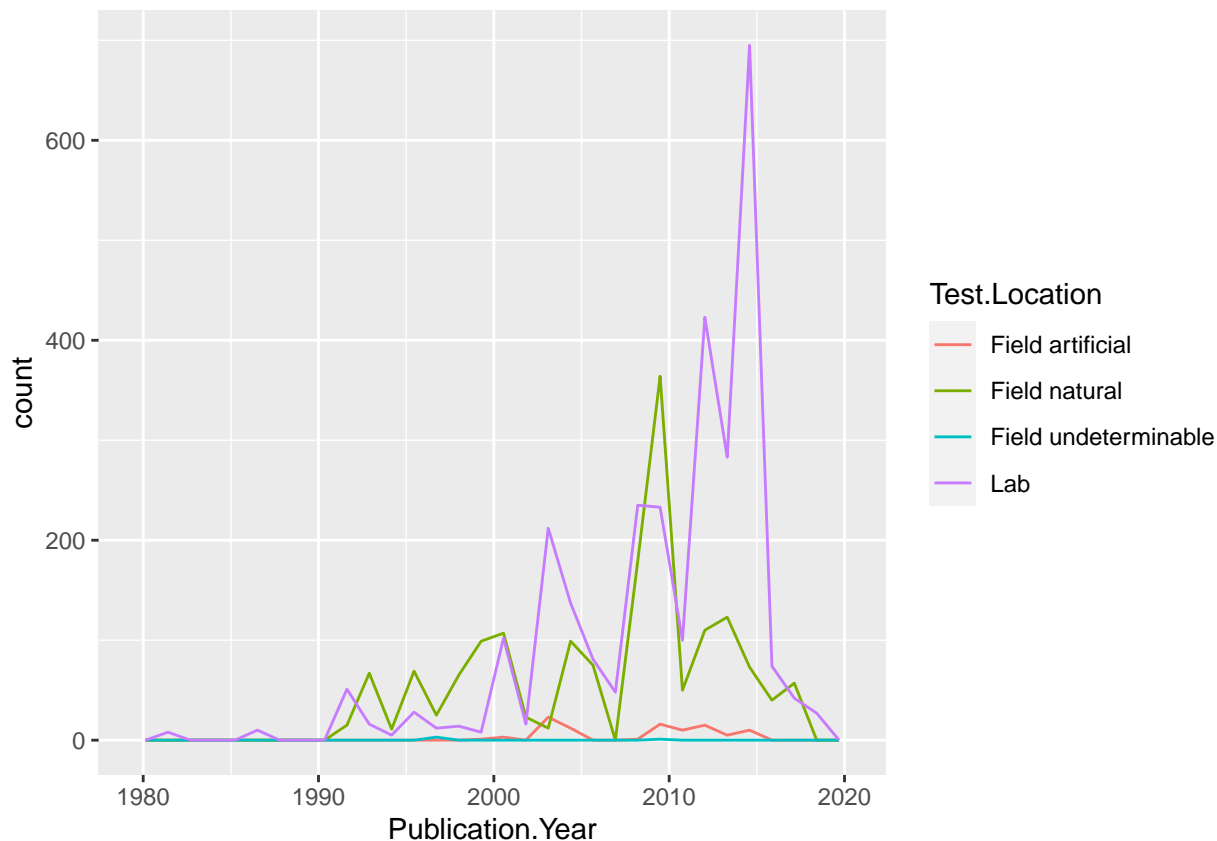
```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



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

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

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

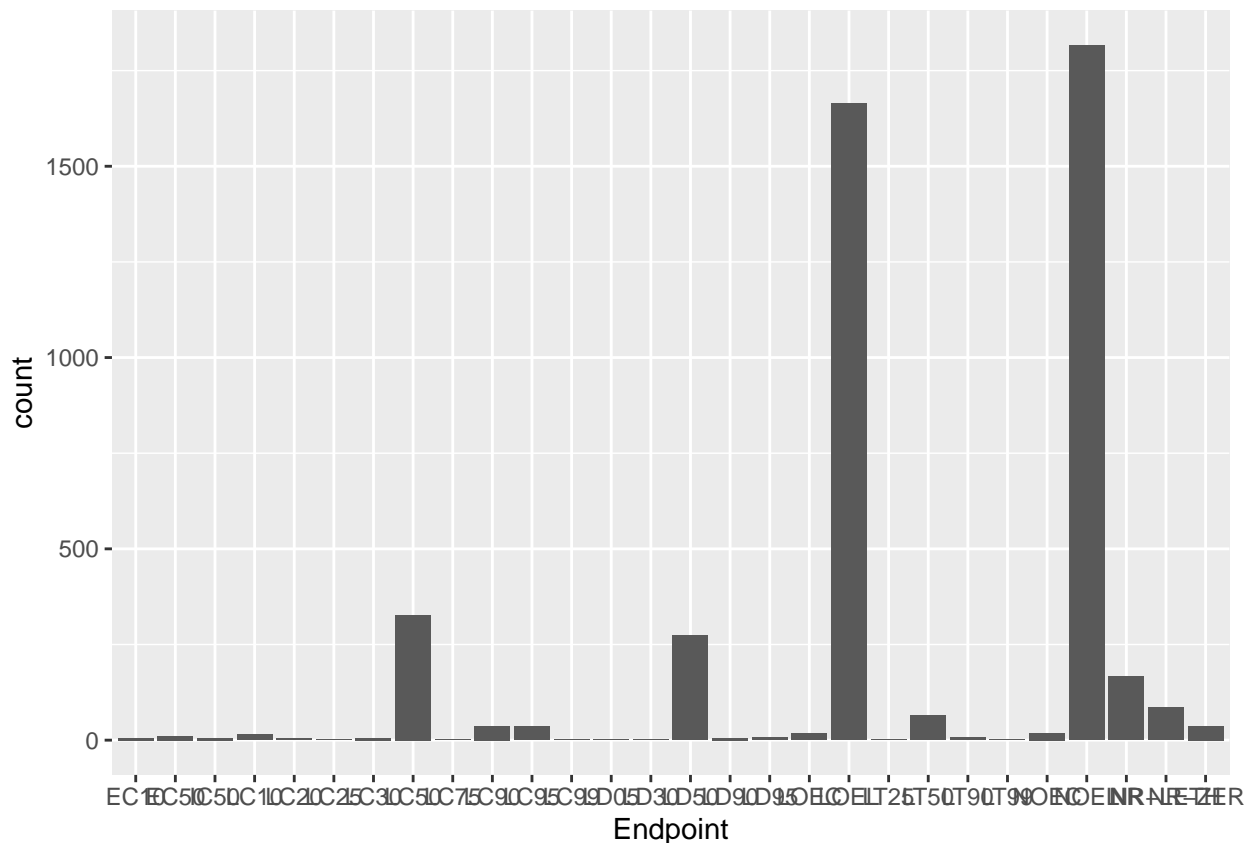


Interpret this graph. What are the most common test locations, and do they differ over time?

Answer: In the 1990s, the most common test location was in the field natural. In the early 2000s, there was a spike in the number of studies about neonicotinoids done in the lab. Then there was a spike in natural field studies around 2009, until finally lab studies dominated in the mid-2010s. By around 2016, all types of studies regarding neonicotinoids declined, which might indicate that the bans on neonicotinoids starting in 2013 in the EU might have made new studies in the region unnecessary.

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.

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



Answer: NOEL and LOEL are the two most common endpoints. NOEL is defines are “no-observable-effect-level: The highest dose producing effects not significantly different from responses of controls according to author’s reported statistical test.” LOEL is defined as “Lowest-observable-effect-level: lowest dose producing effects that were significantly different (as reported by authors) from responses of controls.”

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"
```

```
library(lubridate)
```

##

```
## Attaching package: 'lubridate'
```

```
## The following objects are masked from 'package:base':
```

##

```
##      date, intersect, setdiff, union
```

```
date_collectDate <- ymd(Litter$collectDate)
date_collectDate
```

```
## [1] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [6] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [11] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [16] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [21] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [26] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [31] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [36] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [41] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [46] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [51] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [56] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [61] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [66] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [71] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [76] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [81] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [86] "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02" "2018-08-02"
## [91] "2018-08-02" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [96] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [101] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [106] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [111] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [116] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [121] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [126] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [131] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [136] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [141] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [146] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [151] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [156] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [161] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [166] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [171] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [176] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [181] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
## [186] "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30" "2018-08-30"
```

```
class(date_collectDate)
```

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

The data was initially factor data.

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

```
unique(Litter$plotID, incomparables = FALSE)
```

```
## [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
```

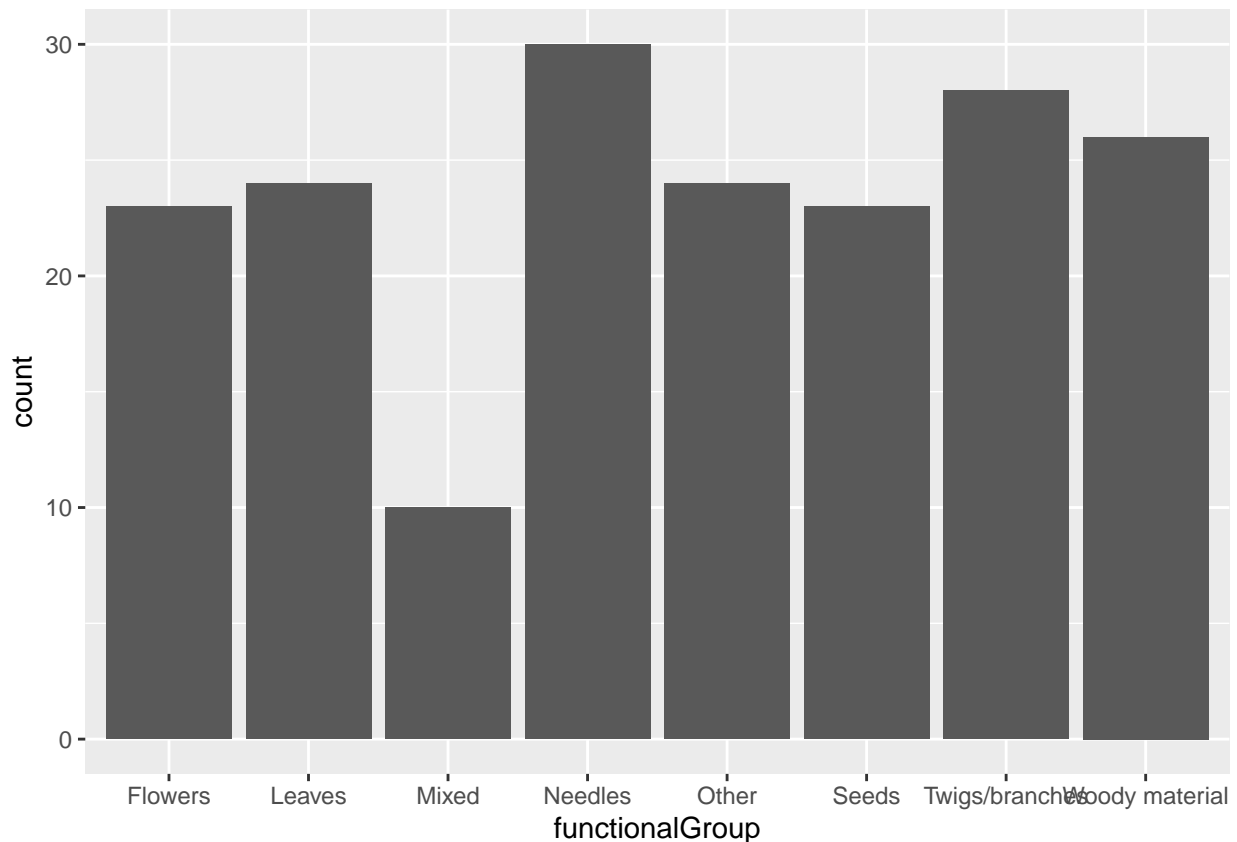
```
summary(Litter$plotID)
```

```
## 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: 12 plots were sampled at Niwot Ridge. The difference between the summary function and the unique function is that the summary function also displays the frequency of each plot sampled, while the unique function just lists all unique values once.

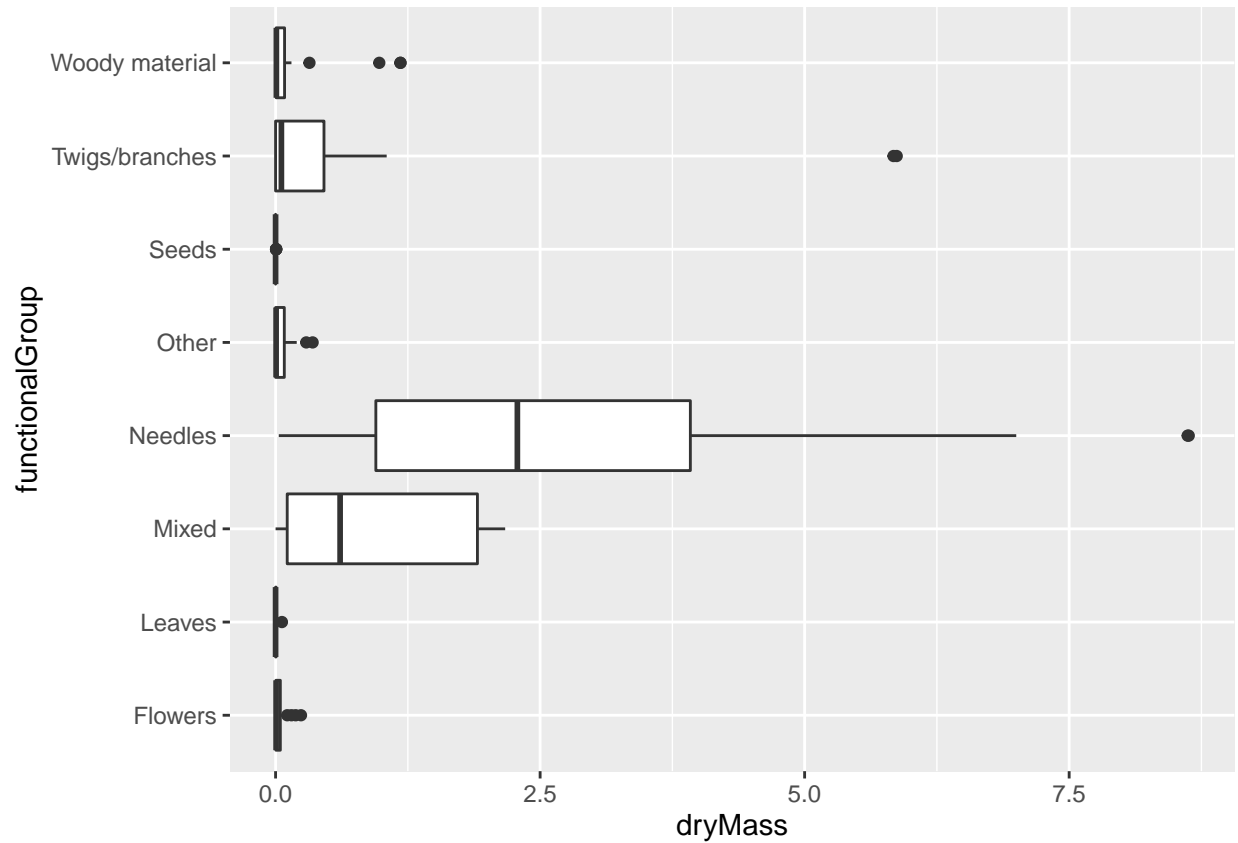
14. Create a bar graph of functionalGroup counts. This shows you what type of litter is collected at the Niwot Ridge sites. Notice that litter types are fairly equally distributed across the Niwot Ridge sites.

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

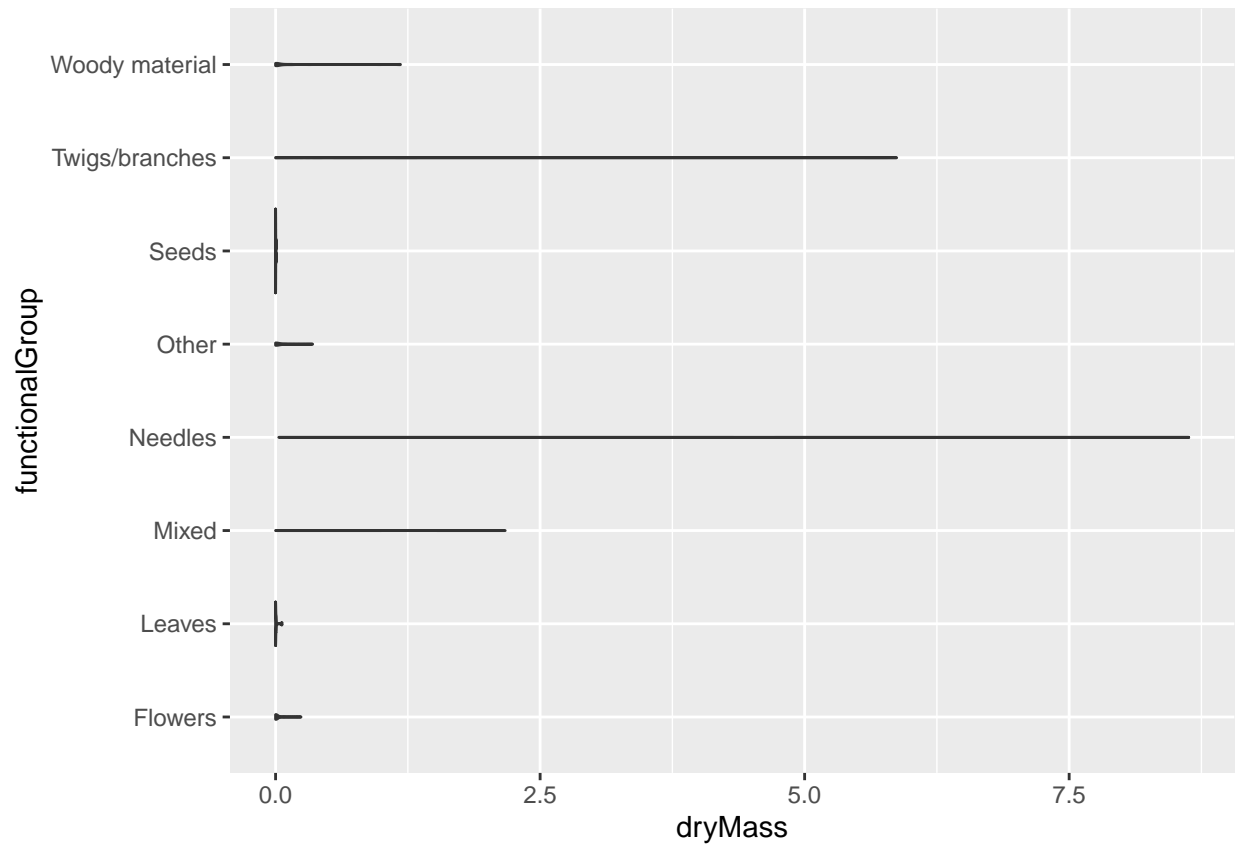


15. Using 'geom_boxplot' and 'geom_violin', create a boxplot and a violin plot of dryMass by functional

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



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



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

Answer: The boxplot is much more effective because the distribution of dry mass within each functional group is relatively uniform so there is no large difference in distribution to visualize.

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

Answer: Needles tend to have the highest biomass at these sites.