# Assignment 3: Data Exploration

## Judy Hua Zhu

## 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.

```
library(tidyverse)
library(lubridate)
library(here)

getwd()
```

```
## [1] "/home/guest/872 Projects/EDE_Fall2024"
```

```r
Neoics <- read.csv(
  file = here('Data','Raw','ECOTOX_Neonicotinoids_Insects_raw.csv'),
  stringsAsFactors = T
)

Litter <- read.csv(
  file = here('Data','Raw','NEON_NIWO_Litter_massdata_2018-08_raw.csv'),
  stringsAsFactors = T
)
```

## 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 study ecotoxicology of neonicotinoids to learn how does it work on various insects, animals, the products itself and the ecosystem (water, soil, air, etc.), not only on the insects we target to kill. Focusing on insects, we need to assess whether they are toxic for beneficial insects, pollinators, invertebrates that we are not intend to harm.

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: The litter and woody debris provides information on the plants functional groups in the area, reflecting the environment they grew, possible incidents (of climate disaster) they experienced, and traits suitable to thrive in the environment. They may also provide data on nutrition content into the ground and be useful to estimate the biomass and productivity of the area and to provide info on its influence to carbon flux and water resource.

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 is defined as materials with diameters <2m and lengths <50m, collected in an elevated 0.5m^2 PVC trap; Woody debris is defined as materials with diameters <2m and length>50m, collected in ground traps. 2. The spacial sampling of data various on the vegetation height and geological conditions of the location, which must be tower plots. 3. Temporal Sampling Design: Ground traps are sampled once per year, while elevated traps are sampled more frequently (once every two weeks/a month) based on the vegetation type (deciduous/evergreen).

## Obtain basic summaries of your data (Neonics)

5. What are the dimensions of the dataset?

```r
str(Neoics)
```

```
## 'data.frame':    4623 obs. of  30 variables:
##  $ CAS.Number                   : int  58842209 58842209 58842209 58842209 58842209 58842209 5884:
##  $ Chemical.Name                : Factor w/ 9 levels "(1E)-N-[(6-Chloro-3-pyridinyl)methyl]-N-eth:
##  $ Chemical.Grade               : Factor w/ 9 levels "Analytical grade",..: 9 9 9 9 9 9 9 9 9 9 .:
##  $ Chemical.Analysis.Method     : Factor w/ 5 levels "Measured","Not coded",..: 4 4 4 4 4 4 4 4 4 4:
##  $ Chemical.Purity              : Factor w/ 80 levels ">=98",">=99.0",..: 69 69 50 50 50 50 50 50 :
##  $ Species.Scientific.Name      : Factor w/ 398 levels "Acalolepta vastator",..: 69 69 248 248 248:
##  $ Species.Common.Name          : Factor w/ 303 levels "Alfalfa Leafcutter Bee",..: 74 74 142 142 :
##  $ Species.Group                : Factor w/ 4 levels "Insects/Spiders",..: 1 1 1 1 1 1 1 1 1 1 1 .:
##  $ Organism.Lifestage           : Factor w/ 20 levels "Adult","Cocoon",..: 1 1 19 19 19 1 19 1 1 :
##  $ Organism.Age                 : Factor w/ 39 levels "<=24","<=48",..: 39 39 39 39 39 36 39 36 3(:
##  $ Organism.Age.Units           : Factor w/ 11 levels "Day(s)","Days post-emergence",..: 9 9 4 4 ·:
##  $ Exposure.Type                : Factor w/ 24 levels "Choice","Dermal",..: 23 23 11 11 11 11 11 :
##  $ Media.Type                   : Factor w/ 10 levels "Agar","Artificial soil",..: 7 7 3 3 3 3 3 3:
##  $ Test.Location                : Factor w/ 4 levels "Field artificial",..: 4 4 4 4 4 4 4 4 4 4 .:
##  $ Number.of.Doses              : Factor w/ 30 levels "' 4-5","' 4-7",..: 30 30 18 18 18 18 18 18 :
##  $ Conc.1.Type..Author.         : Factor w/ 3 levels "Active ingredient",..: 1 1 1 1 1 1 1 1 1 1 1 :
##  $ Conc.1..Author.              : Factor w/ 1006 levels "<0.0004","<0.025",..: 639 510 813 622 44:
##  $ Conc.1.Units..Author.        : Factor w/ 148 levels "%","% v/v","% w/v",..: 132 132 91 91 91 9:
##  $ Effect                       : Factor w/ 19 levels "Accumulation",..: 16 16 16 16 16 16 16 16 :
##  $ Effect.Measurement           : Factor w/ 155 levels "Abundance","Accuracy of learned task, per:
##  $ Endpoint                     : Factor w/ 28 levels "EC10","EC50",..: 15 15 8 8 8 8 8 8 8 8 ... :
##  $ Response.Site                : Factor w/ 19 levels "Abdomen","Brain",..: 14 14 14 14 14 14 14 :
##  $ Observed.Duration..Days.     : Factor w/ 361 levels "<.0002","<.0021",..: 145 145 145 145 145 :
##  $ Observed.Duration.Units..Days.: Factor w/ 17 levels "Day(s)","Day(s) post-emergence",..: 1 1 1 :
##  $ Author                       : Factor w/ 433 levels "Abbott,V.A., J.L. Nadeau, H.A. Higo, and :
##  $ Reference.Number             : int  107388 107388 103312 103312 103312 103312 103312 103312 10:
##  $ Title                        : Factor w/ 458 levels "A Common Pesticide Decreases Foraging Suc:
##  $ Source                       : Factor w/ 456 levels "Acta Hortic.1094:451-456",..: 295 295 296 :
##  $ Publication.Year             : int  1982 1982 1986 1986 1986 1986 1986 1986 1986 1986 ... :
##  $ Summary.of.Additional.Parameters: Factor w/ 943 levels "Purity: \xca NC - NC | Organism Age: \xca:
```

Answer:4623 objetcs(rows) of 30 variables (columns)

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...]

```r
NeoEffect <- summary(Neoics$Effect)
sort(NeoEffect)
```

```
##      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
```

3

Answer: The most common effect is Population (1803 out of 4623), followed by Mortality (1493) and Bahavior (360). We are interested in how does various Neonicotinoids influence insects the most, and population and mortality effect rank on the top shows that it influences the most with insects population sizes.

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...]

```
summary(Neoics$Species.Common.Name, maxsum = 7)
```

```
##           Honey Bee      Parasitic Wasp Buff Tailed Bumblebee
##                 667                   285                 183
##   Carniolan Honey Bee         Bumble Bee      Italian Honeybee
##                 152                   140                 113
##             (Other)
##                3083
```

Answer: Honey Bee(667); Parasitic Wasp Buff Tailed Bumblebee(285); Carniolan Honey Bee(183); Bumble Bee(152); Italian Honeybee (3083). The most commonly studied species are all bees, which are critical pollinators for plants. We are interested in whether and how Neoicotinoids affects their population, in order to minimize the negative effects.

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(Neoics$Conc.1..Author.)
```

```
## [1] "factor"
```
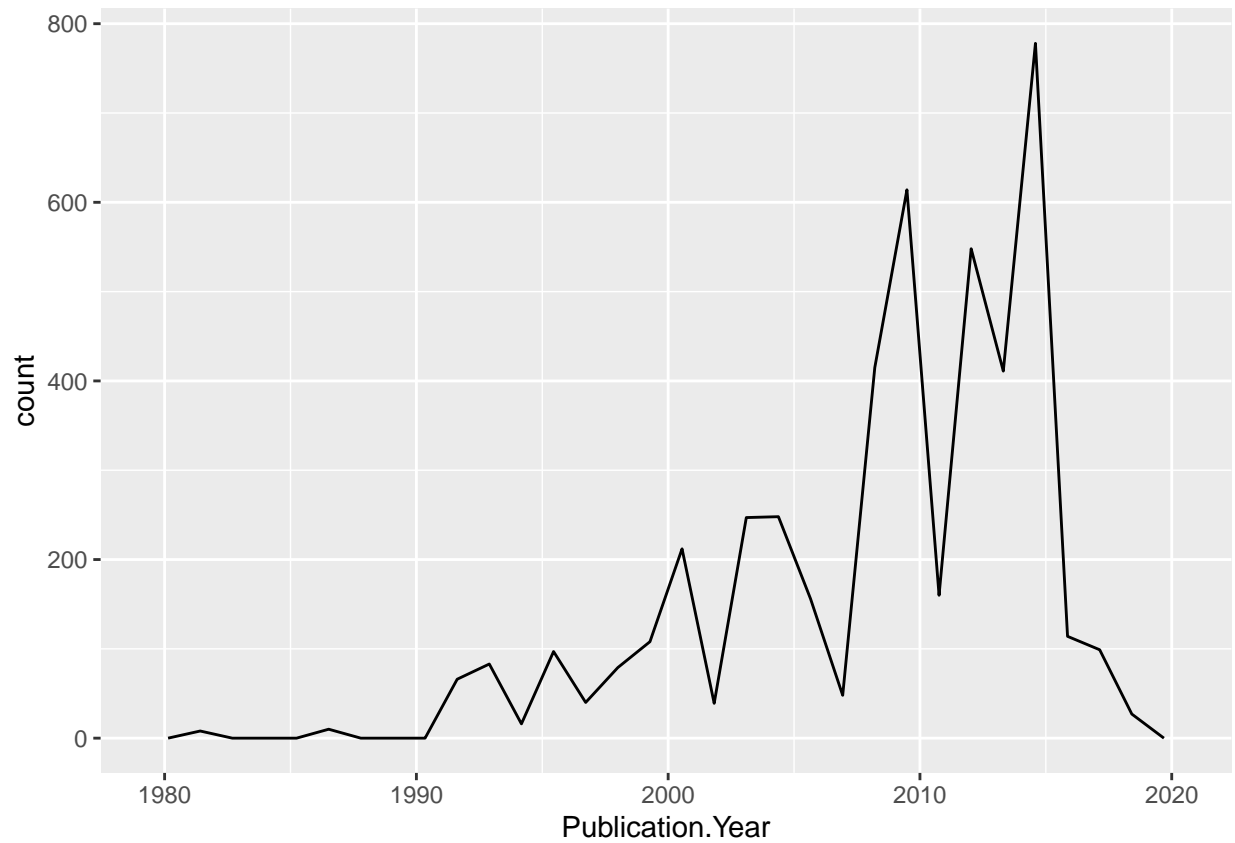
```
view(Neoics$Conc.1..Author.)
```

Answer: Data in 'Conc.1..Author' are factors, because some data include '/' or '~' with the numbers, which make the data not entirely numerical.

## Explore your data graphically (Neonics)

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

```
ggplot(Neoics, aes(x=Publication.Year)) + geom_freqpoly()
```
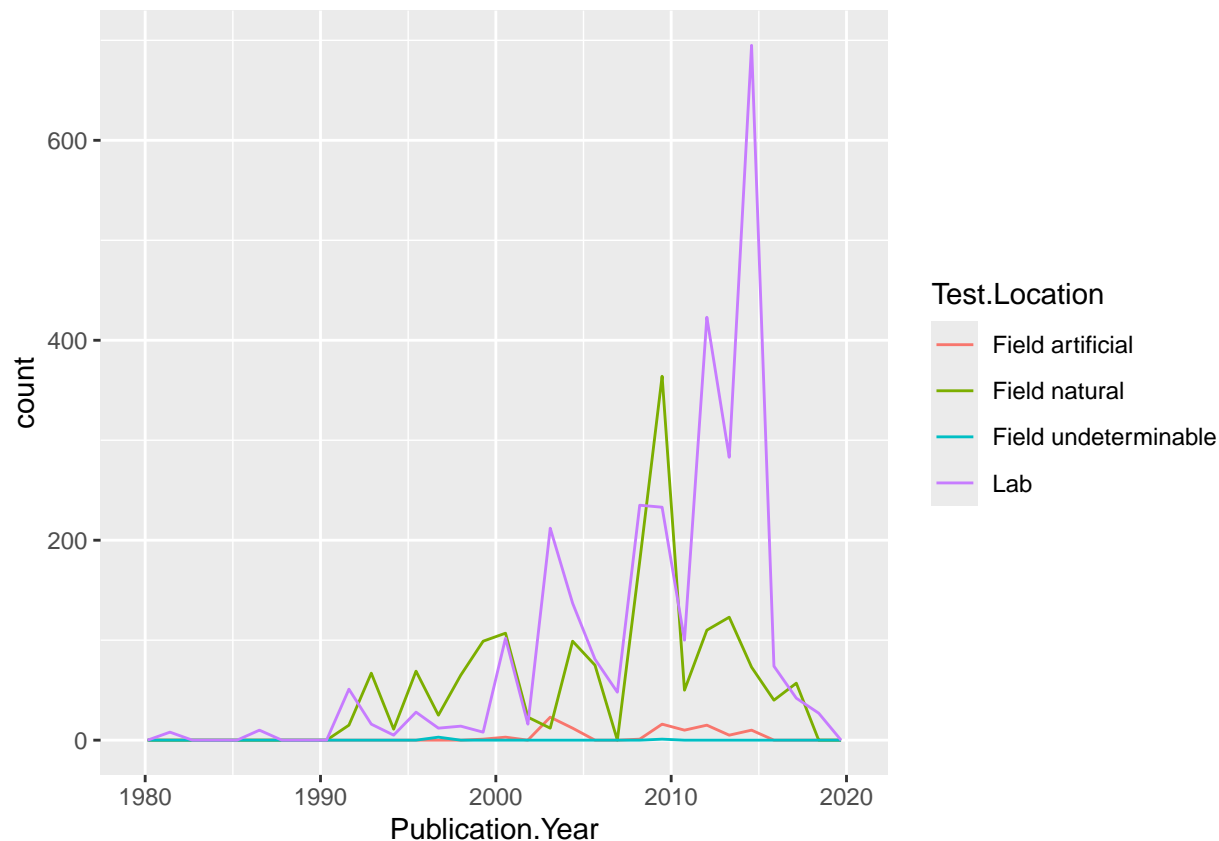
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

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

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

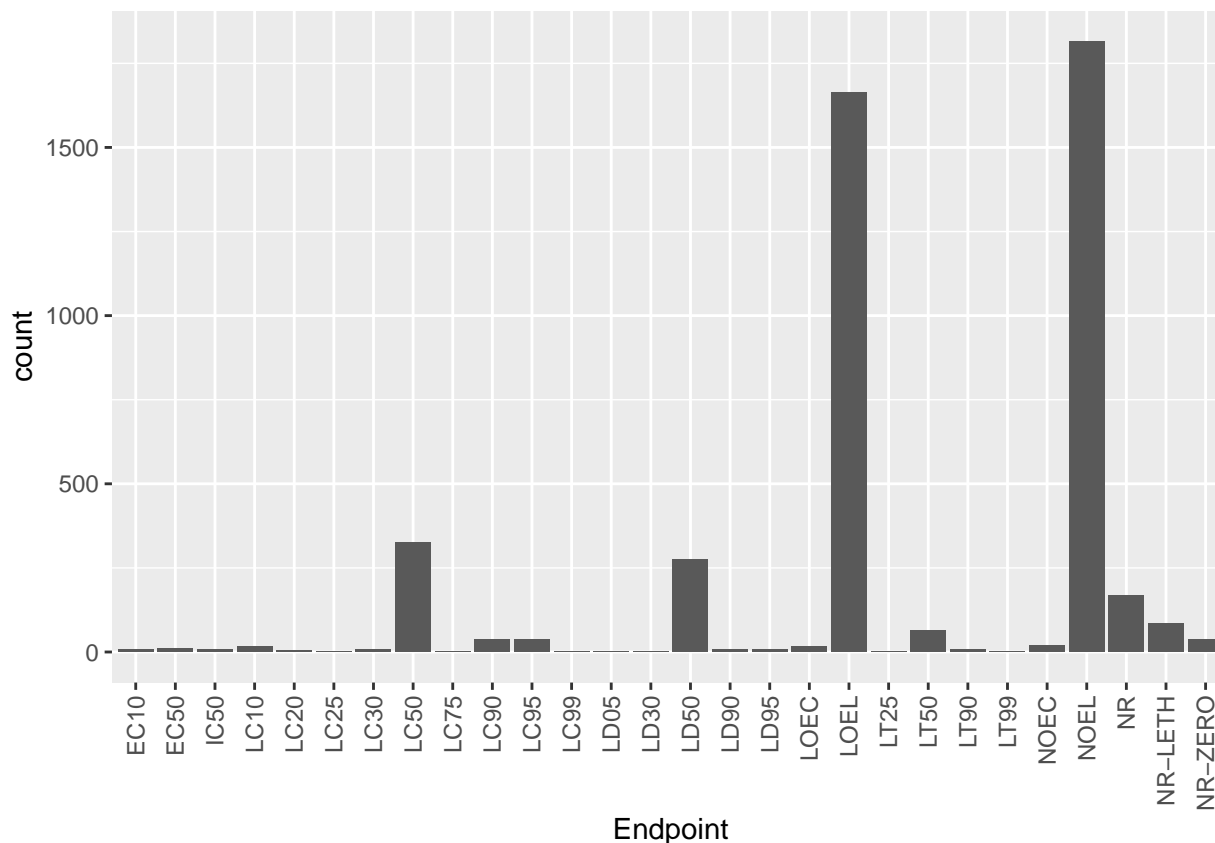## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.

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

> Answer: Most common test locations are in the lab, then in natural fields. Around 1992-2000 and 2008 there was a short time that natural field became most common test location, but other times Labs are the most common.

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(data = Neoics, aes(x = Endpoint)) +
  geom_bar() +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust=1))
```

Answer: The two most common end points are "NOEL" and "LOEL". "NOEL" is defined as "No-observable-effect-level: highest dose (concentration) 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 (concentration) producing effects that were significantly different (as reported by authors) from responses of controls (LOEAL/LOEC), both used in terrestrial database.

## 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$collectDate <- as.Date(Litter$collectDate,format = "%Y-%m-%d")
```

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

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

```
Litter$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"
```

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

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

Answer: Litter was sampled in 02 and 30 of Aug 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`?

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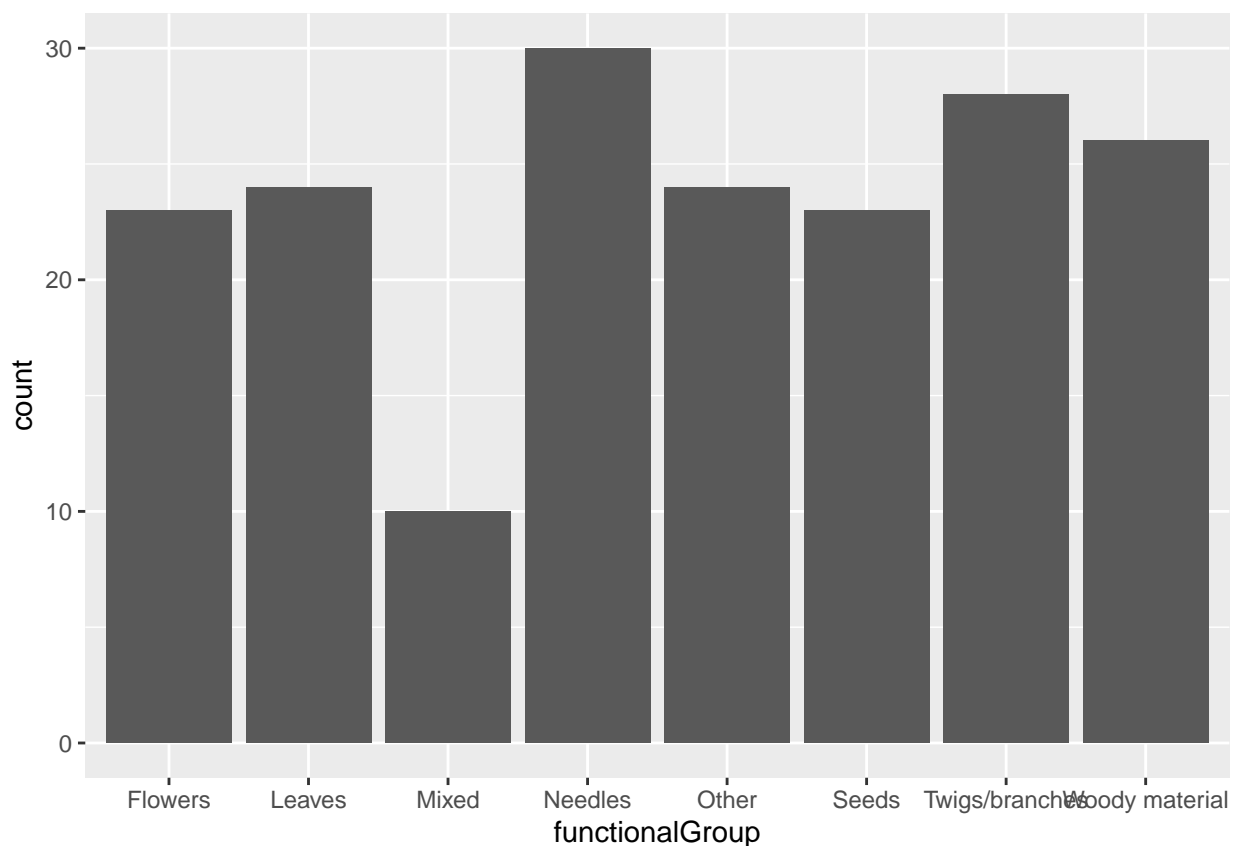
```
unique(Litter$plotID)
```

```
##  [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
```

Answer: Information from summary includes the number of samples of each plots, while the unique does not.
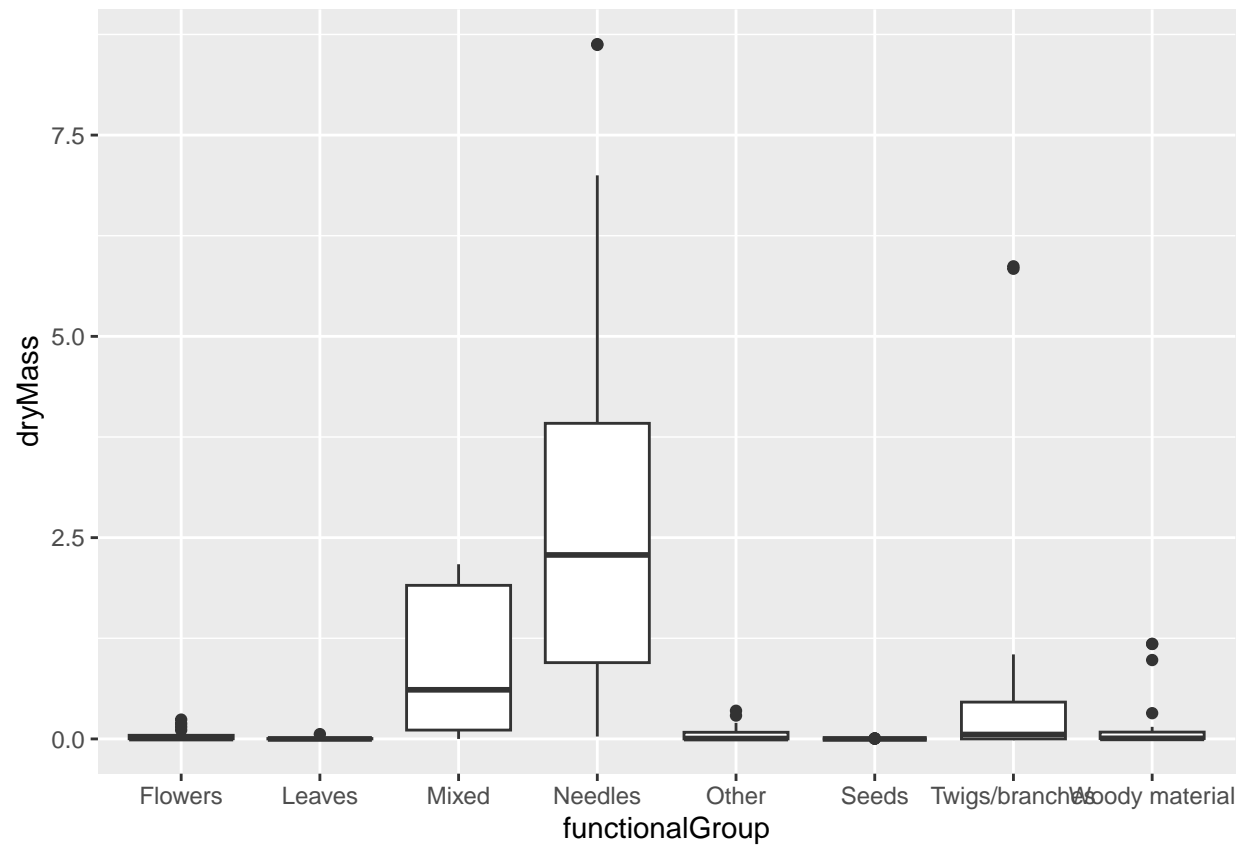
14. Create a bar graph of functionalGroup counts. This shows you what type of litter is collected at the Niwot Ridge sites. Notice that litter types are fairly equally distributed across the Niwot Ridge sites.

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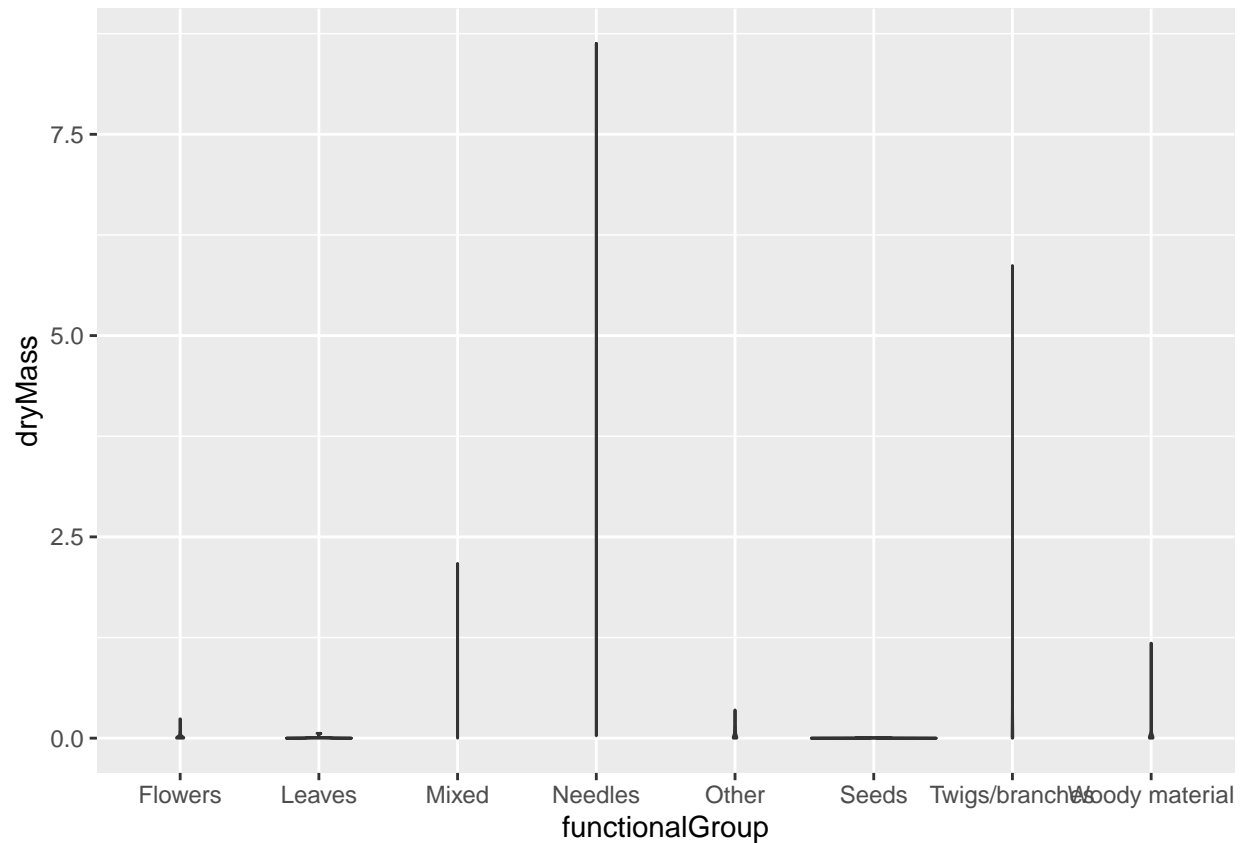


15. Using `geom_boxplot` and `geom_violin`, create a boxplot and a violin plot of dryMass by functional-Group.

```
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: We can see the outlier dots from the boxplots as well as the interquatile range; while the violin cannot show since the dry mass data is dispersed, not much data share the exact same value.

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

Answer: Needles