# Assignment 3: Data Exploration

## Leonardo Rueda

## 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. Be sure to **answer the questions** in this assignment document.
5. When you have completed the assignment, **Knit** the text and code into a single PDF file.
6. After Knitting, submit the completed exercise (PDF file) to the dropbox in Sakai.

The completed exercise is due on Sept 30th.

## 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 include the subcommand to read strings in as factors.

```
getwd()
```

```
## [1] "C:/Users/leor9/OneDrive/Leonardo/MIDP Courses Fall 2022/R Class/EDA-Fall2022/Assignments"
```

```
setwd("C:/Users/leor9/OneDrive/Leonardo/MIDP Courses Fall 2022/R Class/EDA-Fall2022")

library(tidyverse)

Neonics_dataset <- read.csv("./Data/Raw/ECOTOX_Neonicotinoids_Insects_raw.csv",
    stringsAsFactors = TRUE)  # ECOTOX neonicotinoid dataset

Litter_dataset <- read.csv("./Data/Raw/NEON_NIWO_Litter_massdata_2018-08_raw.csv",
    stringsAsFactors = TRUE)  # Niwot Ridge NEON dataset for litter and woody debris
```

## 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:** Yes, for instance some studies have found that some concentrations of this insectiside can be harmful for the population of bees.

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:** We might be interested in studying litter and woody debris in the forest ecosystem because they play a role in carbon budgets and nutrient cycling.

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:** This are three pieces of salient information about the sampling methods: 1) Spatial Sampling Design: Litter and fine woody debris sampling is executed at terrestrial NEON sites that contain woody vegetation >2m tall. 2) Along with most of NEON's plant productivity measurements, sampling for this product occurs only in tower plots 3) Temporal Sampling Design: Ground traps are sampled once per year.

## Obtain basic summaries of your data (Neonics)

5. What are the dimensions of the dataset?

```
dim(Neonics_dataset)
```

```
## [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_dataset$Effect)
```

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

2

**Answer:** The most common effects of the study are Mortality, Population and Feeding behavior. These effects could be of interest to understand how the population of insects evolve when they are exposed to this insecticide.

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.

```
summary(Neonics_dataset$Species.Common.Name)
```

```
##                       Honey Bee              Parasitic Wasp
##                             667                         285
##               Buff Tailed Bumblebee         Carniolan Honey Bee
##                             183                         152
##                       Bumble Bee              Italian Honeybee
##                             140                         113
##                   Japanese Beetle             Asian Lady Beetle
##                              94                          76
##                    Euonymus Scale                    Wireworm
##                              75                          69
##                  European Dark Bee           Minute Pirate Bug
##                              66                          62
##               Asian Citrus Psyllid              Parastic Wasp
##                              60                          58
##             Colorado Potato Beetle            Parasitoid Wasp
##                              57                          51
##                Erythrina Gall Wasp               Beetle Order
##                              49                          47
##         Snout Beetle Family, Weevil   Sevenspotted Lady Beetle
##                              47                          46
##                     True Bug Order          Buff-tailed Bumblebee
##                              45                          39
##                       Aphid Family              Cabbage Looper
##                              38                          38
##                Sweetpotato Whitefly              Braconid Wasp
##                              37                          33
##                       Cotton Aphid              Predatory Mite
##                              33                          33
##               Ladybird Beetle Family              Parasitoid
##                              30                          30
##                     Scarab Beetle               Spring Tiphia
##                              29                          29
##                       Thrip Order          Ground Beetle Family
##                              29                          27
##                 Rove Beetle Family              Tobacco Aphid
##                              27                          27
##                       Chalcid Wasp        Convergent Lady Beetle
##                              25                          25
##                     Stingless Bee             Spider/Mite Class
##                              25                          24
##                 Tobacco Flea Beetle             Citrus Leafminer
##                              24                          23
##                     Ladybird Beetle                   Mason Bee
```

3

| | | |
|---|---:|---:|
| ## | 23 | 22 |
| ## | Mosquito | Argentine Ant |
| ## | 22 | 21 |
| ## | Beetle | Flatheaded Appletree Borer |
| ## | 21 | 20 |
| ## | Horned Oak Gall Wasp | Leaf Beetle Family |
| ## | 20 | 20 |
| ## | Potato Leafhopper | Tooth-necked Fungus Beetle |
| ## | 20 | 20 |
| ## | Codling Moth | Black-spotted Lady Beetle |
| ## | 19 | 18 |
| ## | Calico Scale | Fairyfly Parasitoid |
| ## | 18 | 18 |
| ## | Lady Beetle | Minute Parasitic Wasps |
| ## | 18 | 18 |
| ## | Mirid Bug | Mulberry Pyralid |
| ## | 18 | 18 |
| ## | Silkworm | Vedalia Beetle |
| ## | 18 | 18 |
| ## | Araneoid Spider Order | Bee Order |
| ## | 17 | 17 |
| ## | Egg Parasitoid | Insect Class |
| ## | 17 | 17 |
| ## | Moth And Butterfly Order | Oystershell Scale Parasitoid |
| ## | 17 | 17 |
| ## Hemlock Woolly Adelgid Lady Beetle | | Hemlock Wooly Adelgid |
| ## | 16 | 16 |
| ## | Mite | Onion Thrip |
| ## | 16 | 16 |
| ## | Western Flower Thrips | Corn Earworm |
| ## | 15 | 14 |
| ## | Green Peach Aphid | House Fly |
| ## | 14 | 14 |
| ## | Ox Beetle | Red Scale Parasite |
| ## | 14 | 14 |
| ## | Spined Soldier Bug | Armoured Scale Family |
| ## | 14 | 13 |
| ## | Diamondback Moth | Eulophid Wasp |
| ## | 13 | 13 |
| ## | Monarch Butterfly | Predatory Bug |
| ## | 13 | 13 |
| ## | Yellow Fever Mosquito | Braconid Parasitoid |
| ## | 13 | 12 |
| ## | Common Thrip | Eastern Subterranean Termite |
| ## | 12 | 12 |
| ## | Jassid | Mite Order |
| ## | 12 | 12 |
| ## | Pea Aphid | Pond Wolf Spider |
| ## | 12 | 12 |
| ## | Spotless Ladybird Beetle | Glasshouse Potato Wasp |
| ## | 11 | 10 |
| ## | Lacewing | Southern House Mosquito |
| ## | 10 | 10 |
| ## | Two Spotted Lady Beetle | Ant Family |

```
##                             10                                   9
##                     Apple Maggot                            (Other)
##                              9                                  670
```

**Answer:** The most studied species of insects in the dataset are the Honeybee, the Parasitic Wasp, and the Buff Tailed Bumblebee. They are of interest because they either play a key role in the politization of several plants, or they control the population of other species in the ecosystem.

8. 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_dataset$Conc.1..Author.)
```

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

**Answer:** The variable Conc.1..Author. in the dataset Neonics_dataset is a factor because some of the observations contain non numeric values such as the symbols /, NR/ and ~

## Explore your data graphically (Neonics)

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

```
ggplot(Neonics_dataset) + geom_freqpoly(aes(x = Publication.Year),
    bins = 50)
```

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

```
ggplot(Neonics_dataset) + geom_freqpoly(aes(x = Publication.Year,
    colour = Test.Location), bins = 50)
```
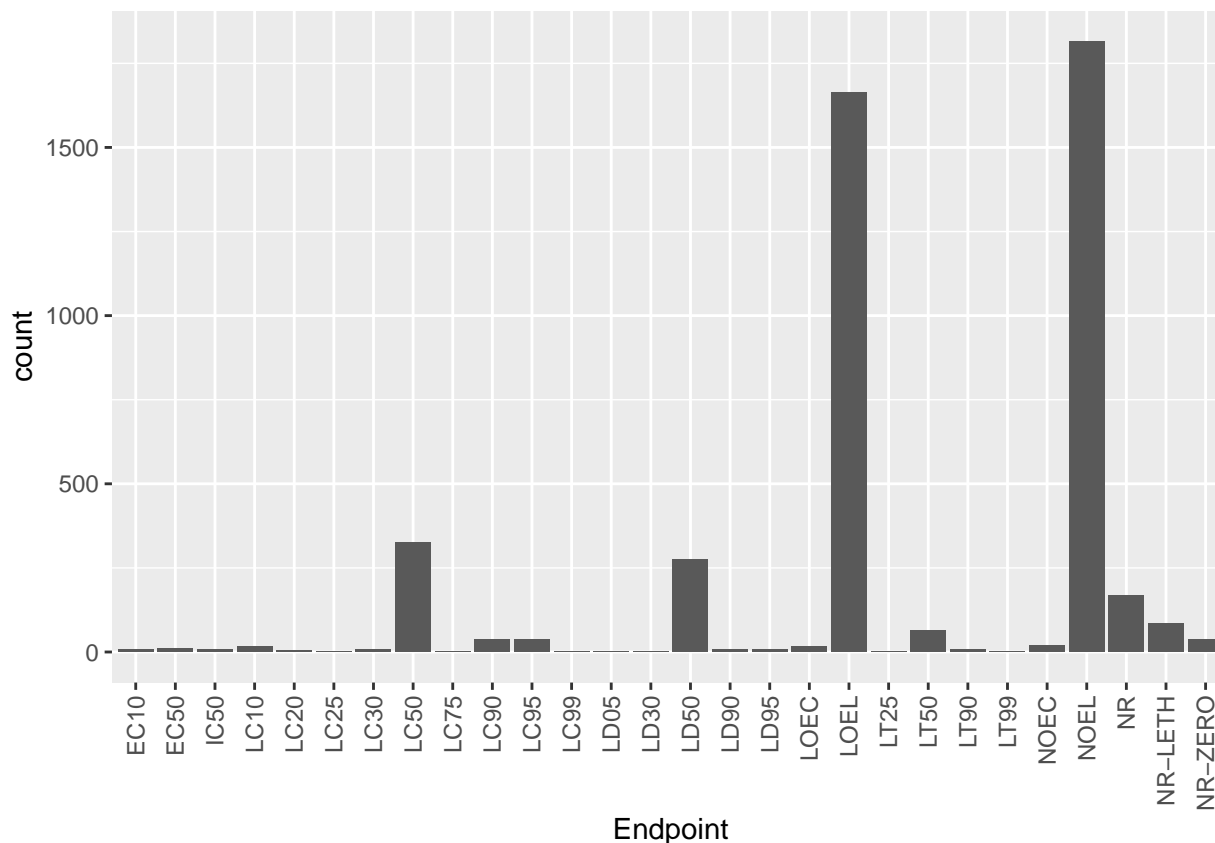


Interpret this graph. What are the most common test locations, and do they differ over time?

**Answer:** The most common test locations in the period analyzed are "Lab" and "Field natural" and the relative importance of each one has changed over time. For instance, between 1990 and 2000, "Field natural" was the most common, while between 2000 and 2020 "Lab" predominated most of the time (except for a couple of years before 2010).

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

**Answer:** The two most common endpoints are "NOEL" and "LOEL". "NOEL" is defined as No-observable-effect-level, with the highest dose (concentration) producing effects not significantly different from responses of controls according to the author's reported statistical test (NOEAL/NOEC), and "LOEL" is defined as the lowest-observable-effect-level, with the lowest dose (concentration) producing effects that were significantly different (as reported by authors) from responses of controls (LOEAL/LOEC).

## 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_dataset$collectDate)
```

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

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

class(Litter_dataset$collectDate)
```

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

7

```
unique(Litter_dataset$collectDate)
```

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

13. 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_dataset$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
```
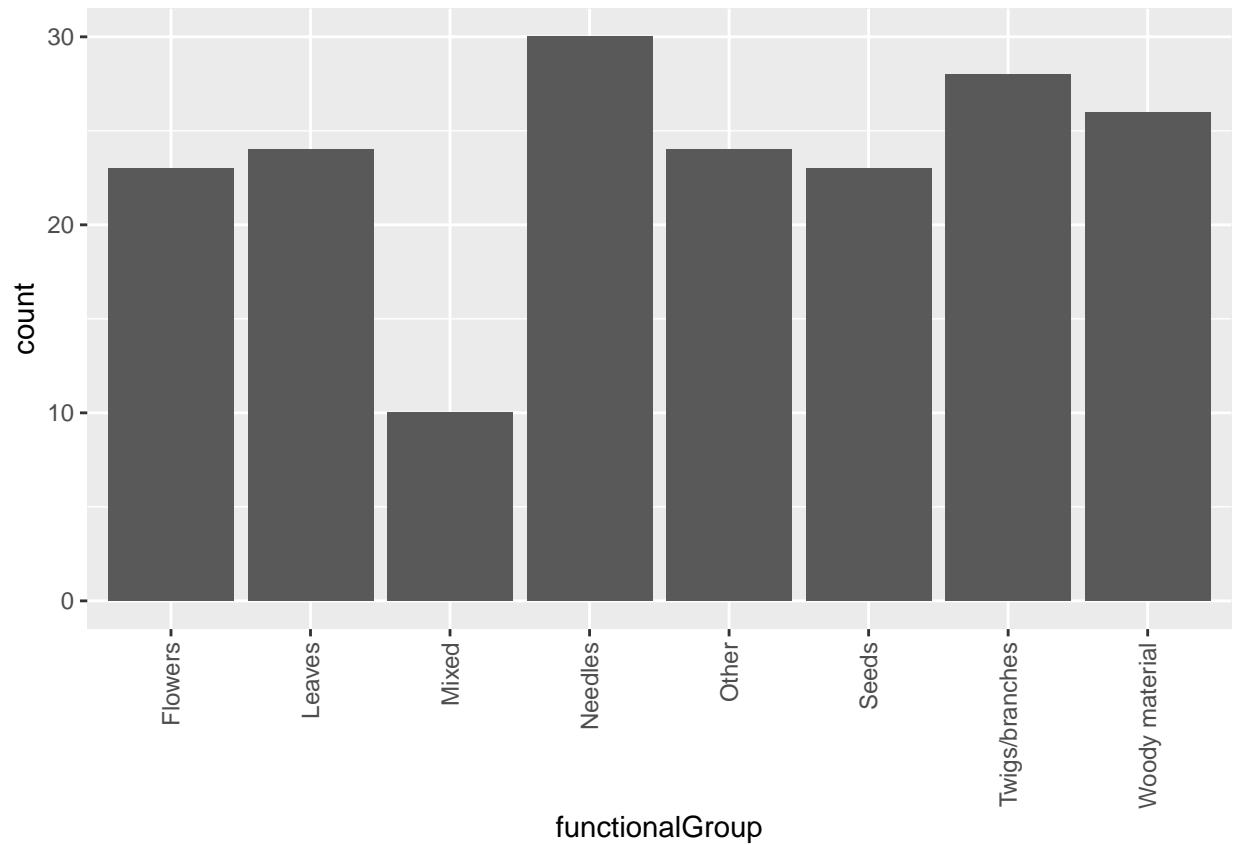
```
summary(Litter_dataset$plotID)
```

```
## 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:** The "unique" function returns the number of plots that were sampled at least once, in this case 12. On the other hand, the function "summary" returns the total number of plots that were sampled for each category.
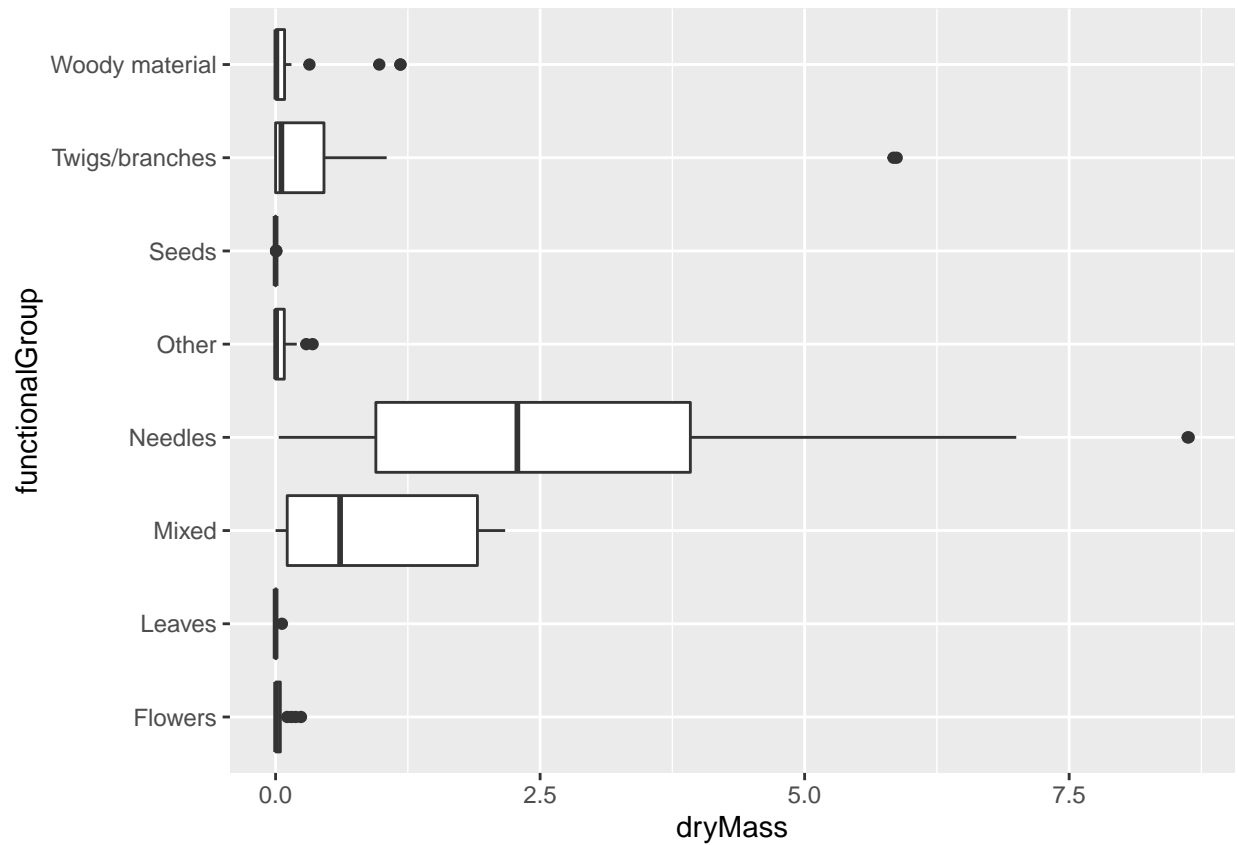
14. Create a bar graph of functionalGroup counts. This shows you what type of litter is collected at the Niwot Ridge sites. Notice that litter types are fairly equally distributed across the Niwot Ridge sites.

```
ggplot(Litter_dataset, aes(x = functionalGroup)) + geom_bar() +
    theme(axis.text.x = element_text(angle = 90, vjust = 0.5,
        hjust = 1))
```

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

```
ggplot(Litter_dataset) + geom_boxplot(aes(x = dryMass, y = functionalGroup))
```
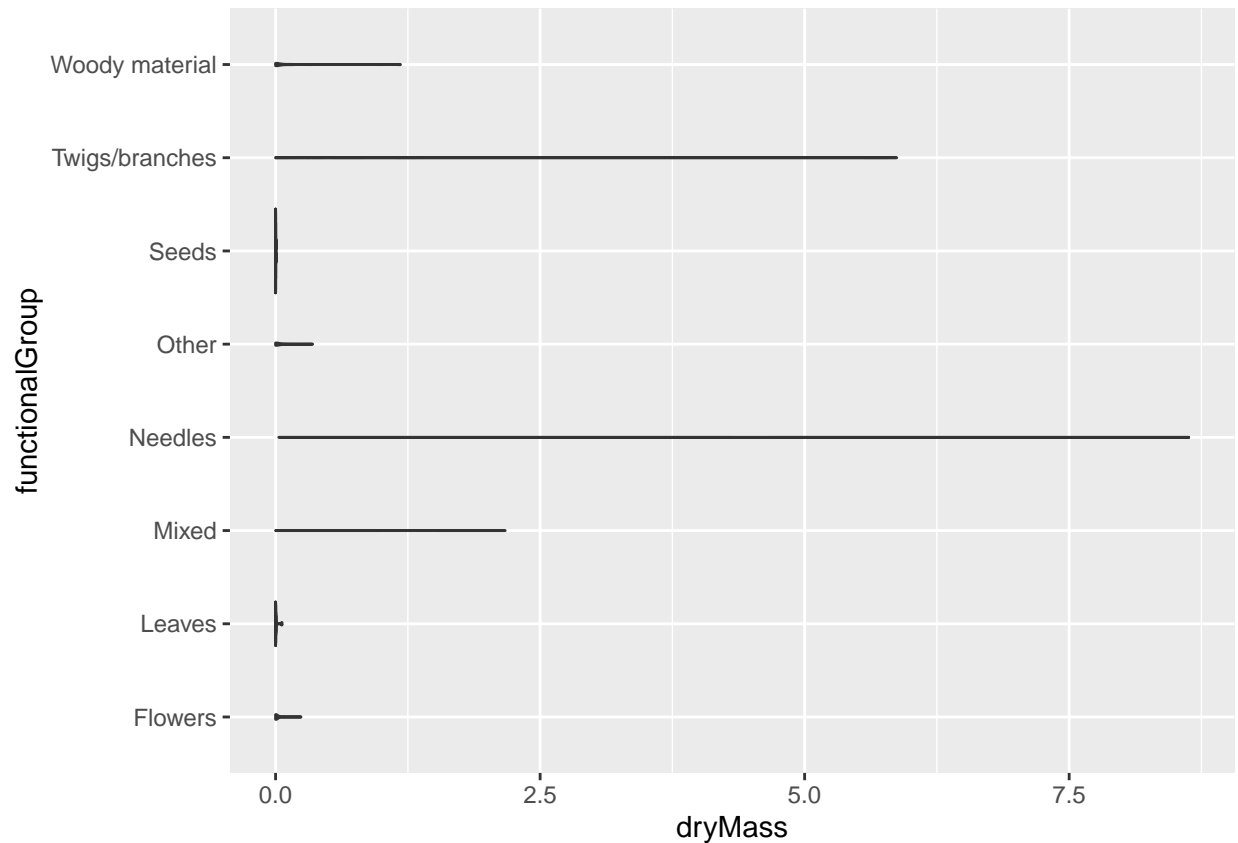
```
ggplot(Litter_dataset) + geom_violin(aes(x = dryMass, y = functionalGroup),
    draw_quantiles = c(0.25, 0.5, 0.75))
```

```
## Warning in regularize.values(x, y, ties, missing(ties), na.rm = na.rm):
## collapsing to unique 'x' values

## Warning in regularize.values(x, y, ties, missing(ties), na.rm = na.rm):
## collapsing to unique 'x' values

## Warning in regularize.values(x, y, ties, missing(ties), na.rm = na.rm):
## collapsing to unique 'x' values
```

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

**Answer:** The boxplot is a more effective visualization because this type of graph shows a box in which we find most of the observations disregarding the distribution inside the box, while the violin plot shows how they are effectively distributed. Given that litter types are distributed fairly equally, the violin plot only shows a line.

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

**Answer:** The types of litter that tend to have the highest biomass are "needles" and "mixed".