

Assignment 3: Data Exploration

Christopher Starr

Spring 2025

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.

```
#Read packages
library(tidyverse); library(lubridate); library(here); library(ggplot2)

#Read in data, don't convert strings to factors
Neonics <- read.csv(here('Data/Raw/ECOTOX_Neonicotinoids_Insects_raw.csv'), stringsAsFactors = TRUE)
# reading in as .csv
str(Neonics)
```

```
## 'data.frame':   4623 obs. of  30 variables:
## $ CAS.Number      : int  58842209 58842209 58842209 58842209 58842209 58842209 58842209 58842209
```

```
## $ Chemical.Name : Factor w/ 9 levels "(1E)-N-[(6-Chloro-3-pyridinyl)methyl]-N-ethy
## $ Chemical.Grade : Factor w/ 9 levels "Analytical grade",...: 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
## $ 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 ..
## $ 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
## $ Test.Location : Factor w/ 4 levels "Field artificial",...: 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
## $ 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 ...
## $ 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 103
## $ 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
```

```
#taking a look at the dataset
```

```
Litter <- read.csv(here("Data/Raw/NEON_NIWO_Litter_massdata_2018-08_raw.csv"), stringsAsFactors = TRUE)
# reading in as .csv
str(Litter)
```

```
## 'data.frame': 188 obs. of 19 variables:
## $ uid : Factor w/ 188 levels "028eea3d-5c20-4afc-bb7e-a05bab305152",...: 84 96 85
## $ namedLocation : Factor w/ 12 levels "NIWO_040.basePlot.ltr",...: 8 8 8 8 8 8 8 11 11 .
## $ domainID : Factor w/ 1 level "D13": 1 1 1 1 1 1 1 1 1 ...
## $ siteID : Factor w/ 1 level "NIWO": 1 1 1 1 1 1 1 1 1 ...
## $ plotID : Factor w/ 12 levels "NIWO_040","NIWO_041",...: 8 8 8 8 8 8 8 11 11 ...
## $ trapID : Factor w/ 12 levels "NIWO_040_205",...: 8 8 8 8 8 8 8 11 11 ...
## $ weighDate : Factor w/ 2 levels "2018-08-06","2018-09-05": 1 1 1 1 1 1 1 1 1 ...
## $ setDate : Factor w/ 2 levels "2018-07-05","2018-08-02": 1 1 1 1 1 1 1 1 1 ...
## $ collectDate : Factor w/ 2 levels "2018-08-02","2018-08-30": 1 1 1 1 1 1 1 1 1 ...
## $ ovenStartDate : Factor w/ 2 levels "2018-08-02T21:00Z",...: 1 1 1 1 1 1 1 1 1 ...
## $ ovenEndDate : Factor w/ 2 levels "2018-08-06T18:02Z",...: 1 1 1 1 1 1 1 1 1 ...
## $ fieldSampleID : Factor w/ 23 levels "NEON.LTR.NIWO040205.20180802",...: 14 14 14 14 14 14
## $ massSampleID : Factor w/ 168 levels "NEON.LTR.NIWO040205.20180802.FLR",...: 102 101 103
## $ samplingProtocolVersion: Factor w/ 1 level "NEON.DOC.001710vE": 1 1 1 1 1 1 1 1 1 ...
## $ functionalGroup : Factor w/ 8 levels "Flowers","Leaves",...: 7 6 8 1 8 4 5 2 1 8 ...
## $ dryMass : num 0.4 0.005 0.04 0.005 0.07 1 0.2 0.005 0.19 1.18 ...
## $ qaDryMass : Factor w/ 2 levels "N","Y": 1 1 2 1 1 1 1 1 2 ...
```

```
## $ remarks : logi NA NA NA NA NA NA ...
## $ measuredBy : Factor w/ 2 levels "kstyers@battelleecology.org",...: 1 1 1 1 1 1 1 1 1 1
```

```
#taking a look at the dataset
```

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: Likely we are looking for links between populations of certain insects and use of insecticides. Ideally we want to find a strong correlation between declines and use of insecticides so we can understand what insects are being killed by their use and how dramatically the population is harmed. Maybe then we can match this information with what we know about different classes of affected insects and their role in other areas of plant/animal life.

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: Likely we want to understand what trees and other forest plants are dropping to understand if their overall health is declining. For instance are they dropping more leaves and branches than before? Is there a decline in seed or other "reproductive" litter count or density?

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. The team set ground-level and above ground traps to collect material in areas with vegetation that is more than 2 meters tall. They measured the materials at different times depending on the collection method (above ground every 1-2 weeks, on ground once per year). 2. The ground and above-ground traps were paired together and were either placed randomly (in thick forest areas) or targeted (where the forest was not as thick). The sampling occurred in something called Tower Plots. 3. They measured up to .01g accuracy 8 types of vegetation (Leaves, needles, seeds, etc.).

Obtain basic summaries of your data (Neonics)

5. What are the dimensions of the dataset?

```
str(Neonics)
```

```
## 'data.frame': 4623 obs. of 30 variables:
## $ CAS.Number : int 58842209 58842209 58842209 58842209 58842209 58842209 58842209 58842209
## $ Chemical.Name : Factor w/ 9 levels "(1E)-N-[(6-Chloro-3-pyridinyl)methyl]-N-ethy
```

```
## $ Chemical.Grade : Factor w/ 9 levels "Analytical grade",...: 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
## $ 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
## $ 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 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
## $ Test.Location : Factor w/ 4 levels "Field artificial",...: 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
## $ 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
## $ 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 103
## $ 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
```

```
# Taking a look at the Neonics dataset
# 4632 observations with 30 variables each
```

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), decreasing = TRUE)
```

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

```
# using R to put the data in order for the Neonics file, Effect vector, decreasing by
# frequency.
```

Answer: Population (1803), Mortality (1493), and Behavior (360) are the largest categories. By measuring the litter we can see if there are more types of litter so we can understand if the population of a certain type of litter is increasing or decreasing. We can also check mortality, which might be checking to see if there are fully dead trees in the sample.

- 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), decreasing = TRUE)
```

##	(Other)	Honey Bee
##	670	667
##	Parasitic Wasp	Buff Tailed Bumblebee
##	285	183
##	Carniolan Honey Bee	Bumble Bee
##	152	140
##	Italian Honeybee	Japanese Beetle
##	113	94
##	Asian Lady Beetle	Euonymus Scale
##	76	75
##	Wireworm	European Dark Bee
##	69	66
##	Minute Pirate Bug	Asian Citrus Psyllid
##	62	60
##	Parastic Wasp	Colorado Potato Beetle
##	58	57
##	Parasitoid Wasp	Erythrina Gall Wasp
##	51	49
##	Beetle Order	Snout Beetle Family, Weevil
##	47	47
##	Sevenspotted Lady Beetle	True Bug Order
##	46	45
##	Buff-tailed Bumblebee	Aphid Family
##	39	38
##	Cabbage Looper	Sweetpotato Whitefly
##	38	37
##	Braconid Wasp	Cotton Aphid
##	33	33
##	Predatory Mite	Ladybird Beetle Family
##	33	30
##	Parasitoid	Scarab Beetle
##	30	29
##	Spring Tiphia	Thrip Order
##	29	29
##	Ground Beetle Family	Rove Beetle Family
##	27	27
##	Tobacco Aphid	Chalcid Wasp
##	27	25
##	Convergent Lady Beetle	Stingless Bee
##	25	25
##	Spider/Mite Class	Tobacco Flea Beetle

##	24	24
##	Citrus Leafminer	Ladybird Beetle
##	23	23
##	Mason Bee	Mosquito
##	22	22
##	Argentine Ant	Beetle
##	21	21
##	Flatheaded Appletree Borer	Horned Oak Gall Wasp
##	20	20
##	Leaf Beetle Family	Potato Leafhopper
##	20	20
##	Tooth-necked Fungus Beetle	Codling Moth
##	20	19
##	Black-spotted Lady Beetle	Calico Scale
##	18	18
##	Fairyfly Parasitoid	Lady Beetle
##	18	18
##	Minute Parasitic Wasps	Mirid Bug
##	18	18
##	Mulberry Pyralid	Silkworm
##	18	18
##	Vedalia Beetle	Araneoid Spider Order
##	18	17
##	Bee Order	Egg Parasitoid
##	17	17
##	Insect Class	Moth And Butterfly Order
##	17	17
##	Oystershell Scale Parasitoid	Hemlock Woolly Adelgid Lady Beetle
##	17	16
##	Hemlock Woolly Adelgid	Mite
##	16	16
##	Onion Thrip	Western Flower Thrips
##	16	15
##	Corn Earworm	Green Peach Aphid
##	14	14
##	House Fly	Ox Beetle
##	14	14
##	Red Scale Parasite	Spined Soldier Bug
##	14	14
##	Armoured Scale Family	Diamondback Moth
##	13	13
##	Eulophid Wasp	Monarch Butterfly
##	13	13
##	Predatory Bug	Yellow Fever Mosquito
##	13	13
##	Braconid Parasitoid	Common Thrip
##	12	12
##	Eastern Subterranean Termite	Jassid
##	12	12
##	Mite Order	Pea Aphid
##	12	12
##	Pond Wolf Spider	Spotless Ladybird Beetle
##	12	11
##	Glasshouse Potato Wasp	Lacewing

```
##              10              10
## Southern House Mosquito      Two Spotted Lady Beetle
##              10              10
##              Ant Family      Apple Maggot
##              9              9
```

asking R to put the data in order for the Neonics file, Species.Common.Name vector, decreasing by frequency.

```
summary(Neonics$Species.Common.Name, maxsum = 10)
```

```
## Honey Bee      Parasitic Wasp Buff Tailed Bumblebee
##      667      285      183
## Carniolan Honey Bee      Bumble Bee      Italian Honeybee
##      152      140      113
## Japanese Beetle      Asian Lady Beetle      Euonymus Scale
##      94      76      75
## (Other)
##      2838
```

asking R to show me the ten largest categories.

Answer: Most of the top 10 is made up of bee species. Bees are thought to have an outsized impact on ecosystems and pollination. The other names sound like parasites of interest. Perhaps their numbers are being closely examined to understand how bad their presence is in the area(s) being studied.

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"
```

asking R to tell me the class of Conc.1..Author.

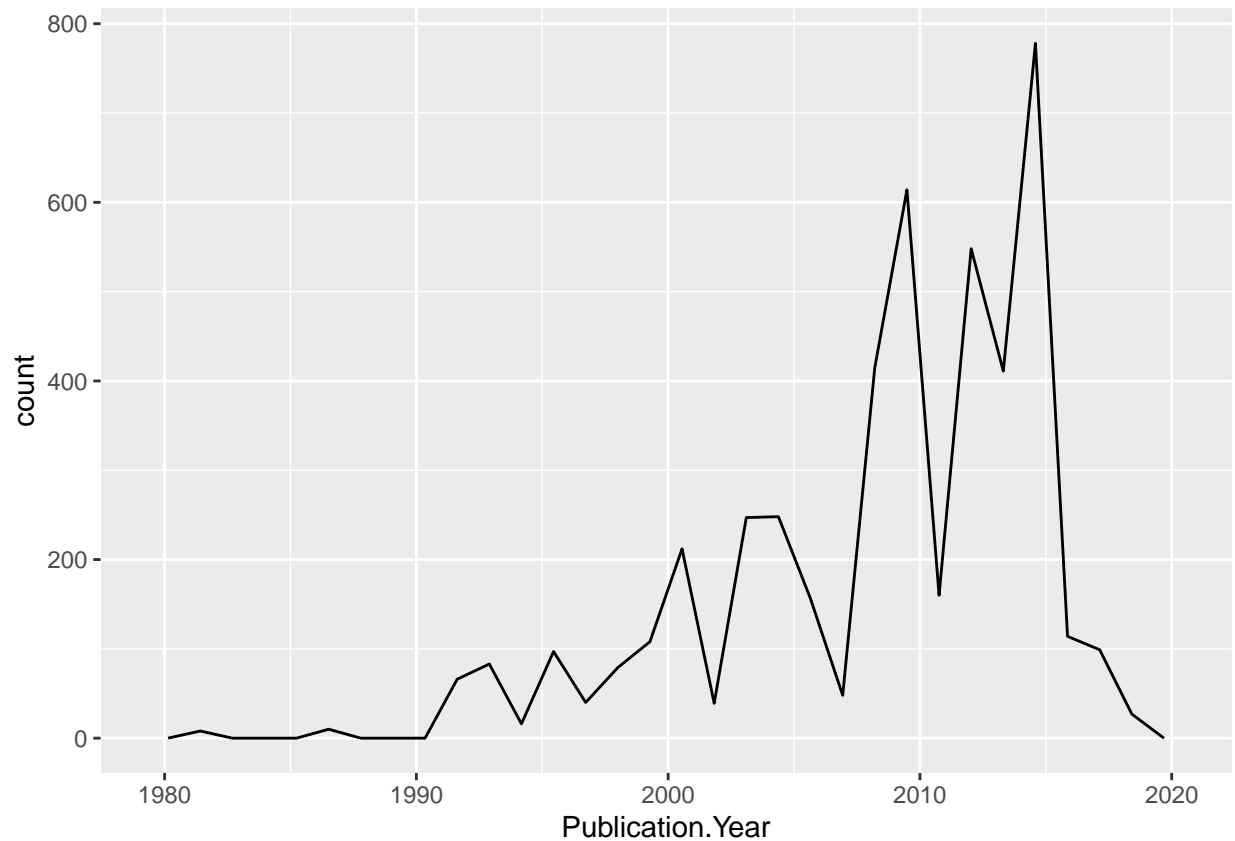
Answer: The vector is a Factor data type. It is not a number because there are a lot of special characters in the vector list (< and >, /, etc) which cannot be processed as numbers.

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()
```

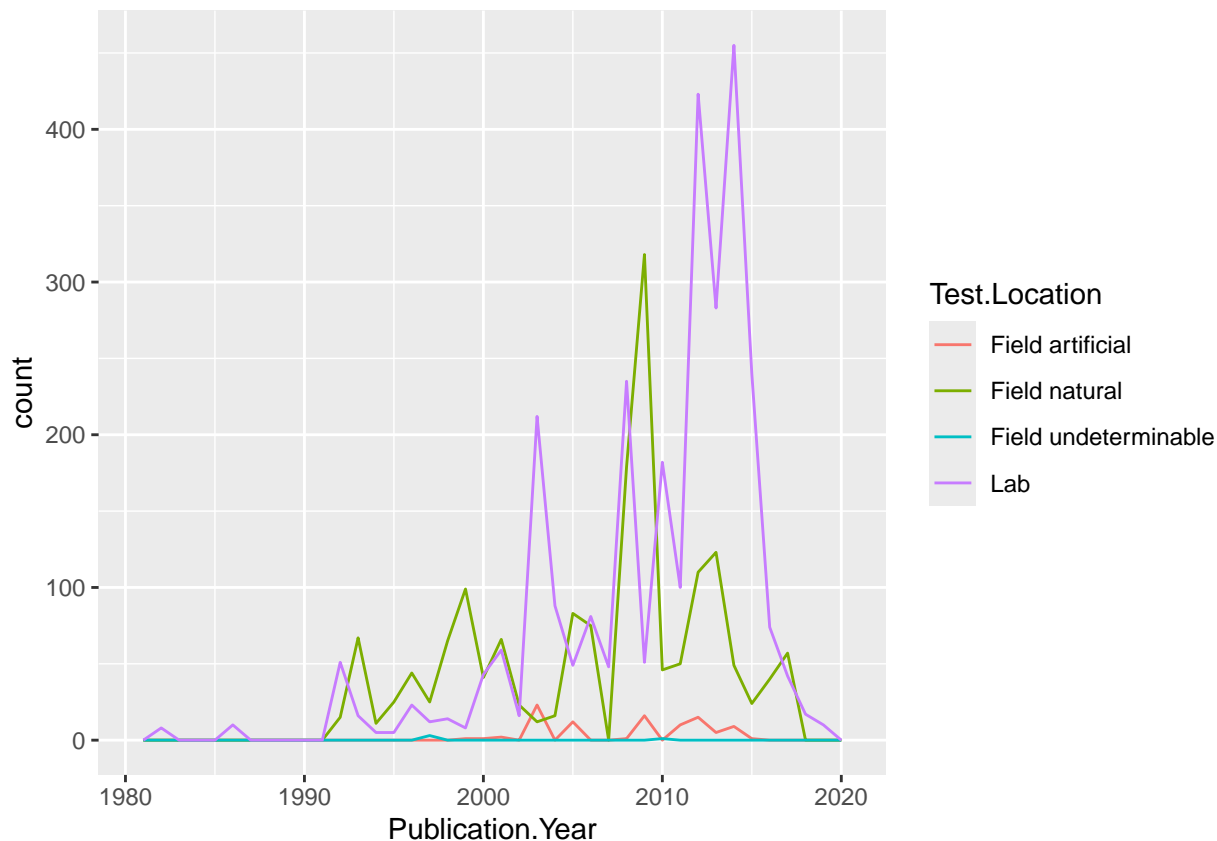
```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



asking R to generate a graph showing x as years and y as number of publications per year

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(binwidth = 1)
```

*# asking R to generate a graph showing years/studies per year with a different color line
to show different locations of where the data was gathered.*

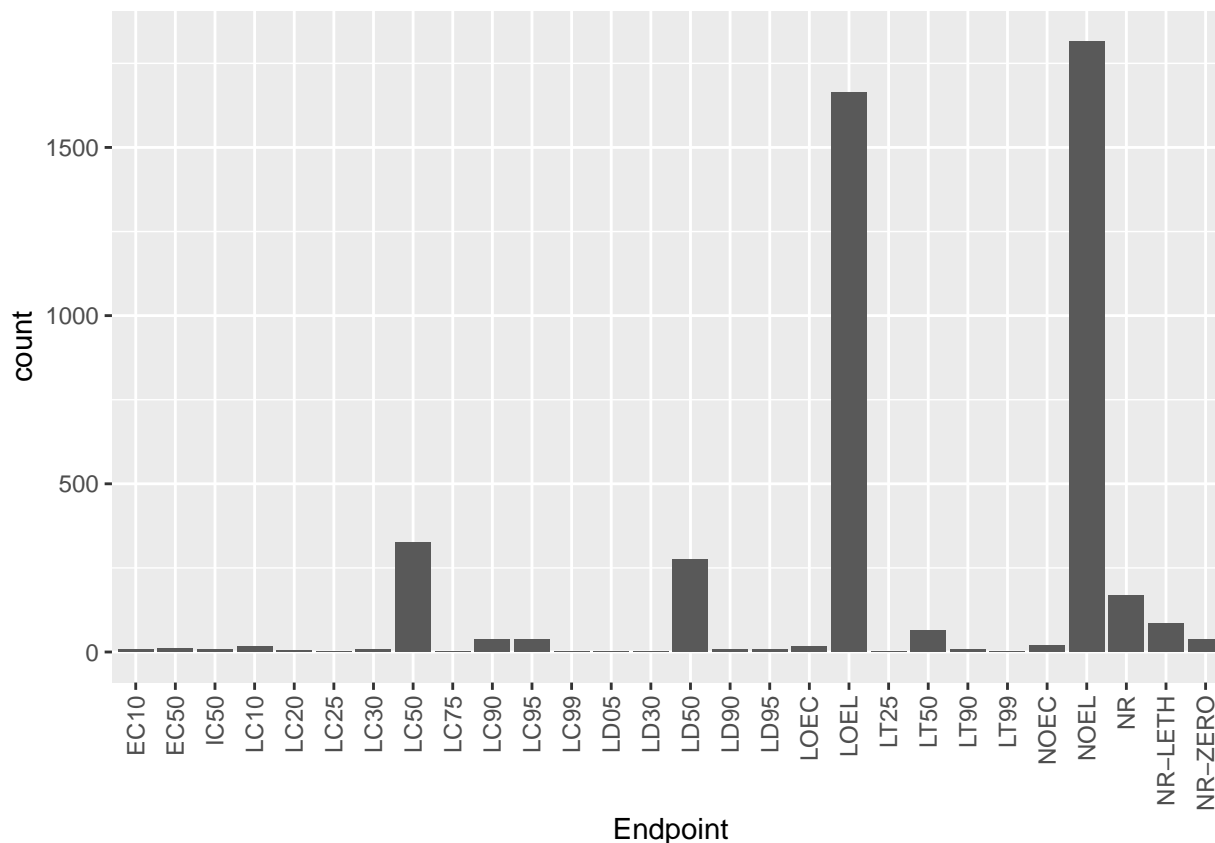
Interpret this graph. What are the most common test locations, and do they differ over time?

Answer: In the early 90s it was a close match between Field natural and Lab. In the late 90s Field natural became dominant. Throughout the 2000s it has been most Lab as the dominant location with the exception of 2009 when Field natural had a big spike.

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))
```



*# asking R to create a bar chart of Endpoints to show frequency. Adding theme elements
to help with aesthetics since the data is not easily visible without some adjustments.*

Answer: the most common endpoints are NOEL and LOEL. No-observable-effect-level: highest dose (concentration) producing effects not significantly different from responses of controls according to author's reported statistical test (NOEL/NOEC) Lowest Observed Effects Residue: The lowest residue concentration producing effects that were significantly different from responses of controls according to author's reported statistical test

Explore your data (Litter)

- 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"
```

asking R what collectDate is. Initial Response = "factor"

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

#asking R to use lubridate to to convert the column to date format.

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

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

```
#verifying the class changed. result = "Date"
```

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
```

```
#asking R to tell me how many plotID values there are  
sort(summary(Litter$plotID), decreasing = TRUE)
```

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

```
#asking R to show me the list of all plotIDs and their count, in descending order.
```

Answer: Unique returns a count of how many different values are in the vector and lists them. Summary returns the list of each vector AND the number of times each value occurs in the vector.

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

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

Answer:

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

Answer: