

# Lab 3

Math 241, Week 3

```
libs <- c('tidyverse','knitr','viridis', 'mosaic','mosaicData','babynames', 'Lahman','nycflights13','rnm
for(l in libs){
  if(!require(l,character.only = TRUE, quietly = TRUE)){
    message( sprintf('Did not have the required package << %s >> installed. Downloading now ... ',l))
    install.packages(l)
  }
  library(l, character.only = TRUE, quietly = TRUE)
}
```

**Due: Friday, February 16th at 8:30am**

## Goals of this lab

1. Practice creating functions.
2. Practice refactoring your code to make it better! Therefore for each problem, make sure to test your functions.

## Problem 1: Subset that R Object

Here are the R objects we will use in this problem (dats, pdxTreesSmall and ht).

```
library(pdxTrees)
library(mosaicData)

pdxTrees <- get_pdxTrees_parks()
# Creating the objects
dats <- list(pdxTrees = head(pdxTrees),
            Births2015 = head(Births2015),
            HELPrct = head(HELPrct),
            sets = c("pdxTrees", "Births2015",
                    "HELPrct"))

pdxTreesSmall <- head(pdxTrees)

ht <- head(pdxTrees$Tree_Height, n = 15)
```

- a. What are the classes of dats, pdxTreesSmall and ht?

Dats: List pdxTreesSmall: Dataframe ht: Vector (of numerics)

- b. Find the 10th, 11th, and 12th values of ht.

112, 112, 48

- c. Provide the Species column of pdxTrees as a data frame with one column.

I'm not sure whether I'm supposed to provide the whole column (including repeats), or whether I should get it and filter for uniqueness. From my understanding, I understood it as including repeats, but I want to make my intention clear for grading. I have also made this eval = F, due to how large this dataset is.

```
data.frame(pdxTrees$Species)
```

d. Provide the `Species` column of `pdxTrees` as a character vector.

I'm not sure whether I'm supposed to provide the whole column (including repeats), or whether I should get it and filter for uniqueness. From my understanding, I understood it as including repeats, but I want to make my intention clear for grading. I have also made this `eval = F`, due to how large this dataset is.

```
as.vector(pdxTrees$Species)
```

e. Provide code that gives us the second entry in `sets` from `dat`s.

```
dat$sets[2]
```

```
## [1] "Births2015"
```

f. Subset `pdxTreesSmall` to only Douglas-fir and then provide the `DBH` and `Condition` of the 4th Douglas-fir in the dataset. (Feel free to mix in some `tidyverse` code if you would like to.)

```
douglasfir <- pdxTreesSmall %>% filter(Common_Name == "Douglas-Fir") %>% select(DBH, Condition)
douglasfir[4,]
```

```
## # A tibble: 1 x 2
##   DBH Condition
##   <dbl> <chr>
## 1  32.1 Fair
```

## Problem 2: Function Creation

Figure out what the following code does and then turn it into a function. For your new function, do the following:

- Test it.
- Provide default values (when appropriate).
- Use clear names for the function and arguments.
- Make sure to appropriately handle missingness.
- Generalize it by allowing the user to specify a confidence level.
- Check the inputs and stop the function if the user provides inappropriate values.

```
library(pdxTrees)
thing1 <- length(pdxTrees$DBH)
thing2 <- mean(pdxTrees$DBH)
thing3 <- sd(pdxTrees$DBH)/sqrt(thing1)
thing4 <- qt(p = .975, df = thing1 - 1)
thing5 <- thing2 - thing4*thing3
thing6 <- thing2 + thing4*thing3
```

```
get_ci <- function(input, p = .975) {
  if (p > 1 | p < 0) {
    stop("Your confidence interval is not between 0 and 1.")
  }
  length <- length(input)
  mean <- mean(input)
  sd <- sd(input)/sqrt(length)
  moe <- qt(p = p, df = thing1 - 1)
  lower_ci <- mean - moe*sd
  upper_ci <- mean + moe*sd
  return(c(lower_ci, upper_ci))
}
```

```
}
get_ci(pdxTrees$DBH, 0.975)
```

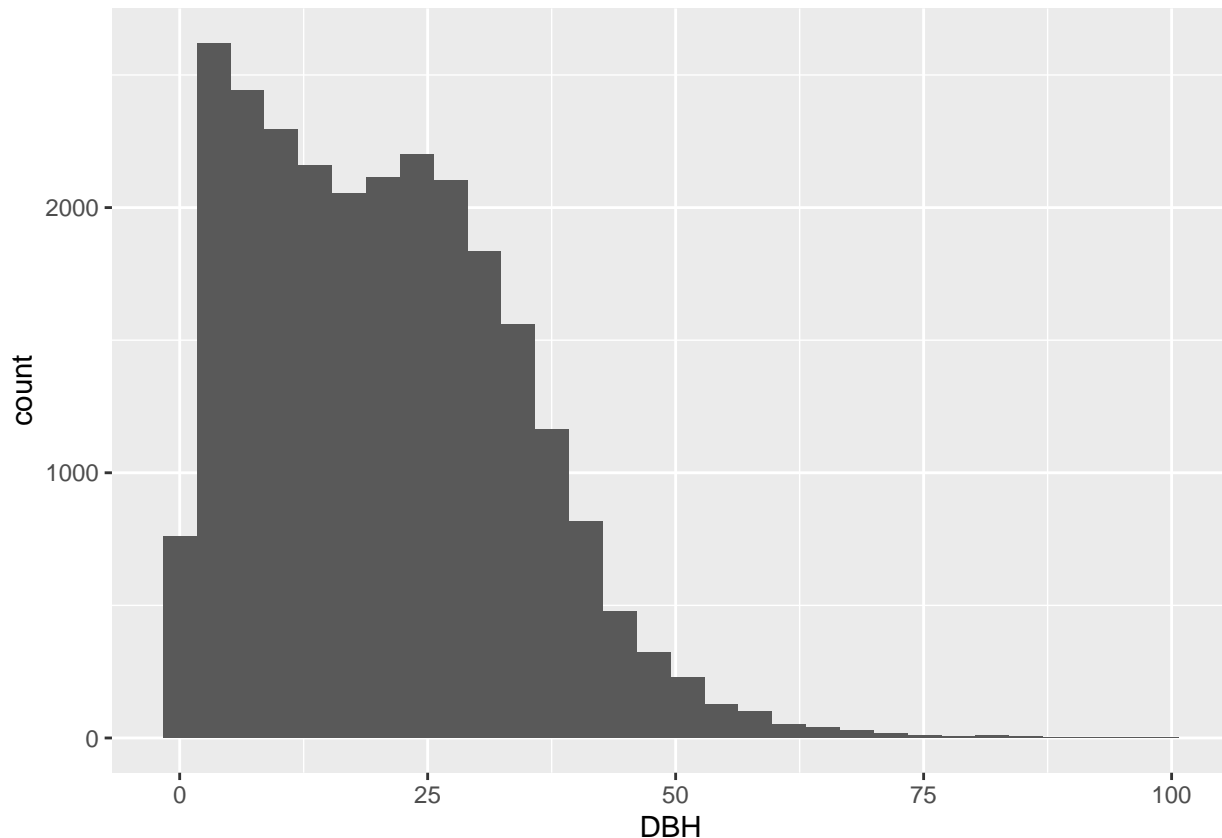
```
## [1] 20.44981 20.77835
```

### Problem 3: Wrapper Function for your ggplot

While we (i.e. Math 241 students) all love the grammar of graphics, not everyone else does. So for this problem, we are going to practice creating wrapper functions for `ggplot2`.

Here's an example of a wrapper for a histogram. Notice that I can't just list the variable name as an argument. The issue has to do with how many of the `tidyverse` functions evaluate the arguments. Therefore we have to quote (`enquo()`) and then unquote (`!!`) the arguments. (If you want to learn more, go [here](#).)

```
# Minimal viable product working code
ggplot(data = pdxTrees, mapping = aes(x = DBH)) +
  geom_histogram()
```



```
# Shorthand histogram function
histo <- function(data, x, bins, fill, color){
  x <- enquo(x)
  ggplot(data = data, mapping = aes(x = !!x)) +
    geom_histogram(bins = bins, fill = as.character(fill), color = as.character(color))
}

scatter <- function(data, x, y, color_var, alpha) {
  x <- enquo(x)
  y <- enquo(y)
```

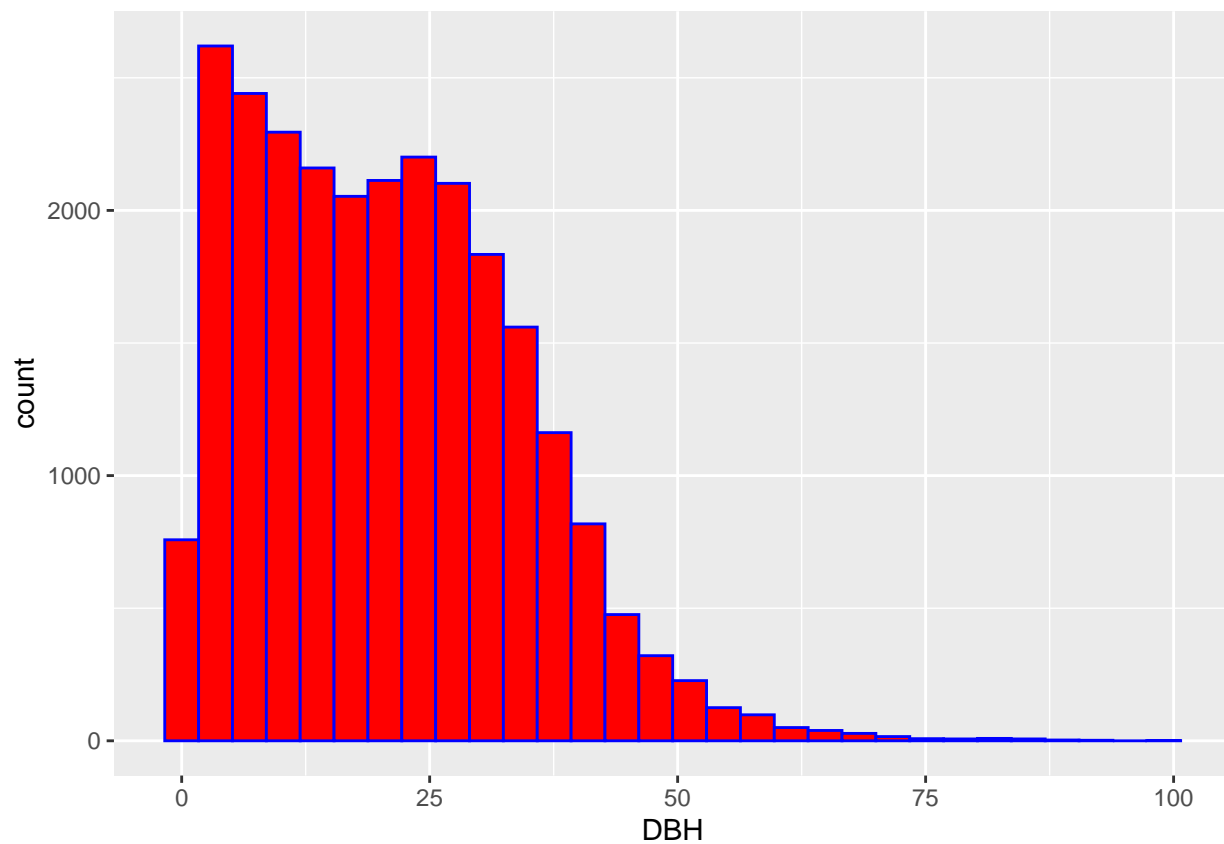
```

color_var <- enquos(color_var)
ggplot(data = data, mapping = aes(x = !!x, y = !!y, color = !!color_var)) +
  geom_point(alpha = alpha)
}

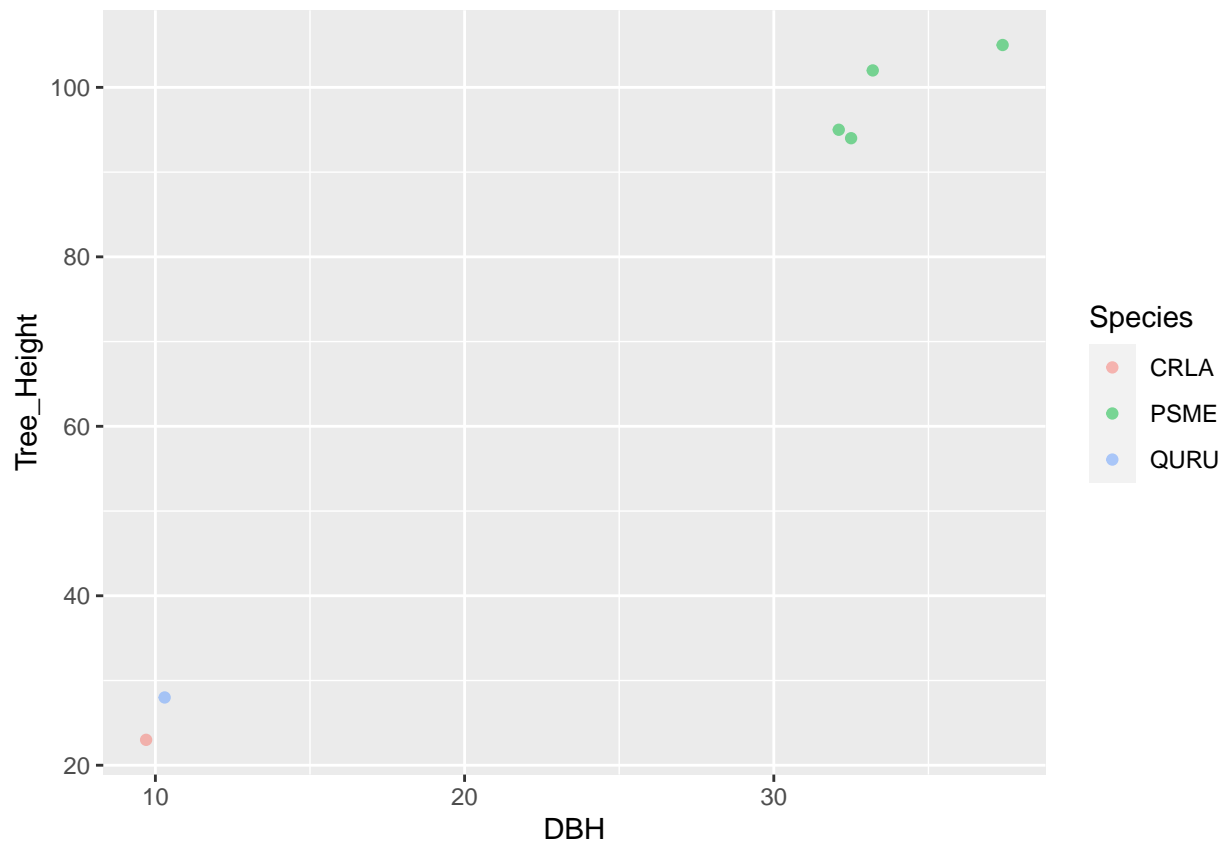
boxplot <- function(data, x, y, fill_var) {
  x <- enquos(x)
  y <- enquos(y)
  fill_var <- enquos(fill_var)
  ggplot(data = data, mapping = aes(x = !!x, y = !!y, fill = !!fill_var)) +
    geom_boxplot()
}

# Test it
histo(pdxTrees, DBH, 30, "red", "blue")

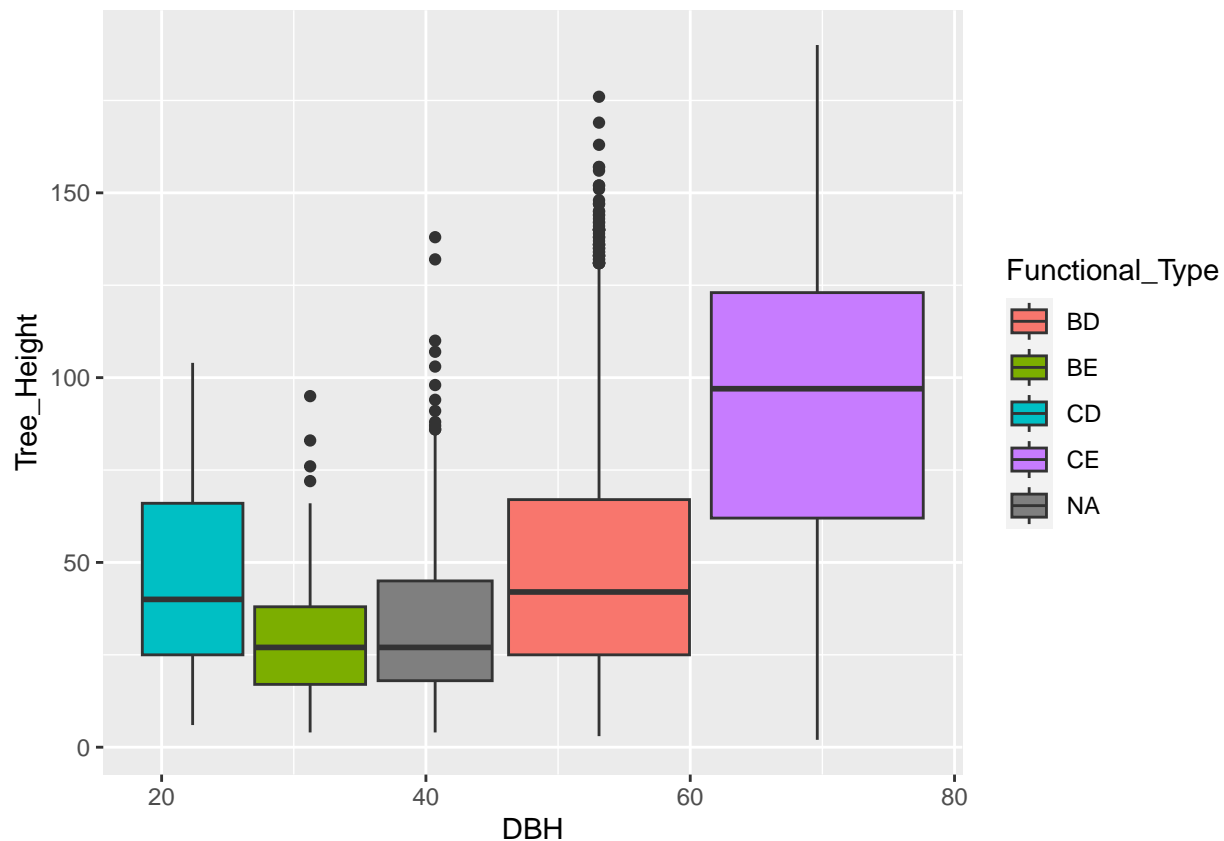
```



```
scatter(pdxTreesSmall, DBH, Tree_Height, Species, 0.5)
```



```
boxplot(pdxTrees, DBH, Tree_Height, Functional_Type)
```



- Edit `histo()` so that the user can set
  - The number of bins
  - The fill color for the bars
  - The color outlining the bars
- Write code to create a basic scatterplot with `ggplot2`. Then write and test a function to create a basic scatterplot.
- Modify your scatterplot function to allow the user to ...
  - Color the points by another variable.
  - Set the transparency.
- Write and test a function for your favorite `ggplot2` graph.

#### Problem 4: Functioning dplyr

- Take the following code and turn it into an R function to create a **conditional proportions** table. Similar to `ggplot2`, you will need to quote and unquote the variable names. Make sure to test your function!

```
pdxTrees %>%
  count(Native, Condition) %>%
  group_by(Native) %>%
  mutate(prop = n/sum(n)) %>%
  ungroup()
```

```
## # A tibble: 10 x 4
##   Native Condition     n  prop
##   <chr>   <chr>   <int> <dbl>
```

```
## 1 No Fair 12284 0.865
## 2 No Good 1043 0.0734
## 3 No Poor 875 0.0616
## 4 Yes Fair 9877 0.904
## 5 Yes Good 600 0.0549
## 6 Yes Poor 454 0.0415
## 7 <NA> Dead 264 0.658
## 8 <NA> Fair 118 0.294
## 9 <NA> Good 3 0.00748
## 10 <NA> Poor 16 0.0399
```

```
condition_prop <- function(data, input, condition) {
  input <- enquos(input)
  condition <- enquos(condition)
  data %>%
    count(!!input, !!condition) %>%
    group_by(!!input) %>%
    mutate(prop = n/sum(n)) %>%
    ungroup()
}
```

```
condition_prop(pdxTrees, Native, Condition)
```

```
## # A tibble: 10 x 4
##   Native Condition      n    prop
##   <chr>   <chr>   <int>  <dbl>
## 1 No     Fair    12284 0.865
## 2 No     Good     1043 0.0734
## 3 No     Poor      875 0.0616
## 4 Yes    Fair    9877 0.904
## 5 Yes    Good      600 0.0549
## 6 Yes    Poor      454 0.0415
## 7 <NA>    Dead     264 0.658
## 8 <NA>    Fair     118 0.294
## 9 <NA>    Good        3 0.00748
## 10 <NA>   Poor      16 0.0399
```

- b. Write a function to compute the mean, median, sd, min, max, sample size, and number of missing values of a quantitative variable by the categories of another variable. Make sure the output is a data frame (or tibble). Don't forget to test your function.

```
get_stats <- function(data, quant_var, category) {
  quant_var <- enquos(quant_var)
  category <- enquos(category)
  df <- data %>%
    group_by(!!category) %>%
    summarize(mean = mean(!!quant_var, na.rm = T),
              median = median(!!quant_var, na.rm = T),
              sd = sd(!!quant_var, na.rm = T),
              min = min(!!quant_var, na.rm = T),
              max = max(!!quant_var, na.rm = T),
              count = count(!!category),
              na_count = sum(is.na(!!quant_var)))
  print(df)
}
```

```
get_stats(pdxTrees, DBH, Common_Name)
```

```
## # A tibble: 304 x 8
##   Common_Name      mean median      sd   min   max count na_count
##   <chr>          <dbl>  <dbl>  <dbl> <dbl> <dbl> <int>   <int>
## 1 Accolade Elm      3.41    3.3   1.10   2.2    6.4    18      0
## 2 Alaska Yellow-Cedar 5.05    2.5   5.81   0.5   29.2    81      0
## 3 Aleppo Pine     12.0    11.4   2.94    9    15.7     6      0
## 4 Allegheny Serviceberry 4.22    4.1   0.933  2.7    5.7    21      0
## 5 American Beech   23.8    23.5   9.64   8.3   41.9    17      0
## 6 American Elm     29.6    30.1   9.89   1.5   58.7   379      0
## 7 American Hophornbeam 10.4     9.5   7.26   2.5   23.1     9      0
## 8 American Hornbeam, Blue Beech 11.2    11.9   6.19    1    20.7    53      0
## 9 American Linden  20.8    20.0   3.24  16.7   30.8    18      0
## 10 American Persimmon 2.9      2.9   NA      2.9    2.9     1      0
## # i 294 more rows
```

### Problem 5: another babynames exercise

Write a function called `grab_name` that, when given a **name** *and a year* as an argument, returns the rows from the `babynames` data frame in the `babynames` package that match that name for that year (and returns an error if that name and year combination does not match any rows). Run the function once with the arguments **Ezekiel** and **1883** and once with **Ezekiel** and **1983**.

```
#' Make sure to switch eval = FALSE to eval = TRUE before knitting!!
```

```
grab_name <- function(myname, myyear){
  # Code your function here
  myname <- enquo(myname)
  myyear <- enquo(myyear)
  temp_data <- filter(babynames, name == !!myname, year == !!myyear)
  if (nrow(temp_data) < 1)
    {stop("At the year ", quo_name(myyear), ", there are no matches for the name: ", quo_name(myname))}
  else
    {return(temp_data)}
}

grab_name("Ezekiel", 1983)
```

```
## # A tibble: 1 x 5
##   year sex  name      n      prop
##   <dbl> <chr> <chr>  <int>   <dbl>
## 1  1983 M    Ezekiel  149 0.0000800
```