

manipulating and analyzing data with dplyr



the tidyverse: an "umbrella" package

- **ggplot2**: a "grammar of graphics" by Hadley Wickham. Divide the data and the aesthetics. Create and modify the plots layer by layer
- **dplyr**: a way to lead with data frames, sql external data bases, written in **C++**
- **readr**: read data
- **tidyr**: format data frames
- **stringr**: deals with strings
- additional packages for other tasks: **tibble**, **lubridate** and many more

Most of R is still **base**-based and both philosophies communicate well with each other

reading data with readr

```
library(dplyr)
library(readr)
```

```
surveys <- readr::read_csv("data_raw/portal_data_joined.csv")
```

```
## Parsed with column specification:
## cols(
##   record_id = col_double(),
##   month = col_double(),
##   day = col_double(),
##   year = col_double(),
##   plot_id = col_double(),
##   species_id = col_character(),
##   sex = col_character(),
##   hindfoot_length = col_double(),
```

```
## inspect the data  
str(surveys)
```

View(surveys)

surveys

some principal functions in dplyr

- **select** (columns)
- **filter** (rows)
- **rename** (columns)
- **mutate** (create new columns or modify existing columns)
- **arrange** to sort according to a column
- **count** cases of one or many columns

select columns

```
select(surveys, plot_id, species_id, weight)
```

1. there is no need to put quotes
2. there is no need to put variables between `c()`

base R still works in a tibble

```
surveys[, c("plot_id", "species_id", "weight")]
```

removing columns

```
select(surveys, -record_id, -species_id)
```


additional functions

```
select(surveys, -ends_with("id"))
```

filter rows

logical clauses!

```
surv_1995 <- filter(surveys, year == 1995)
```

No need to use \$ or brackets

```
surveys$year == 1995  
surveys[surveys$year == 1995 , ]
```

mutate creates or modifies columns

```
surveys <- mutate(surveys, weight_kg = weight / 1000)
```

```
mutate(surveys,  
       weight_kg = weight / 1000,  
       weight_lb = weight_kg * 2.2)
```

group_by() and summarise()

- if you have a column factor (e.g. sex) and want to apply a function to the levels of this factor

```
surveys_g <- group_by(surveys, sex) #does nothing?
summary_sex <- summarize(surveys_g,
                          mean_weight = mean(weight, na.rm = TRUE))
summary_sex
```

```
## # A tibble: 3 x 2
##   sex    mean_weight
##   <chr>      <dbl>
## 1 F         42.2
## 2 M         43.0
## 3 <NA>      64.7
```

another example:

```
surveys_g2 <- group_by(surveys, sex, species_id)
mean_w <- summarize(surveys_g2,
                     mean_weight = mean(weight, na.rm = TRUE))
```

```
mean_w
```

```
## # A tibble: 92 x 3
## # Groups:   sex [3]
##   sex    species_id mean_weight
##   <chr> <chr>         <dbl>
## 1 F      BA           9.16
## 2 F      DM          41.6
## 3 F      DO          48.5
## 4 F      DS         118.
## 5 F      NL         154.
## 6 F      OL          31.1
## 7 F      OT          24.8
## 8 F      OX           21
## 9 F      PB          30.2
## 10 F     PE          22.8
## # ... with 82 more rows
```

arrange sorts by a column

```
arrange(mean_w, mean_weight)
```

```
arrange(mean_w, desc(mean_weight))
```

the pipe operator



Classic syntax goes like this

```
object1  
object2 <- function1(object1)  
object3 <- function2(object2)
```

The pipe operator allows to apply functions sequentially:

```
object3 <- object1 %>% function1() %>% function2()
```

functions in the tidyverse work very well with pipes

select and filter

```
surveys2 <- filter(surveys, weight < 5)  
surveys_sml <- select(surveys2, species_id, sex, weight)
```

```
surveys %>%  
  filter(weight < 5) %>%  
  select(species_id, sex, weight)
```

group_by() and summarize()

```
surveys_g    <- group_by(surveys, sex) #does nothing?  
summary_sex <- summarize(surveys_g,  
                          mean_weight = mean(weight, na.rm = TRUE))
```

```
summary_sex <- surveys %>%  
  group_by(sex) %>%  
  summarize(mean_weight = mean(weight, na.rm = TRUE))
```

count

```
surveys %>%  
  count(sex)
```

```
surveys %>%  
  count(sex, species)
```

```
surveys %>%  
  count(sex, species) %>%  
  arrange(species, desc(n))
```

challenge

- How many animals were caught in each `plot_type` surveyed?
- Use `group_by()` and `summarize()` to find the mean, min, and max hindfoot length for each species (using `species_id`). Also add the number of observations (hint: see `?n`).

save data!

```
surveys <- readr::read_csv("data_raw/portal_data_joined.csv")

surveys_complete <- surveys %>%
  filter(!is.na(weight),
         !is.na(hindfoot_length),
         !is.na(sex))

species_counts <- surveys_complete %>%
  count(species_id) %>%
  filter(n >= 50)

surveys_complete <- surveys_complete %>%
  filter(species_id %in% species_counts$species_id)

write_csv(surveys_complete, path = "data/surveys_complete.csv")
```

data visualization with ggplot2

ggplot2

- **ggplot2** separates the data from the aesthetics part and allows layers of information to be added sequentially with **+**

```
ggplot(data = <data>,  
       mapping = aes(<mappings>)) +  
  geom_xxx()
```

- **data**
- **mappings**: the specific variables (x, y, z, group...)
- **geom_xxx()**: functions for plotting options **geom_point()**, **geom_line()**

[cheat sheet link](#)

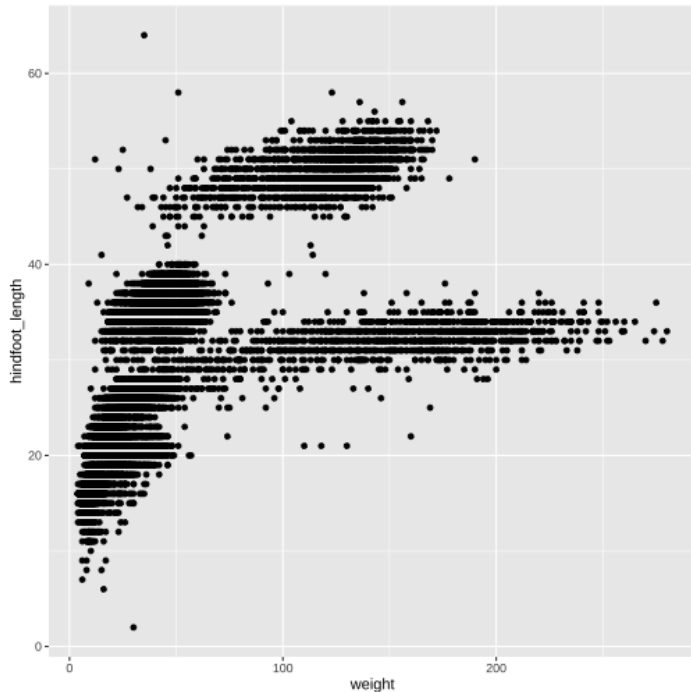
ggplot2 plots are built sequentially in layers

```
library(ggplot2)
library(readr)

surveys_complete <- read_csv("data/surveys_complete.csv")
```


ggplot2 plots are built sequentially in layers

```
ggplot(data = surveys_complete,           # data
       mapping = aes(x = weight, y = hindfoot_length)) + # aesthetics
  geom_point()                             # plot function
```

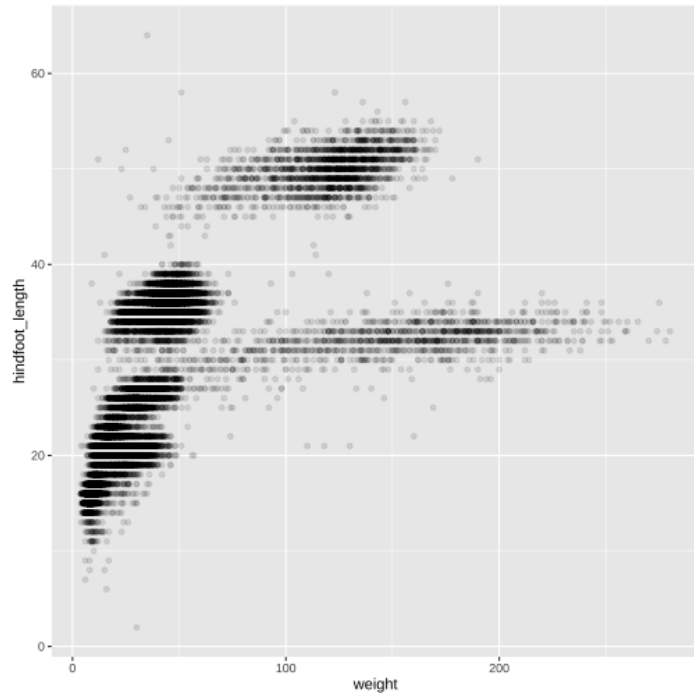


you can assign a plot to an object and build on it

```
surveys_plot <- ggplot(data = surveys_complete,  
                        mapping = aes(x = weight,  
                                      y = hindfoot_length))  
  
surveys_plot +  
  geom_point()
```

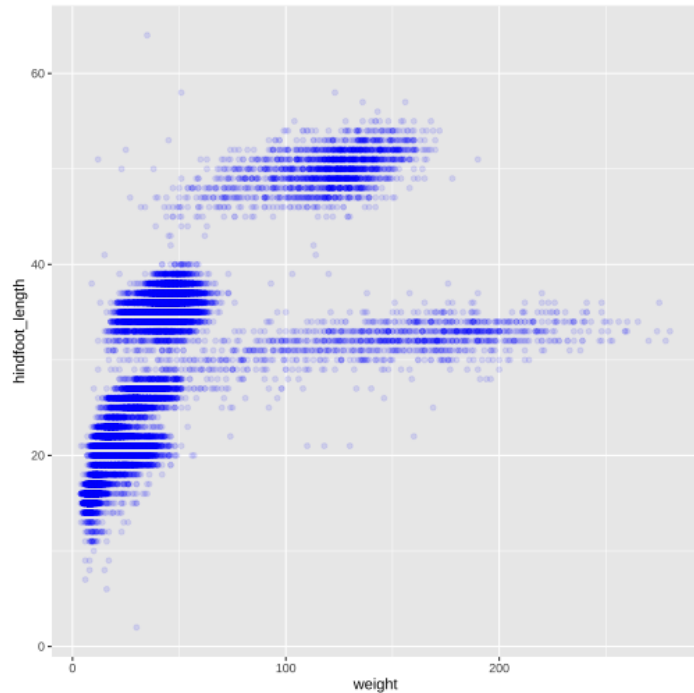
ggplot2 plots are built sequentially in layers

```
surveys_plot +  
  geom_point(alpha = 0.1) #transparency
```



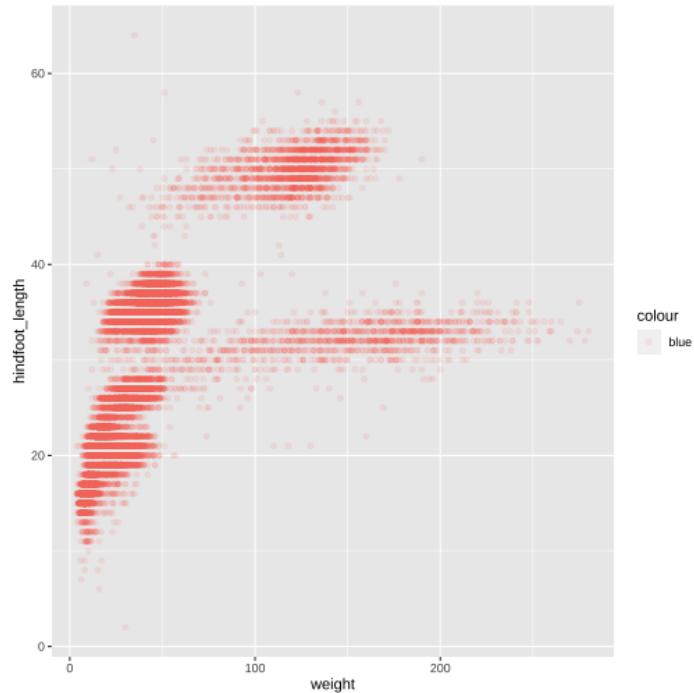
ggplot2 plots are built sequentially in layers

```
surveys_plot +  
  geom_point(alpha = 0.1, color = "blue") #color
```



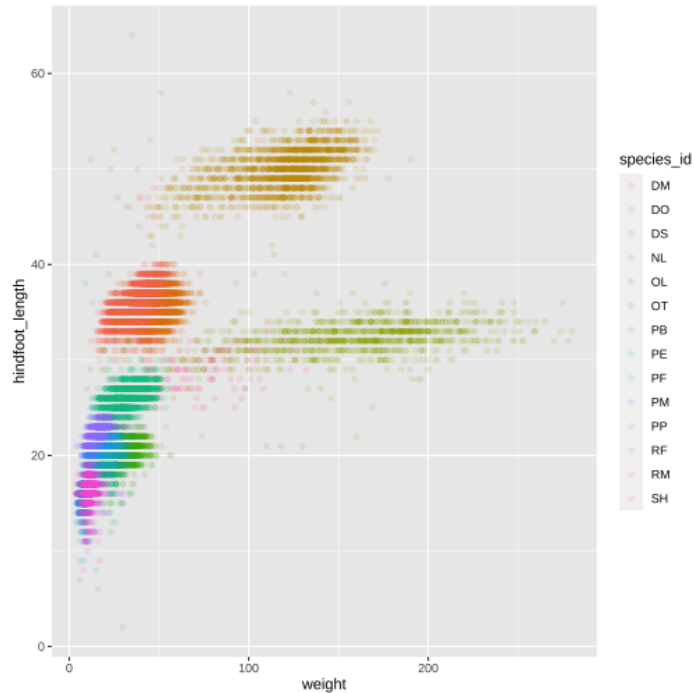
ggplot2 plots are built sequentially in layers

```
surveys_plot +  
  geom_point(alpha = 0.1, aes(color = "blue")) #this is a mistake!
```



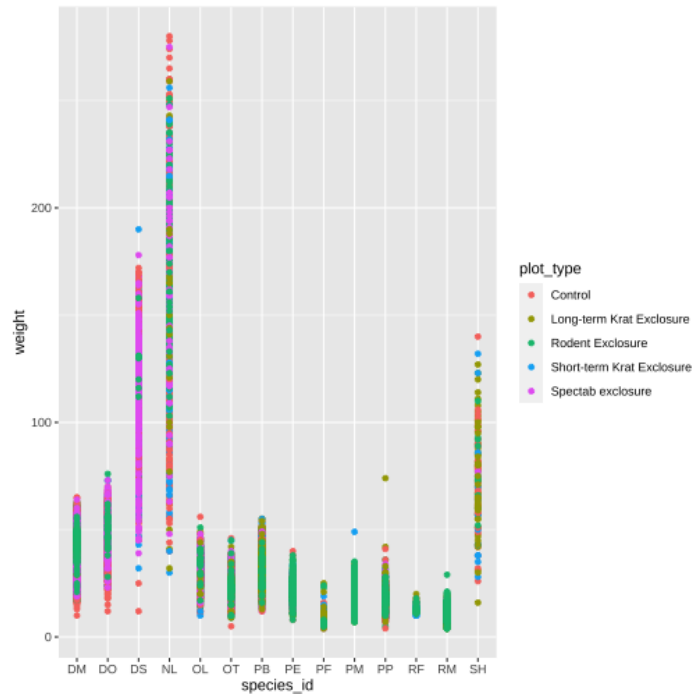
ggplot2 plots are built sequentially in layers

```
surveys_plot +  
  geom_point(alpha = 0.1, aes(color = species_id))
```



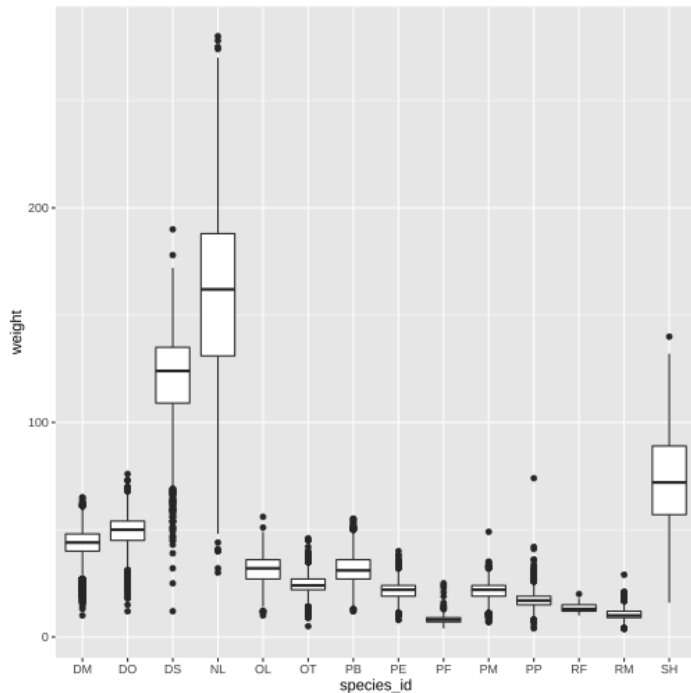
challenge: change x to categorical variable

```
ggplot(data = surveys_complete,  
       mapping = aes(x = species_id, y = weight)) +  
  geom_point(aes(color = plot_type))
```



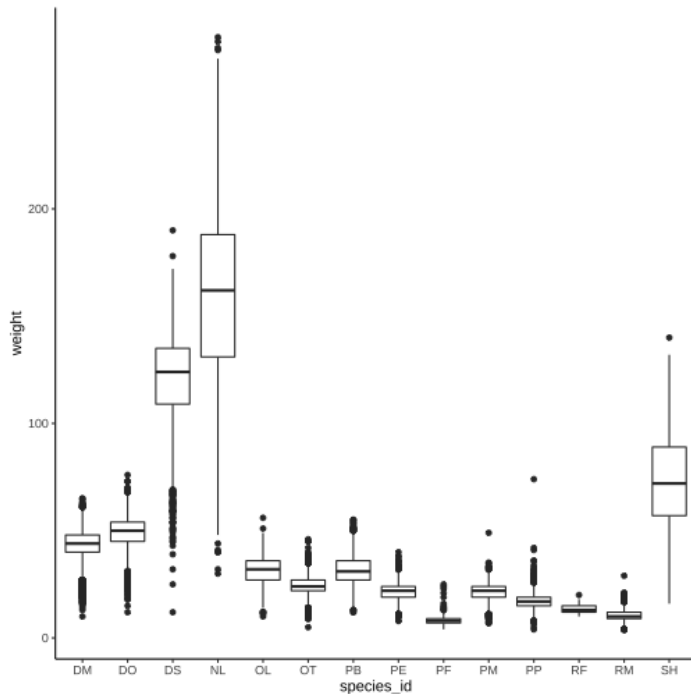
boxplots!

```
# boxplots  
ggplot(data = surveys_complete,  
       mapping = aes(x = species_id, y = weight)) +  
  geom_boxplot()
```



theme options theme_

```
ggplot(data = surveys_complete,  
       mapping = aes(x = species_id, y = weight)) +  
  geom_boxplot() +  
  theme_classic()
```



add jitter layer

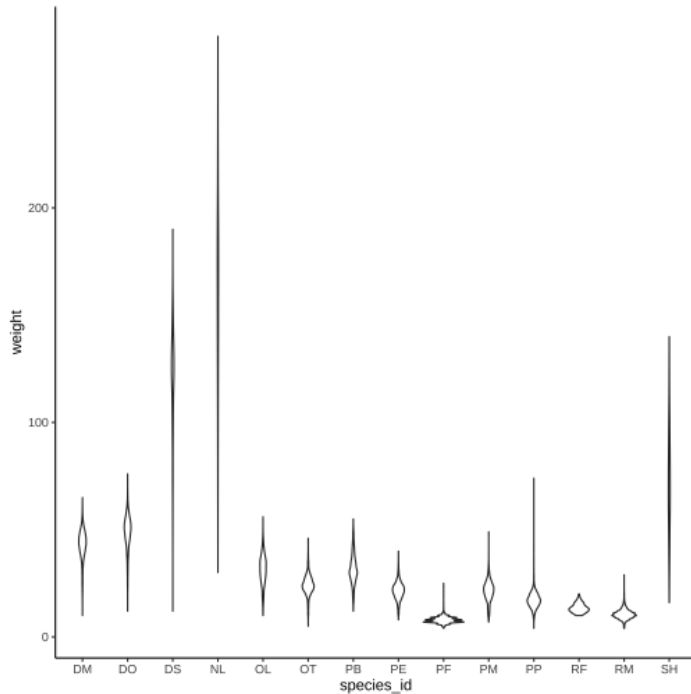
```
ggplot(data = surveys_complete,  
       mapping = aes(x = species_id, y = weight)) +  
  geom_boxplot() +  
  geom_jitter(alpha = 0.3, color = "dodgerblue", width = 0.2) +  
  theme_classic()
```

change plot order

```
ggplot(data = surveys_complete,  
       mapping = aes(x = species_id, y = weight)) +  
  geom_jitter(alpha = 0.3, color = "dodgerblue", width = 0.2) +  
  geom_boxplot() +  
  theme_classic()
```

violin plots

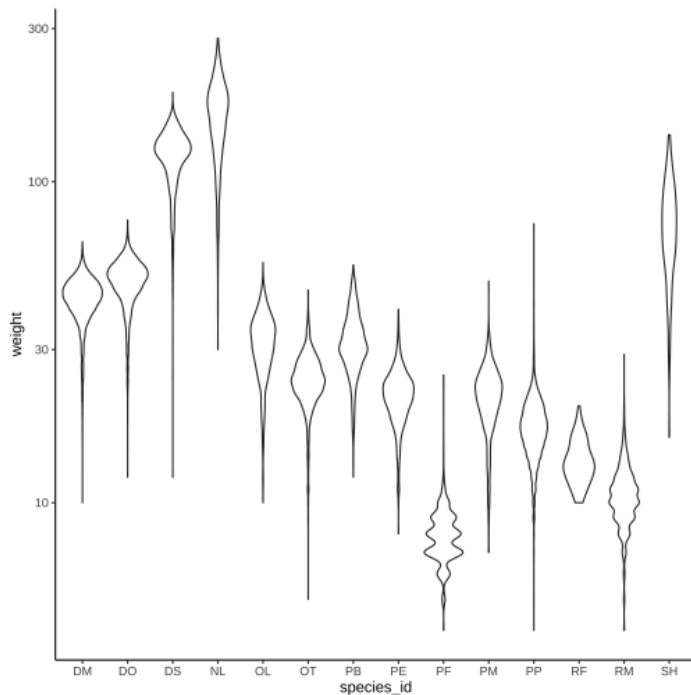
```
ggplot(data = surveys_complete,  
       mapping = aes(x = species_id, y = weight)) +  
  geom_violin() + theme_classic()
```



change scale (scale_xx options)

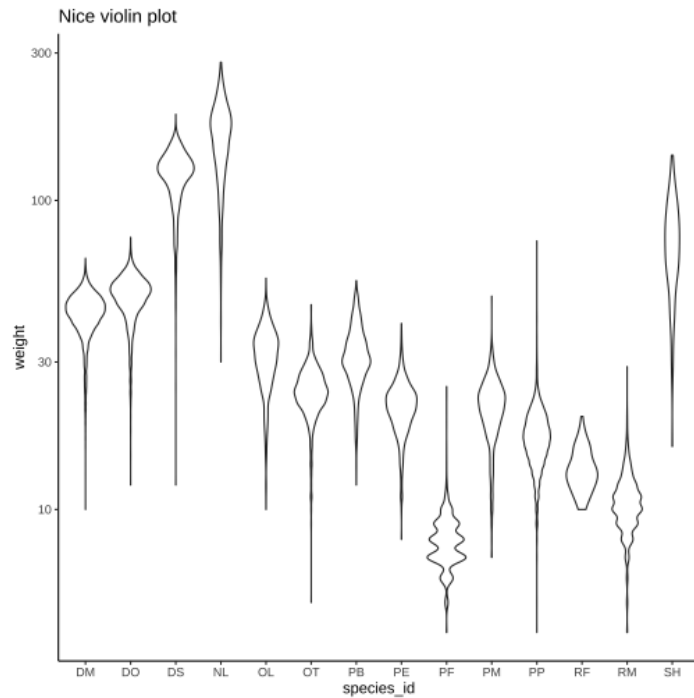
```
p <- ggplot(data = surveys_complete,  
  mapping = aes(x = species_id, y = weight)) +  
  geom_violin() + scale_y_log10() + theme_classic() #nice!
```

p



add title ggtitle()

```
p + #remember the plot can be an object  
  ggtitle("Nice violin plot")
```



ggplot2 plots are built sequentially in layers

```
ggplot(data = surveys_complete,  
       mapping = aes(x = species_id, y = hindfoot_length)) +  
  geom_jitter(size = 0.5, alpha = 0.1, width = 0.2, aes(col = plot_id)) +  
  geom_boxplot() + scale_y_log10() + theme_classic() + ggtitle("Nice violin plot")
```

