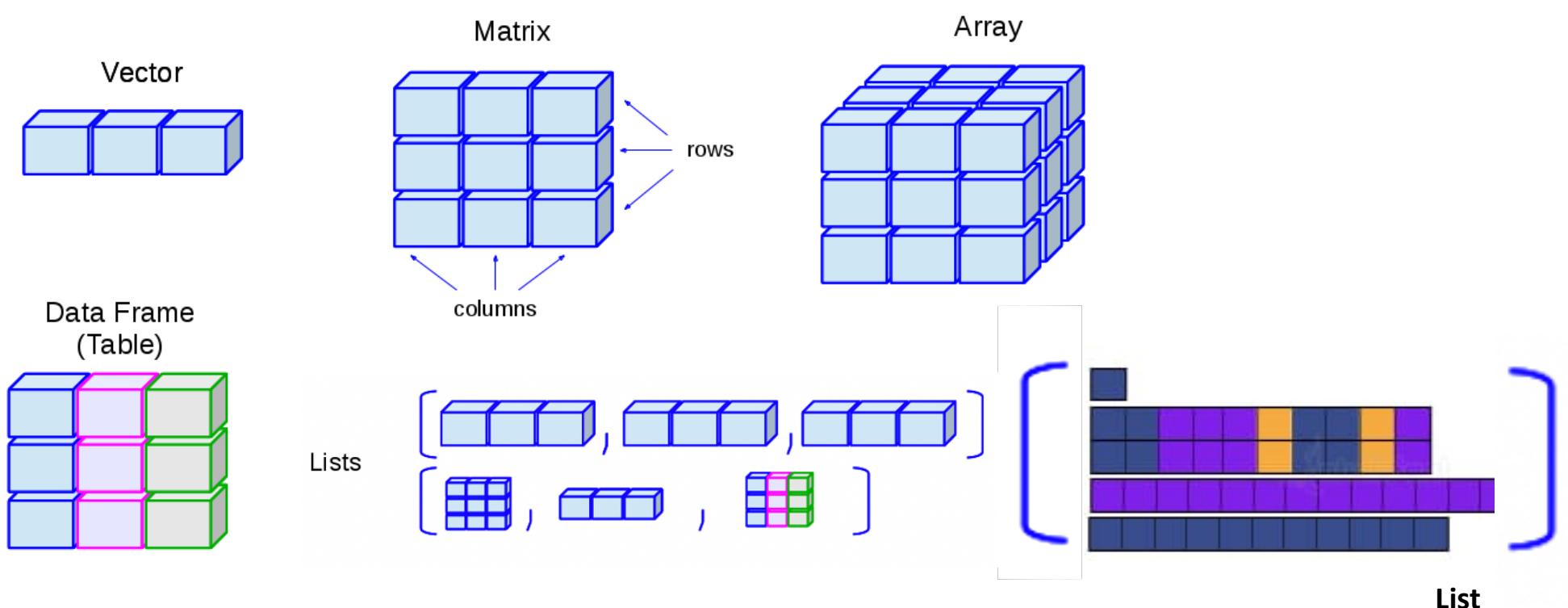


# Introduction to for Biologists

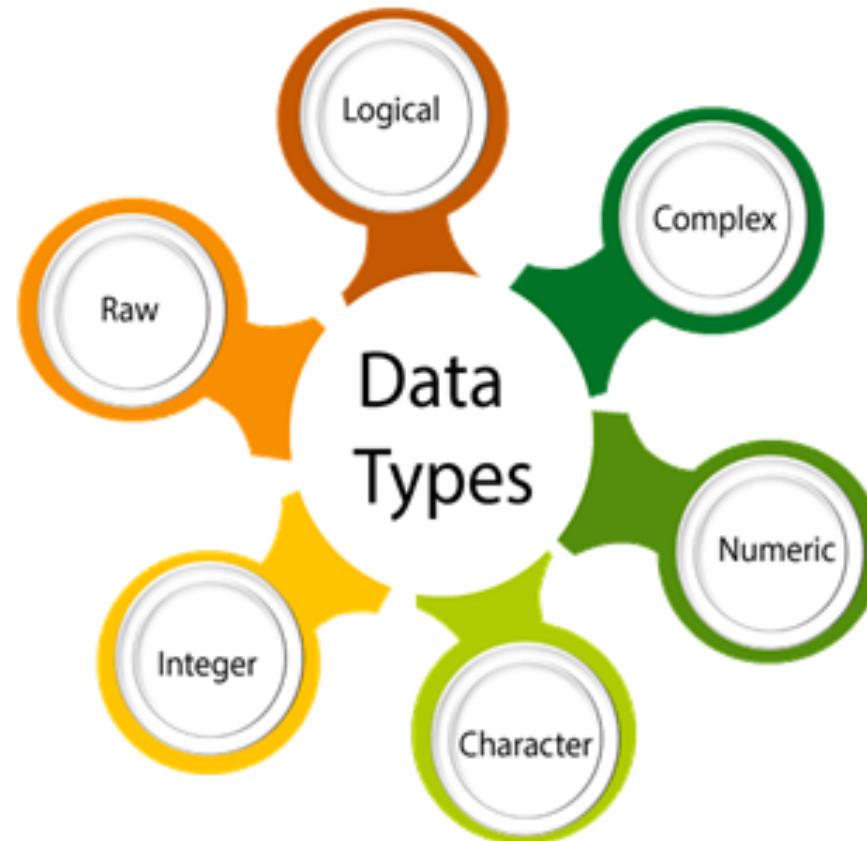
WORKFLOW & VISUALISATION

Author: Gita Yadav

# Data Structures (Recap)

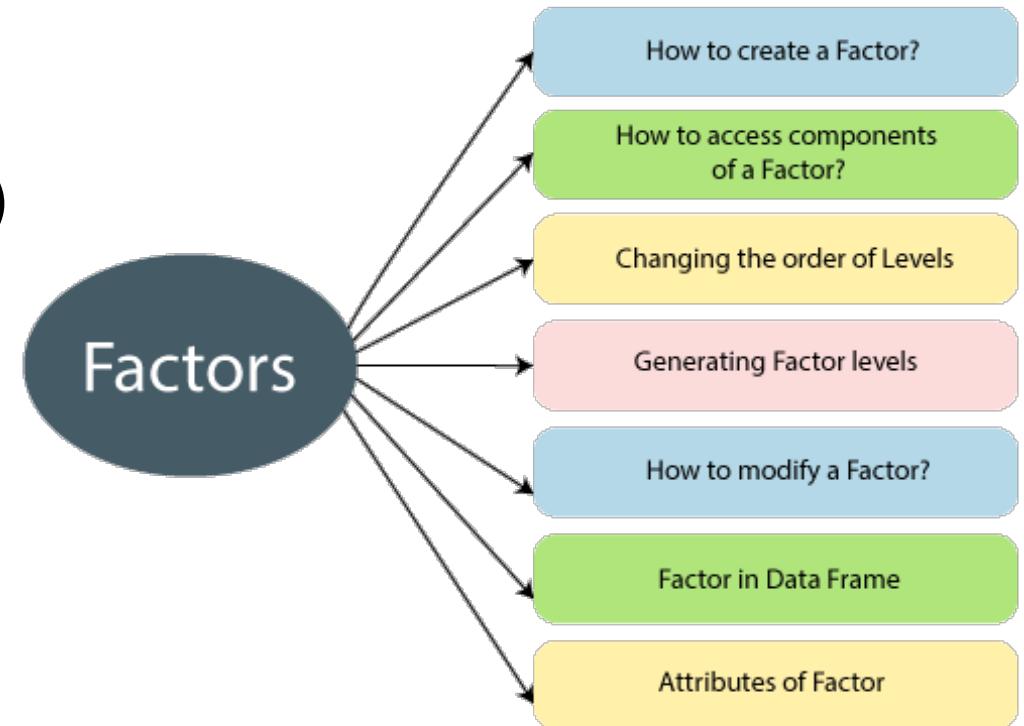


# Data Type (Recap)

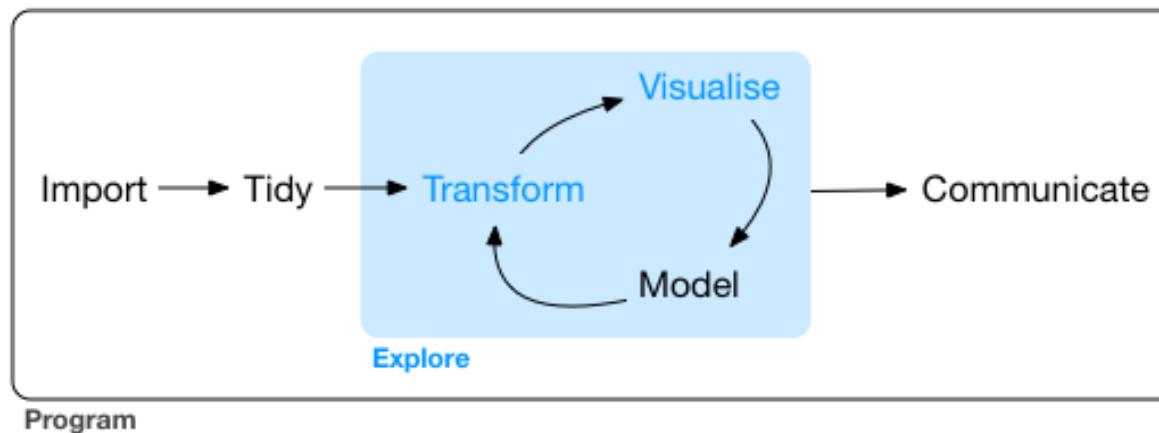


# Factors

- Categorical Variables in Statistics
- Can take Limited set of Values (levels)
  - Example:
    - “Gender” = {Male, Female}
    - “Meal” = {Breakfast, Lunch, Dinner}
  - No Intrinsic Ordering (alphabetical)
  - Order can be changed
    - Why/When would you wish this!?



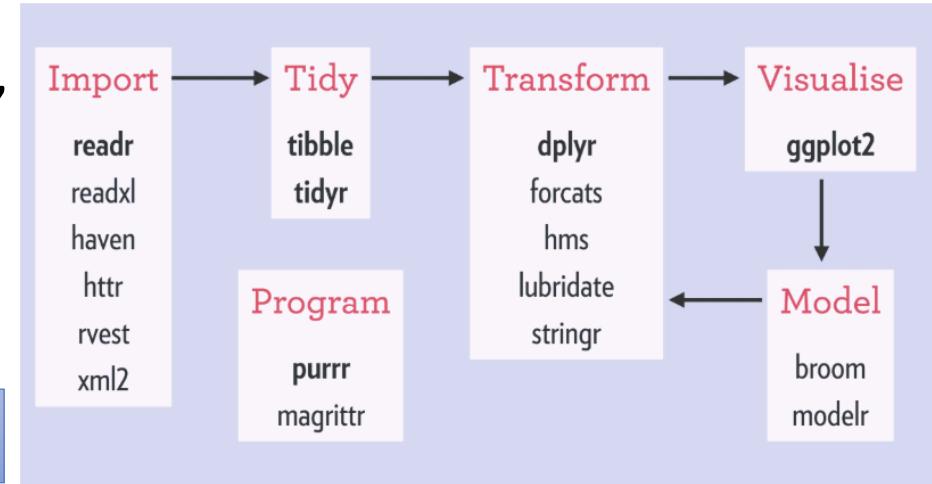
# Data Exploration Workflow



# Packages

- In-built R Functions: R's **R base**.  
Eg. What we've been using till now: str() and head() and so on
- External functions are enabled through **Packages**
- **Packages in R**: Additional functions that let you do more stuff!
- **Tidyverse (2016)** is an “Umbrella Package”  
Lets you perform the entire  
Data exploration Workflow!

**TIDYVERSE**



# Tidyverse

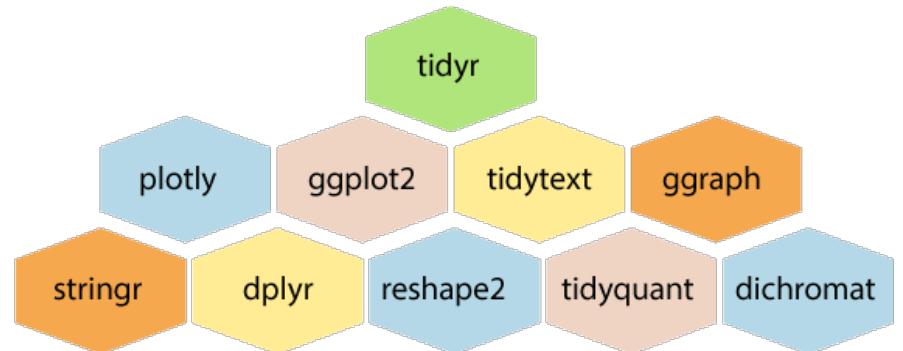
- Easier to read
- 10x Faster
- Strings Not Converted to Factors

```
# R Base
surveys[surveys$species=="albigula" &
          surveys$year==1977, ]

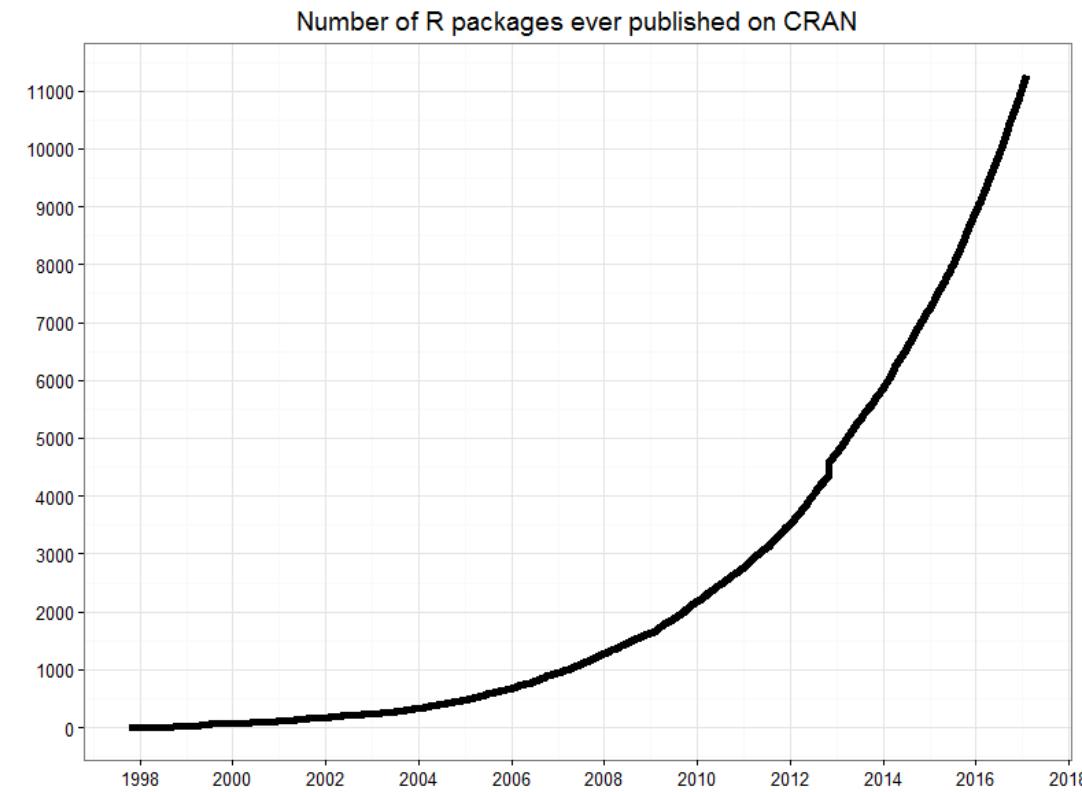
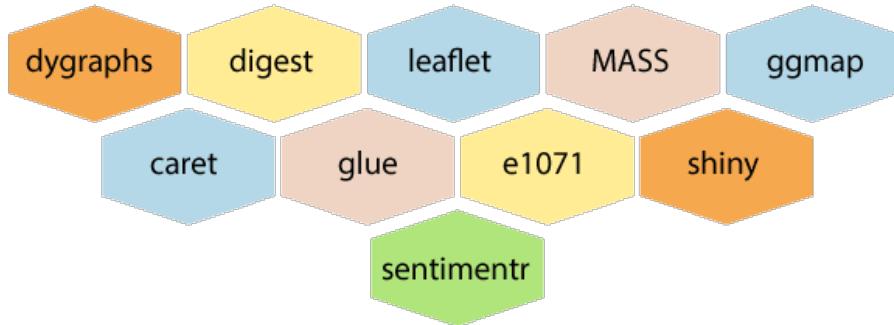
# tidyverse
filter(surveys, species=="albigula" & year==1977)
```



# R Packages

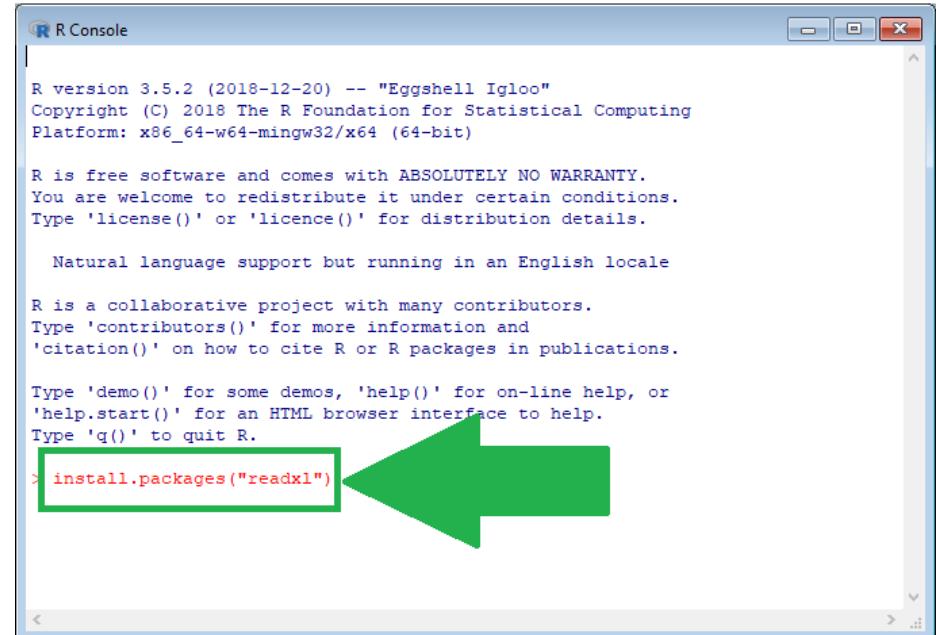


list of Packages



# Installing Packages

- Directly from R console
- No need to Repeat or Re-Install

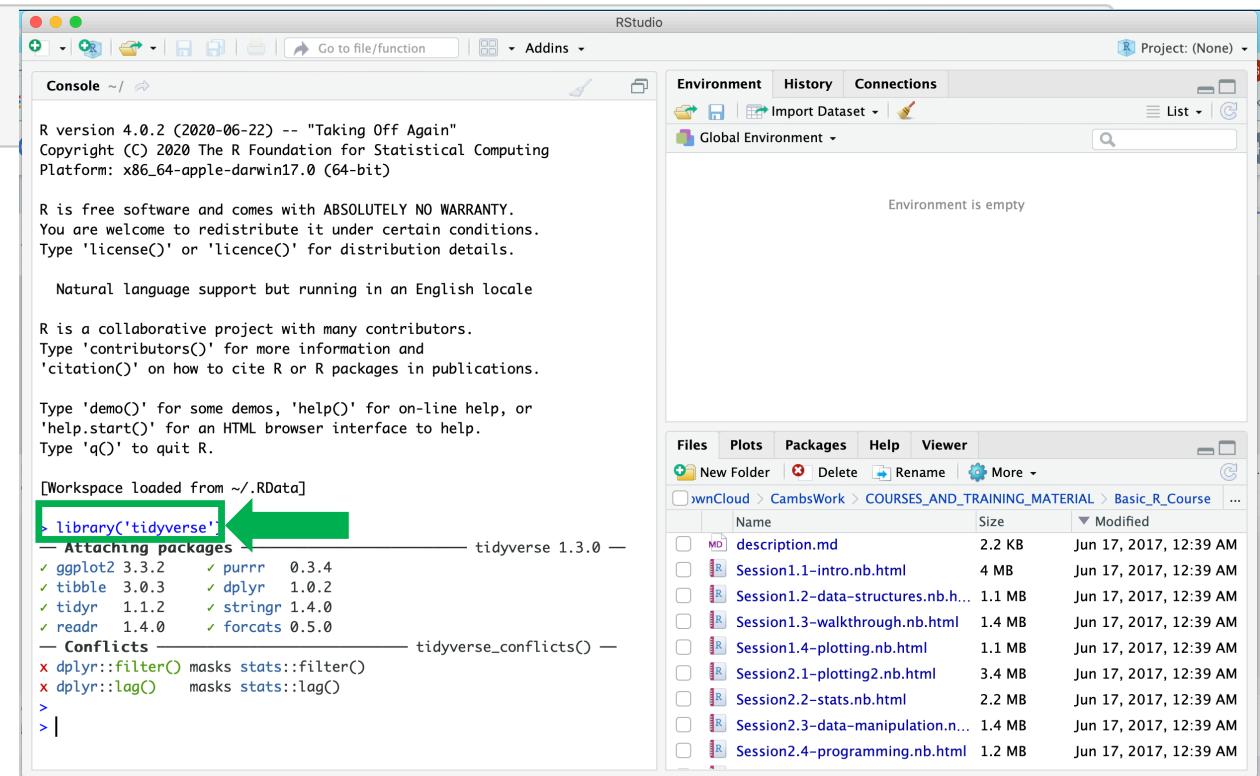


```
#install the tidyverse package  
install.packages("tidyverse")
```

# Loading Packages

```
## load the tidyverse package
library(tidyverse)
```

Try it !



# Importing / Reading Data from Files

```
surveys <- read_csv("data/portal_data_joined.csv")
```

- The **read\_csv** function is from tidyverse package **readr**
- Instead of **read.csv** (R base)

# Tibble

- Modern take on data frames!
- Retain useful features
- Drop Frustrating features  
(i.e. converting character vectors to factors).

Console ~/ownCloud/CambsWork/COURSES\_AND\_TRAINING\_MATERIAL/Basic\_R\_Course/TESTR/

```

>
> surveys <- read_csv("portal_data_joined.csv")

— 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(),
  weight = col_double(),
  genus = col_character(),
  species = col_character(),
  taxa = col_character(),
  plot_type = col_character()
)

> View(surveys)

```



# Tibble

```
## Display first 6 rows
head(surveys)
```

```
## To print the first 15 rows
print(surveys, n=15)
```

- Provides Information on
  - data structure
  - data types of each column
  - Size of the dataset
- Str() function redundant!

```
>
> head(surveys)
# A tibble: 6 x 13
  record_id month   day year plot_id species_id sex  hindfoot_length weight genus species taxa
  <dbl> <dbl> <dbl> <dbl> <dbl> <chr>    <chr>        <dbl> <dbl> <chr> <chr> <chr>
1       1     7    16 1977      2   NL       M            32   NA Neot... albigu... Rode...
2      72     8    19 1977      2   NL       M            31   NA Neot... albigu... Rode...
3     224     9    13 1977      2   NL      NA           NA   NA Neot... albigu... Rode...
4     266    10    16 1977      2   NL      NA           NA   NA Neot... albigu... Rode...
5     349    11    12 1977      2   NL      NA           NA   NA Neot... albigu... Rode...
6     363    11    12 1977      2   NL      NA           NA   NA Neot... albigu... Rode...
# ... with 1 more variable: plot_type <chr>
>
>
> print(surveys)
# A tibble: 34,786 x 13
  record_id month   day year plot_id species_id sex  hindfoot_length weight genus species taxa
  <dbl> <dbl> <dbl> <dbl> <dbl> <chr>    <chr>        <dbl> <dbl> <chr> <chr> <chr>
1       1     7    16 1977      2   NL       M            32   NA Neot... albigu... Rode...
2      72     8    19 1977      2   NL       M            31   NA Neot... albigu... Rode...
3     224     9    13 1977      2   NL      NA           NA   NA Neot... albigu... Rode...
4     266    10    16 1977      2   NL      NA           NA   NA Neot... albigu... Rode...
5     349    11    12 1977      2   NL      NA           NA   NA Neot... albigu... Rode...
6     363    11    12 1977      2   NL      NA           NA   NA Neot... albigu... Rode...
7     435    12    10 1977      2   NL      NA           NA   NA Neot... albigu... Rode...
8     506     1     8 1978      2   NL      NA           NA   NA Neot... albigu... Rode...
9     588     2    18 1978      2   NL       M           NA   NA 218 Neot... albigu... Rode...
10    661     3    11 1978      2   NL      NA           NA   NA Neot... albigu... Rode...
# ... with 34,776 more rows, and 1 more variable: plot_type <chr>
>
```

# DATA VISUALIZATION IN R

# Visualizing data in R (ggplot)

- THREE main elements:

```
ggplot (data = <DATA>, mapping = aes(<MAPPINGS>)) + <GEOM_FUNCTION>()
```

- The function takes TWO arguments:

**data**: Data frame with variables to plot (columns and rows)

**mapping**: Aesthetics (How variables are mapped to visual properties of the plot)

- Using **ggplot** function on its own will not plot anything!

- Add a **geom\_function** as a layer (By using +)

**geom\_function** specifies the type of plot would you like to plot

# Visualizing data in R

Greatest advantage of ggplot:  
➤ Easy to change plot types!  
(New **geom\_function** all else same)

[Order of Steps 2 & 3 does not matter]

## The Basics of ggplot2

There are three main components that get *added* together to build up a plot.

1

**ggplot()**

+

2

**geom\_xxx()**

+

3

**theme()**

Required: "Hey, R, I'm building a plot, and I'm using the data in the `data=` argument\* to do it!"

Optional: "Here's the 'default' mappings (`mapping=aes (...)`) of specific data in that dataset and how I want to plot them – as the 'x' values or 'y' values or something else."

Required: "Give me a 'layer' of that chart as a column chart or a scatterplot, or a line chart, (or something else)."\*

Optional: "If I didn't give it to you above, or if I need to override some part of it, here's the specific mappings (`mapping=aes (...)`) to use for this layer, the specific data (`data=`), and some guidance on what colors to use for what."

Optional: "Here's how I actually want to *format* the visualization: which tick marks and gridlines to show and what color to make them, which axis labels to show and various font properties, etc."

\* Not always the case, and not strictly required, but ggplot2 works best with data in a long format (or, at least, a tidy format). `gather()` from the `tidyverse` package is your friend here.

# Visualizing data in R

*Main aspects:*

**Data source** - A tidy data source in long format

**Geoms** - geometric objects, the type of plot to produce. Line charts, bar charts, tiled plots etc.

**Aesthetics** - this specifies which variables in your data will vary and be plotted.  
(Misnomer – you actually control the look and feel of plots using *themes*!)

*Other aspects (Not discussed today)*

**Themes** - Like styles/CSS in HTML; Tweak / refine font, colors, spacing etc

**Coordinate systems** - Not just the usual x-y, but polar and more exotic systems

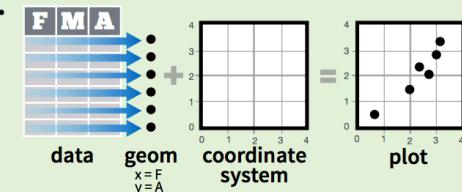
**Scales** - How your variables map onto the coordinate system (e.g. a log scale)

**Statistics** - Applied to data before plotting [eg. Binning for histograms]

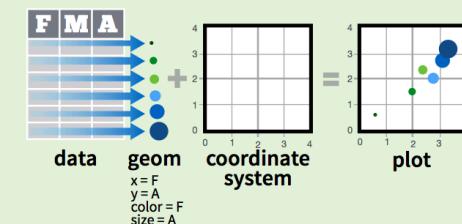
**Thinking about what you want to produce via the components above  
will get you to your desired plot quicker**

## Basics

**ggplot2** is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data** set, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.



To display data values, map variables in the data set to aesthetic properties of the geom like **size**, **color**, and **x** and **y** locations.



# Visualizing data in R

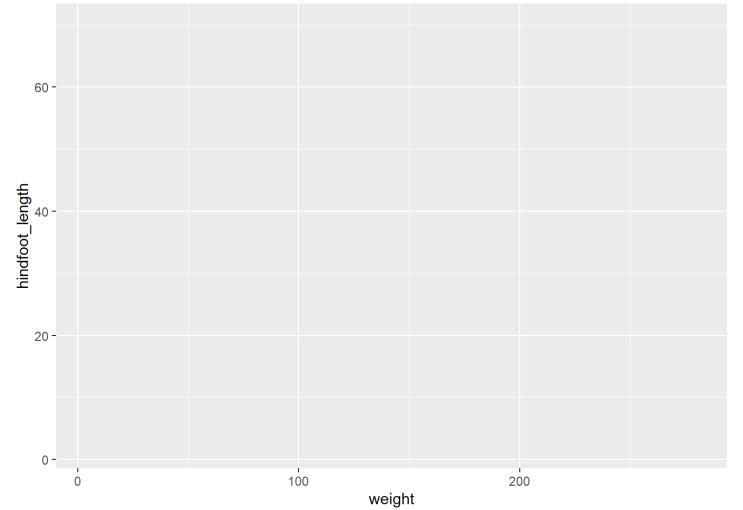
Think of GEOMS as Layers!

Geoms - Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.		
One Variable	Two Variables	Continuous Bivariate Distribution
<b>Continuous</b> <pre>a + geom_area(stat = "bin") x, y, alpha, color, fill, linetype, size b + geom_area(aes(y = ..density..), stat = "bin")</pre>  <pre>a + geom_density(kernel = "gaussian") x, y, alpha, color, fill, linetype, size, weight b + geom_density(aes(y = ..count..))</pre>  <pre>a + geom_dotplot() x, y, alpha, color, fill</pre>  <pre>a + geom_freqpoly() x, y, alpha, color, linetype, size b + geom_freqpoly(aes(y = ..density..))</pre>  <pre>a + geom_histogram(binwidth = 5) x, y, alpha, color, fill, linetype, size, weight b + geom_histogram(aes(y = ..density..))</pre>  <b>Discrete</b> <pre>b &lt;- ggplot(mpg, aes(f1)) b + geom_bar() x, alpha, color, fill, linetype, size, weight</pre> 	<b>Continuous X, Continuous Y</b> <pre>f &lt;- ggplot(mpg, aes(cty, hwy)) f + geom_blank()</pre>  <pre>f + geom_jitter() x, y, alpha, color, fill, shape, size</pre>  <pre>f + geom_point() x, y, alpha, color, fill, shape, size</pre>  <pre>f + geom_quantile() x, y, alpha, color, linetype, size, weight</pre>  <pre>f + geom_rug(sides = "bl") alpha, color, linetype, size</pre>  <pre>f + geom_smooth(model = lm) x, y, alpha, color, fill, linetype, size, weight</pre>  <b>C</b> <b>A</b> <b>B</b> <b>geom_text(aes(label = cty)) x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust</b>	<pre>i &lt;- ggplot(movies, aes(year, rating)) i + geom_bin2d(binwidth = c(5, 0.5)) xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size, weight</pre>  <pre>i + geom_density2d() x, y, alpha, colour, linetype, size</pre>  <pre>i + geom_hex() x, y, alpha, colour, fill, size</pre>  <b>Continuous Function</b> <pre>j &lt;- ggplot(economics, aes(date, unemployed)) j + geom_area() x, y, alpha, color, fill, linetype, size</pre>  <pre>j + geom_line() x, y, alpha, color, linetype, size</pre>  <pre>j + geom_step(direction = "hv") x, y, alpha, color, linetype, size</pre> 
<b>Graphical Primitives</b> <pre>c &lt;- ggplot(map, aes(long, lat)) c + geom_polygon(aes(group = group)) x, y, alpha, color, fill, linetype, size</pre>  <pre>d &lt;- ggplot(economics, aes(date, unemployed)) d + geom_path(lineend = "butt", linejoin = "round", linemitre = 1) x, y, alpha, color, linetype, size</pre>  <pre>d + geom_ribbon(aes(ymin = unemployed - 900, ymax = unemployed + 900)) x, ymax, ymin, alpha, color, fill, linetype, size</pre>  <pre>e &lt;- ggplot(seals, aes(x = long, y = lat)) e + geom_segment(aes( xend = long + delta_long, yend = lat + delta_lat)) x, xend, y, yend, alpha, color, linetype, size</pre>  <pre>e + geom_rect(aes(xmin = long, ymin = lat, xmax = long + delta_long, ymax = lat + delta_lat)) xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size</pre> 	<b>Discrete X, Continuous Y</b> <pre>g &lt;- ggplot(mpg, aes(class, hwy)) g + geom_bar(stat = "identity") x, y, alpha, color, fill, linetype, size, weight</pre>  <pre>g + geom_boxplot() lower, middle, upper, x, ymax, ymin, alpha, color, fill, linetype, shape, size, weight</pre>  <pre>g + geom_dotplot(binaxis = "y", stackdir = "center") x, y, alpha, color, fill</pre>  <pre>g + geom_violin(scale = "area") x, y, alpha, color, fill, linetype, size, weight</pre> 	<b>Discrete X, Discrete Y</b> <pre>h &lt;- ggplot(diamonds, aes(cut, color)) h + geom_jitter() x, y, alpha, color, fill, shape, size</pre> 
		<b>Maps</b> <pre>state &lt;- data.frame(murder = USArrests\$Murder, state = tolower(rrownames(USArrests))) map &lt;- map_data("state") l &lt;- ggplot(data, aes(fill = murder)) l + geom_map(aes(map_id = state), map = map) + expand_limits(x = map\$long, y = map\$lat) map_id, alpha, color, fill, linetype, size</pre> 
		<b>Three Variables</b> <pre>seals\$z &lt;- with(seals, sqrt(delta_long^2 + delta_lat^2)) m &lt;- ggplot(seals, aes(long, lat)) m + geom_raster(aes(fill = z), hijst = 0.5, vjust = 0.5, interpolate = FALSE)</pre>  <pre>m + geom_contour(aes(z = z)) x, y, z, alpha, colour, linetype, size, weight</pre>  <pre>m + geom_tile(aes(fill = z)) x, y, alpha, color, fill, linetype, size</pre> 

# Visualizing data in R

- Let's Try it on the **surveys** data

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length))
```



# Visualizing data in R

- Let's Try it on the **surveys** data

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length))
```

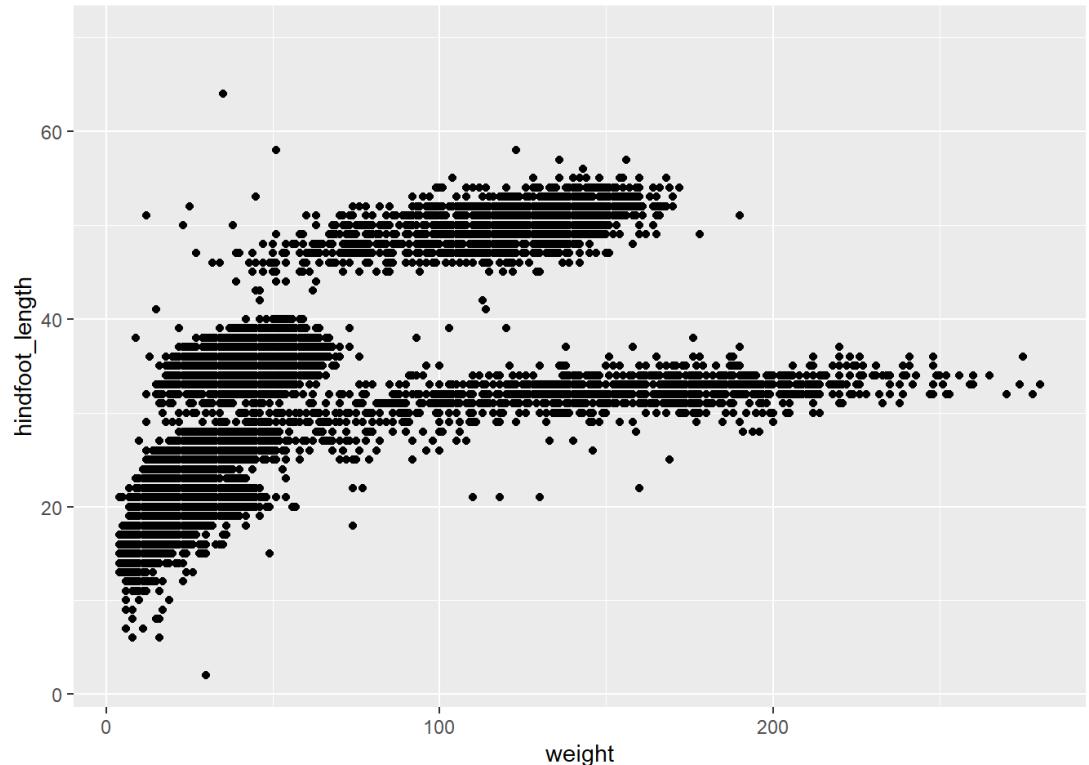
- Add Layers

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point()
```

# Visualizing data in R

➤ Add Layers

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point()
```



# Visualizing data in R

- Let's Try it on the **surveys** data

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length))
```

- Add Layers

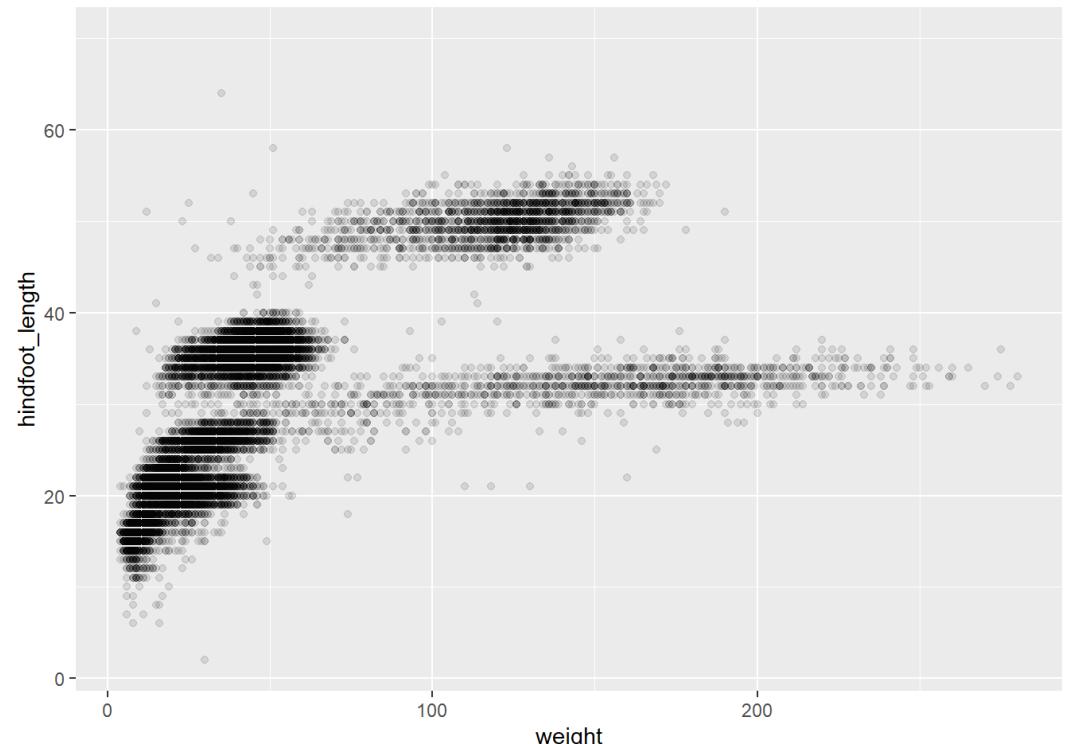
```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point()
```

- Customize: Add Transparency

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point (alpha = 0.1)
```

# Visualizing data in R

- Customize: Add Transparency



```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point (alpha = 0.1)
```

# Visualizing data in R

- Let's Try it on the **surveys** data

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length))
```

- Add Layers

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point()
```

- Customize: Add Transparency

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point (alpha = 0.1)
```

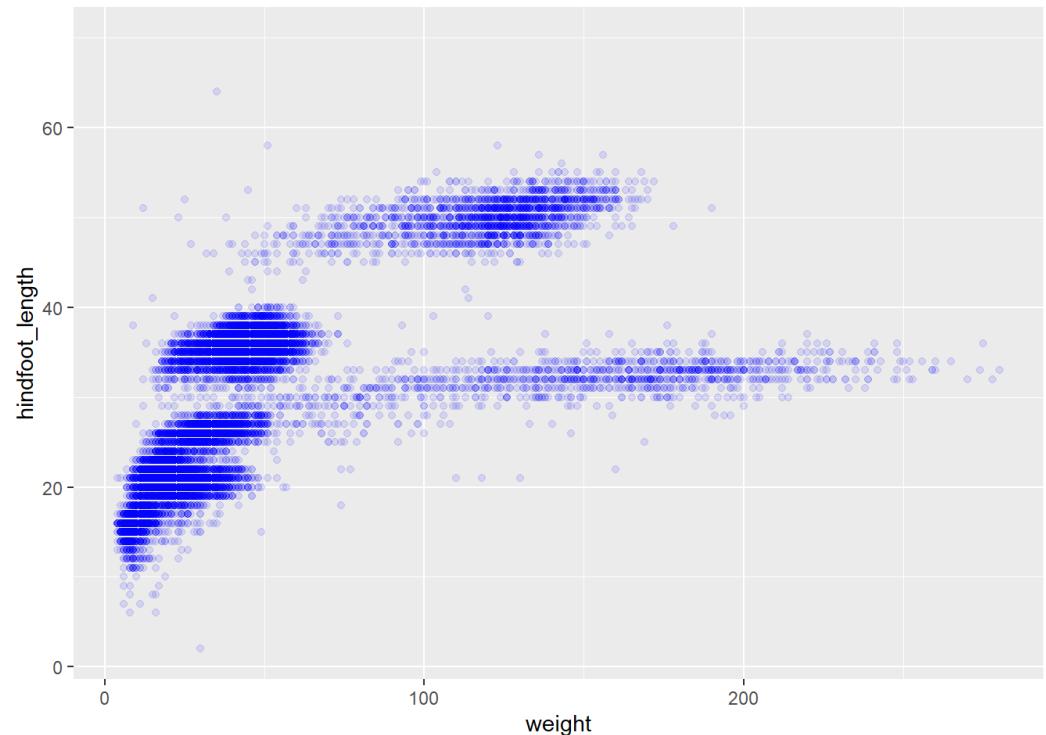
- Add Color

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point (alpha = 0.1, color = 'blue')
```

# Visualizing data in R

➤ Add color

```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point (alpha = 0.1, color = 'blue')
```



# Visualizing data in R

Try other types of plots!

- Save the ggplot into a variable 'surveys\_plot'

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

- Draw a scatter plot

```
surveys_plot + geom_point()
```

- Now draw a geom\_smooth plot

```
surveys_plot + geom_smooth()
```

# Visualizing data in R

Try other types of plots!

- Save the ggplot into a variable ‘surveys\_plot’

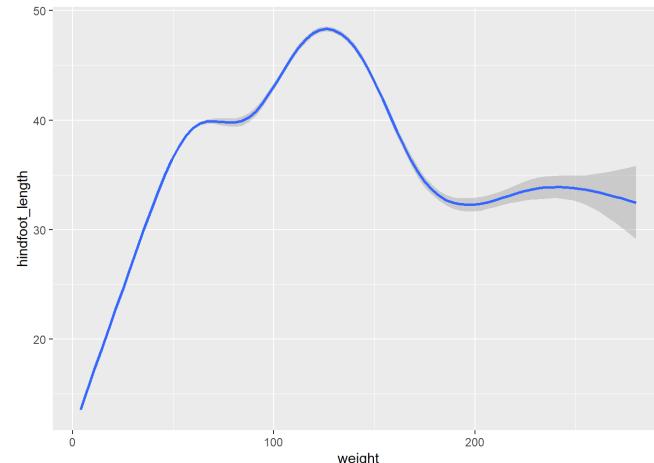
```
surveys_plot <- ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length))
```

- Draw a scatter plot

```
surveys_plot + geom_point()
```

- Now draw a geom\_smooth plot

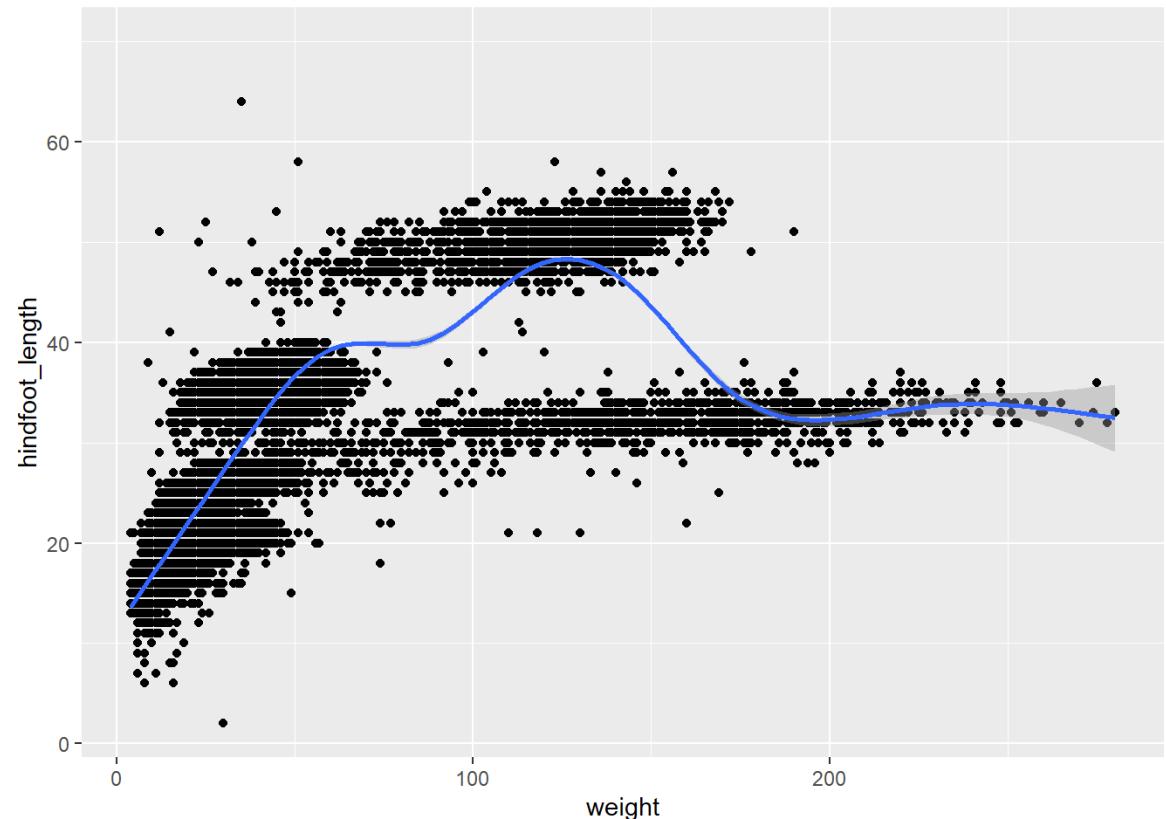
```
surveys_plot +  
geom_smooth()
```



# Visualizing data in R

Multiple Plots as Layers:

```
surveys_plot  
+ geom_point()  
+ geom_smooth()
```



What happens if you change the order of layers added!?

# Visualizing data in R

## Syntax is Important!

+ sign MUST come at the END of line of previous layer

```
# This is the correct syntax for adding layers
surveys_plot +
  geom_point()

# This will not add the new layer and will return an error message
surveys_plot
  + geom_point()
```

Mappings for a given geom can be independent of global mappings defined in the ggplot() function

# Exercise : Data Visualisation

- Scatter plots can be useful exploratory tools for **small** datasets.
- For data sets with **large** numbers of observations, such as the surveys data set, **overplotting** of points can be a limitation of scatter plots.
- We have already seen how we can visualise data better when we have overplotting with the **geom\_smooth** plot.
- Another way for handling overplotting is to display the **density** of the data through **contours**.
- **As this challenge's task, do the following!**
  - Create a script called **plot\_density2d.R** which :
    1. loads the file **data/portal\_data\_joined.csv** into the variable **surveys**.
    2. It then uses this dataset to plot the **weight** on the x-axis and **hindfoot\_length** on the y-axis in a **geom\_density2d** plot

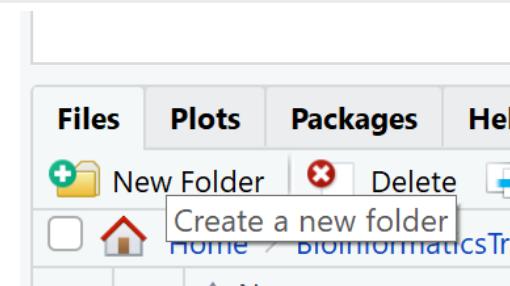
# Saving a plot to a file

```
#save plot that you would like to save into a variable  
out_plot <- surveys_plot + geom_density2d()  
#save plot to file  
ggsave(filename="img_output/plot_weight_hindfoot_density2d.png", plot=out_plot)
```

[Not necessary to save plot first; **ggsave** will auto save the LAST plot plotted by R]

```
#save plot to file  
ggsave(filename="img_output/plot_weight_hindfoot_density2d.png")
```

Make a new folder 'img\_output' beforehand!  
Create from within RStudio



NEXT

## DATA TRANSFORMATION

