

# Introduction to for Biologists

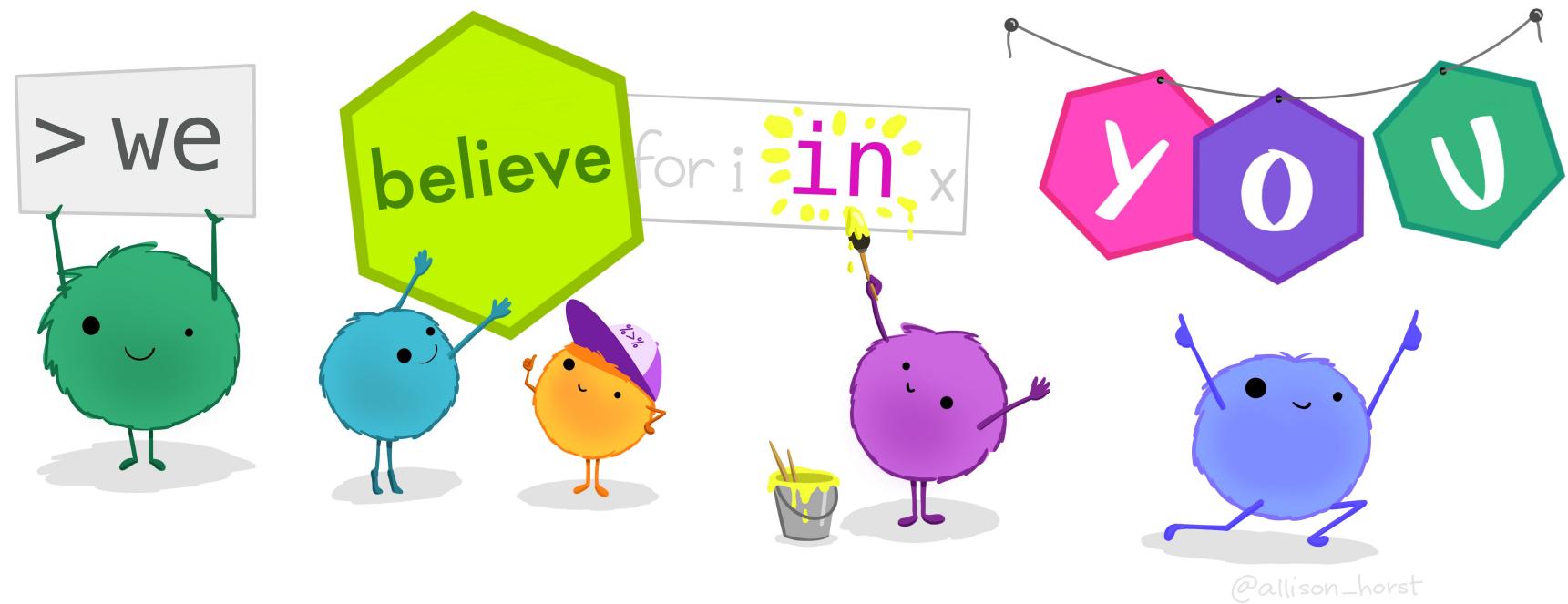


## WORKFLOW & VISUALISATION

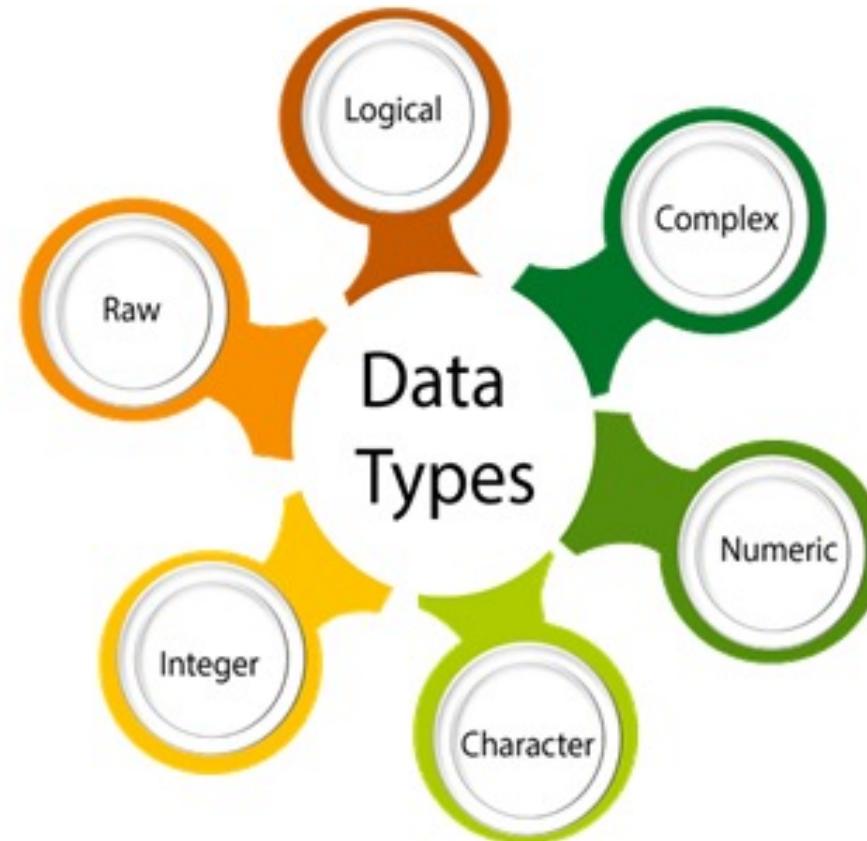
Author: Gita Yadav

Indebted to Alison Horst!

remote R learners,



# Data Type (Recap)



# Recap: Reading in Data from a file

- We've **looked** at the **raw** format of the file (CSV format)
- Let us **Load** the data into R
- Use function `read.csv()` to load the data into an object of class `data.frame`

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

- Look at the data

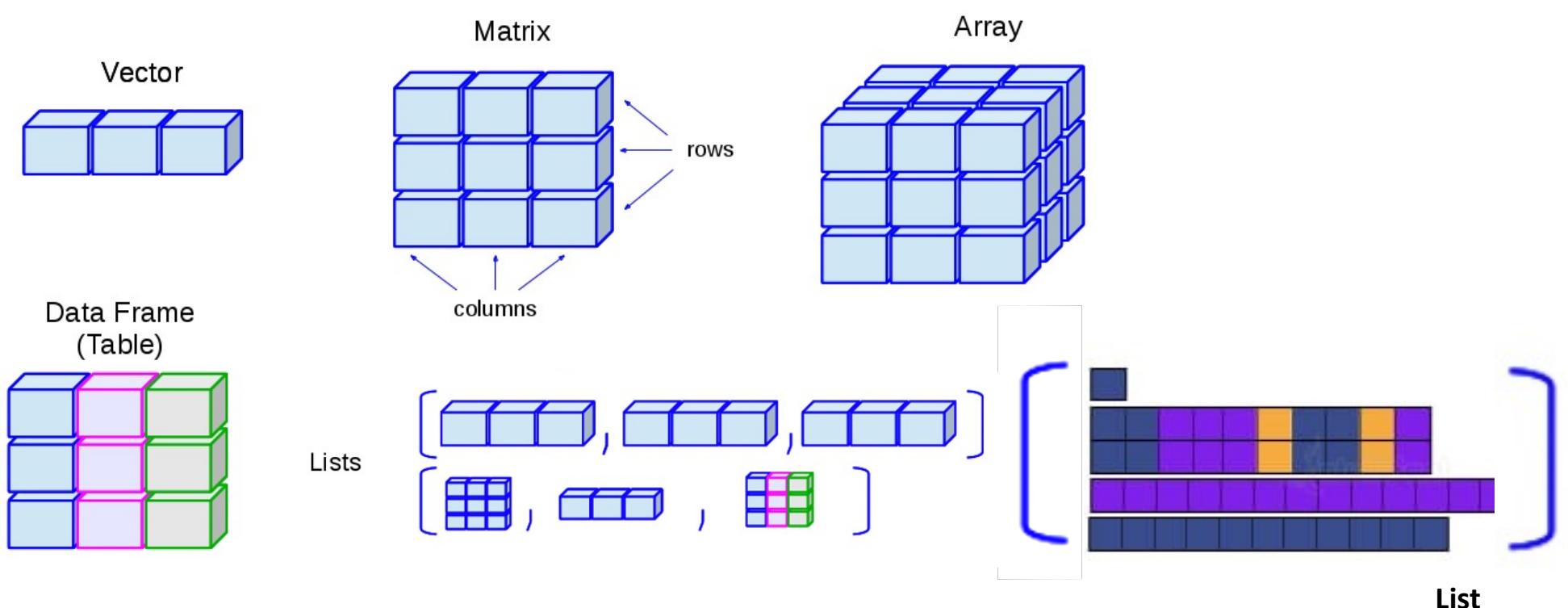
```
head(surveys)
```

```
## Try also  
View(surveys)
```

# Recap Data frame

data frame	1	"S"	TRUE
numeric	7	"A"	FALSE
character	3	"U"	TRUE

# Data Structures (Recap)



# STRUCTURE of a data.frame Object

```
str(surveys)
```

```
#> 'data.frame': 34786 obs. of 13 variables:  
#> $ record_id : int 1 72 224 266 349 363 435 506 588 661 ...  
#> $ month      : int 7 8 9 10 11 11 12 1 2 3 ...  
#> $ day        : int 16 19 13 16 12 12 10 8 18 11 ...  
#> $ year       : int 1977 1977 1977 1977 1977 1977 1977 1978 1978 1978 ...  
#> $ plot_id    : int 2 2 2 2 2 2 2 2 2 2 ...  
#> $ species_id : chr "NL" "NL" "NL" "NL" ...  
#> $ sex         : chr "M" "M" "" "" ...  
#> $ hindfoot_length: int 32 31 NA NA NA NA NA NA NA NA ...  
#> $ weight      : int NA NA NA NA NA NA NA 218 NA ...  
#> $ genus       : chr "Neotoma" "Neotoma" "Neotoma" "Neotoma" ...  
#> $ species     : chr "albigula" "albigula" "albigula" "albigula" ...  
#> $ taxa         : chr "Rodent" "Rodent" "Rodent" "Rodent" ...  
#> $ plot_type   : chr "Control" "Control" "Control" "Control" ...
```



# Inspecting a `data.frame` object

- Size:
  - `dim(surveys)` - returns a vector with the number of rows in the first element, and the number of columns as the second element (the **dimensions** of the object)
  - `nrow(surveys)` - returns the number of rows
  - `ncol(surveys)` - returns the number of columns
- Content:
  - `head(surveys)` - shows the first 6 rows
  - `tail(surveys)` - shows the last 6 rows
- Names:
  - `names(surveys)` - returns the column names (synonym of `colnames()` for `data.frame` objects)
  - `rownames(surveys)` - returns the row names
- Summary:
  - `str(surveys)` - structure of the object and information about the class, length and content of each column
  - `summary(surveys)` - summary statistics for each column



# Indexing and sub-setting data frames

```
> head(surveys)
#> #> #> #> #>
#>   record_id month day year plot_id species_id sex hindfoot_length weight   genus species   taxa plot_type
#> 1       1      7    16 1977       2        NL     M             32    NA Neotoma albigena Rodent Control
#> 2      72      8    19 1977       2        NL     M             31    NA Neotoma albigena Rodent Control
#> 3     224      9   13 1977       2        NL           NA    NA Neotoma albigena Rodent Control
#> 4      ...      ... ...       2        NL           NA    NA Neotoma albigena Rodent Control
#> 5      ...      ... ...       2        NL           NA    NA Neotoma albigena Rodent Control
#> 6     363      ... ...       2        NL           NA    NA Neotoma albigena Rodent Control
#>
```

	record_id	month	day	year	plot_id	species_id	sex	hindfoot_length	weight	genus	species	taxa	plot_type
1	[1,1]	[1,2]	[1,3]										
2	[2,1]	[2,2]	[2,3]										
3	[3,1]	[3,2]	[3,3]										

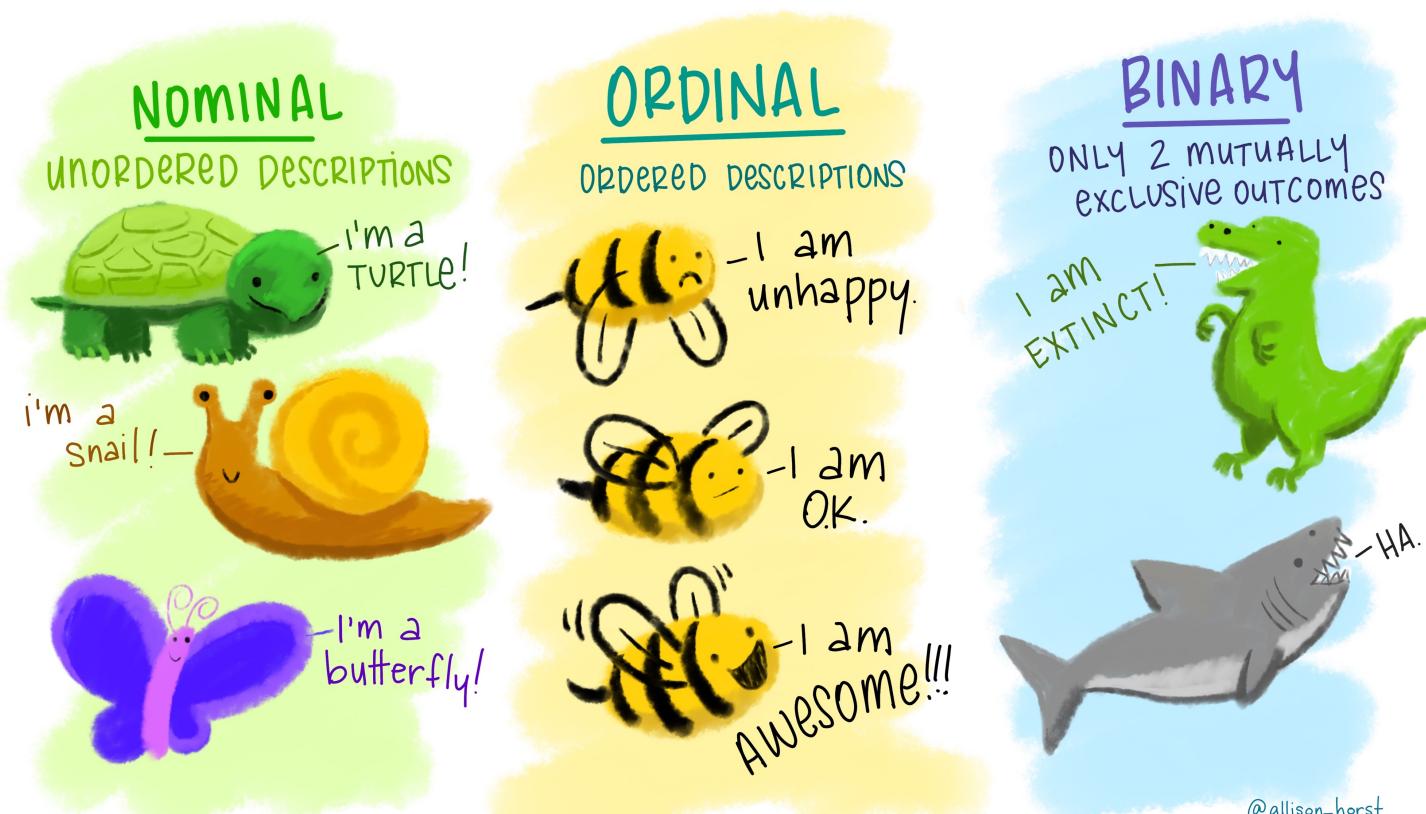
rownames (surveys)

colnames (surveys)

index [R,C]

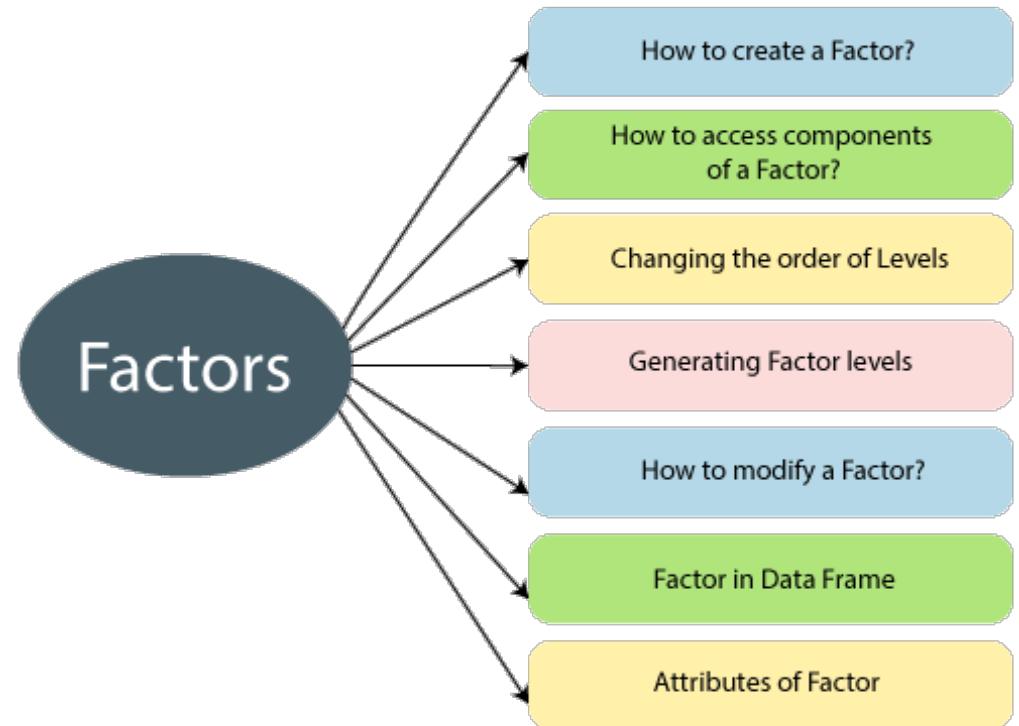
Numeric Indexing [1,2]  
 Name indexing \$  
 Logical indexing == | >=

# The Strange Case of “Categorical Variables”!

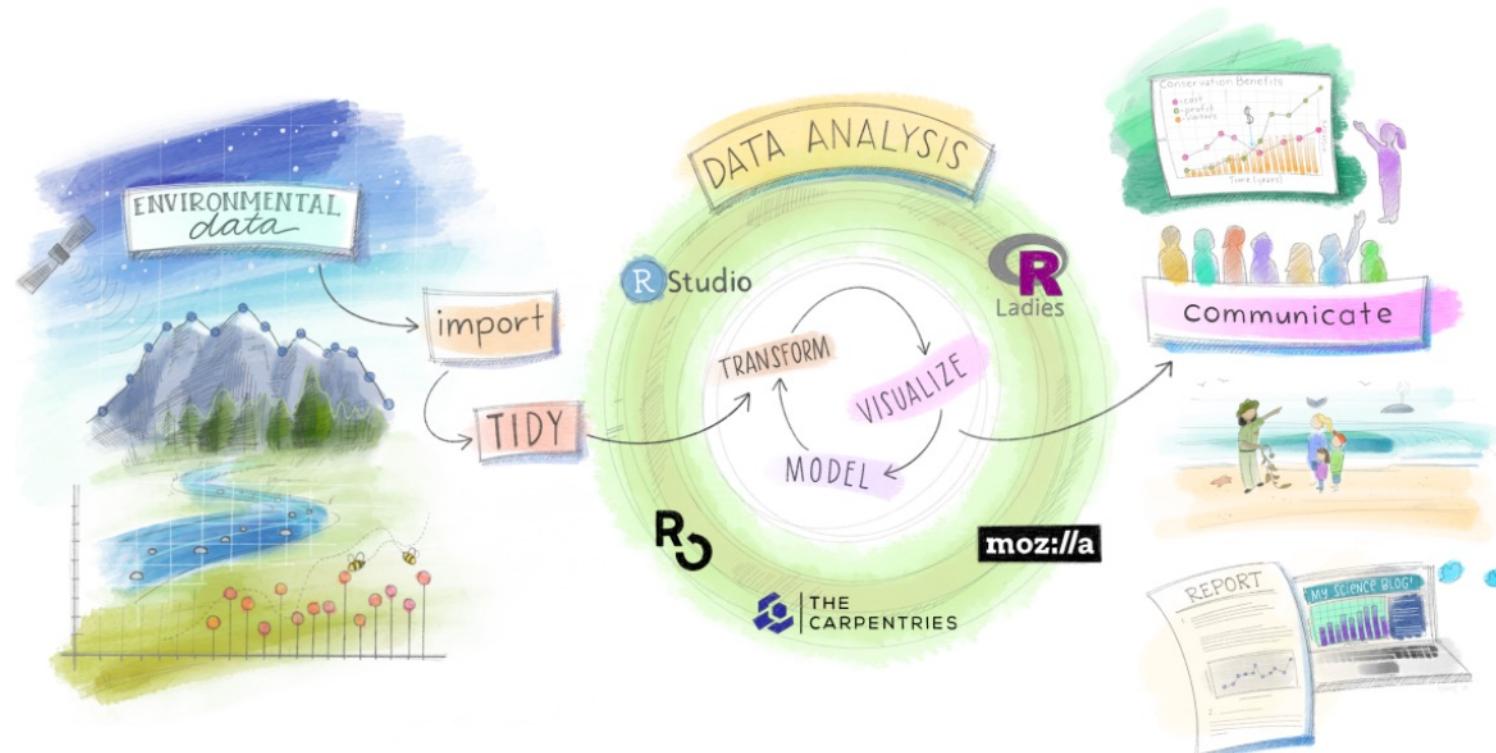


# Factors

- **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!?

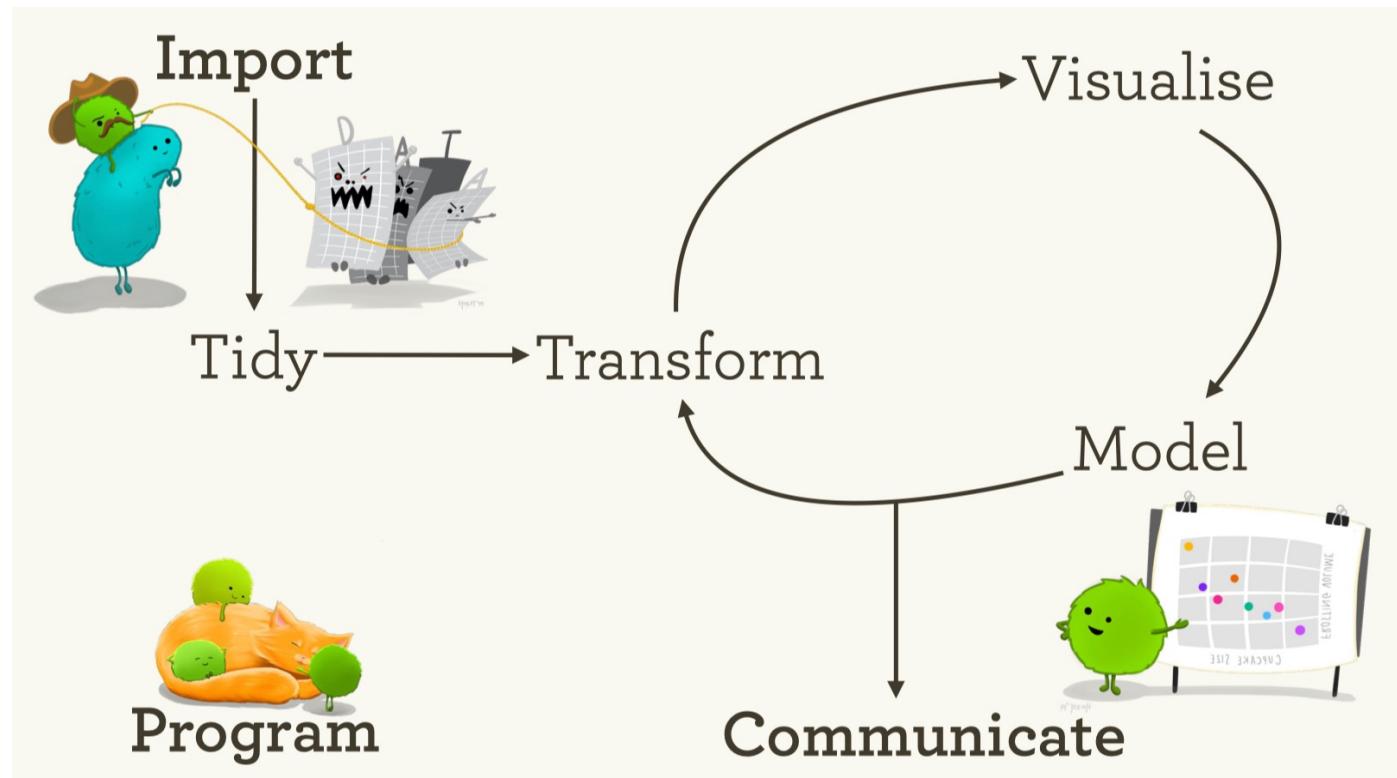


# The Data Exploration Workflow



art: @allison\_horst; updated from Wickham & Grolemund

# Learning Data Science!

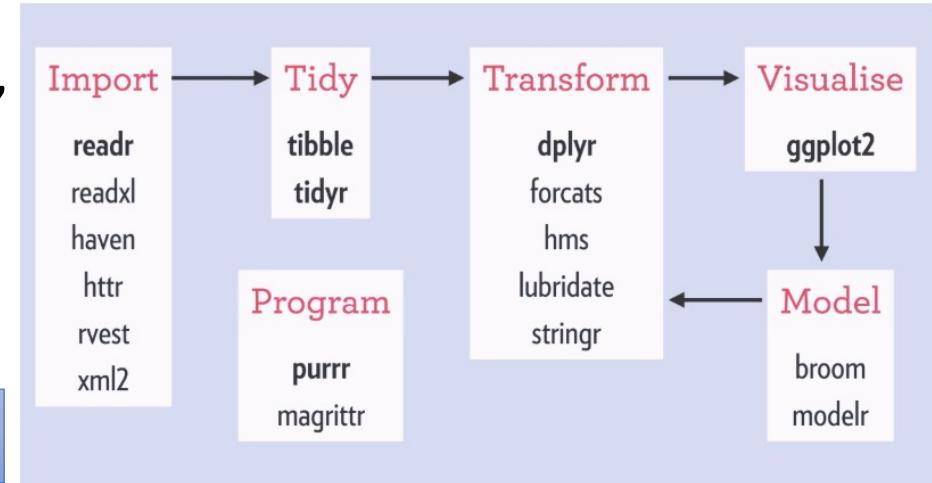


Art: Alison Horst (GitHub)

# 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**



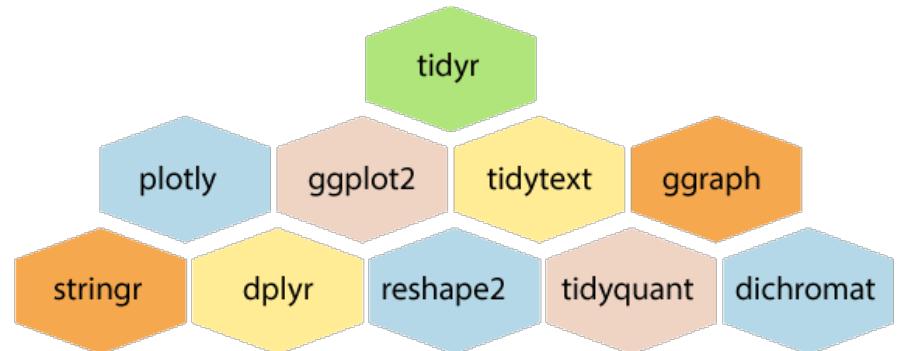
# Tidyverse

- Easier to read
- 10x Faster

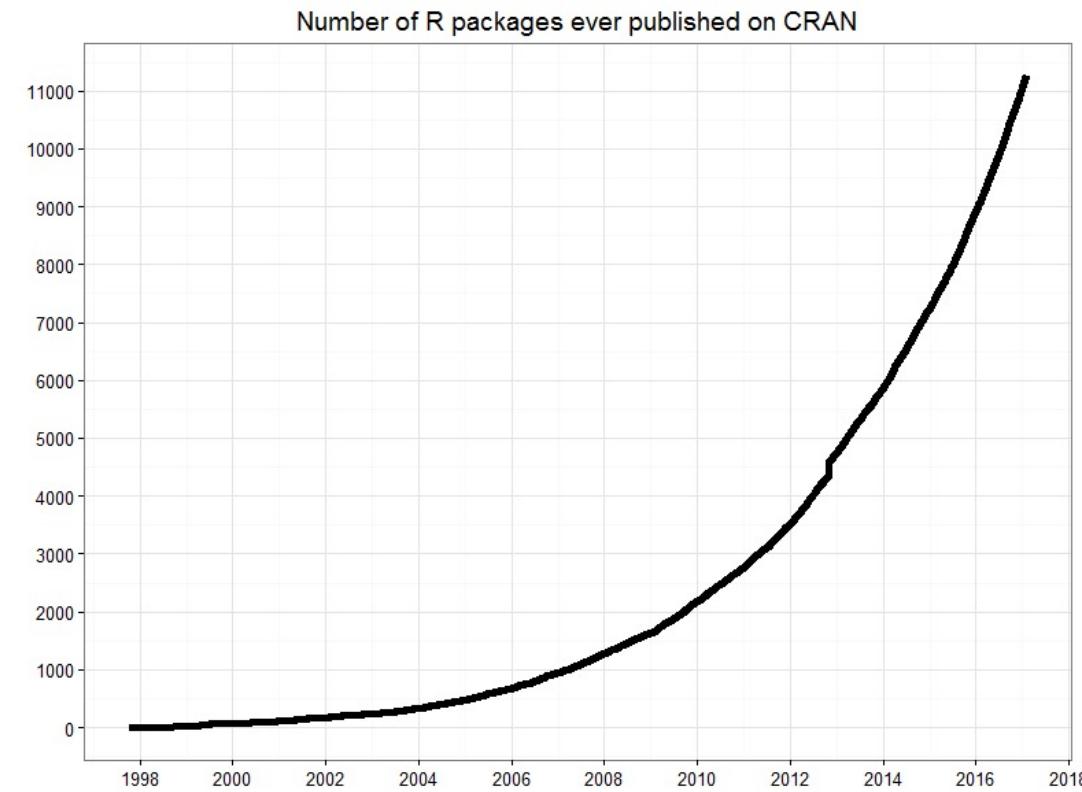
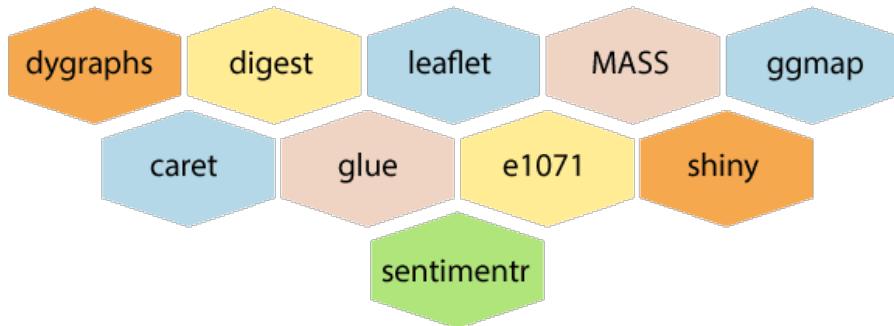
```
# R Base
surveys[surveys$species=="albigula" &
          surveys$year==1977, ]

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

# R Packages

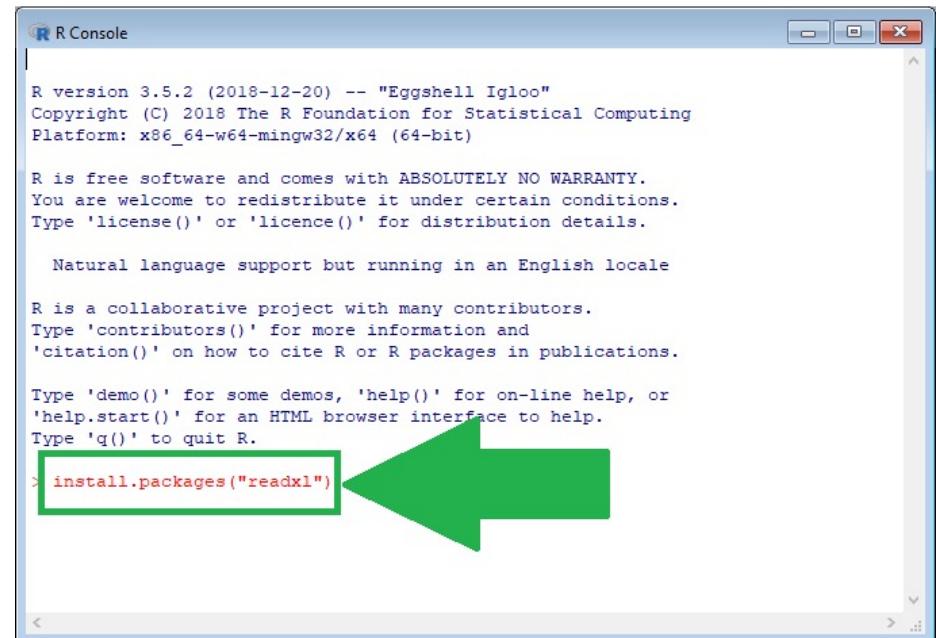


list of Packages



# Installing Packages

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

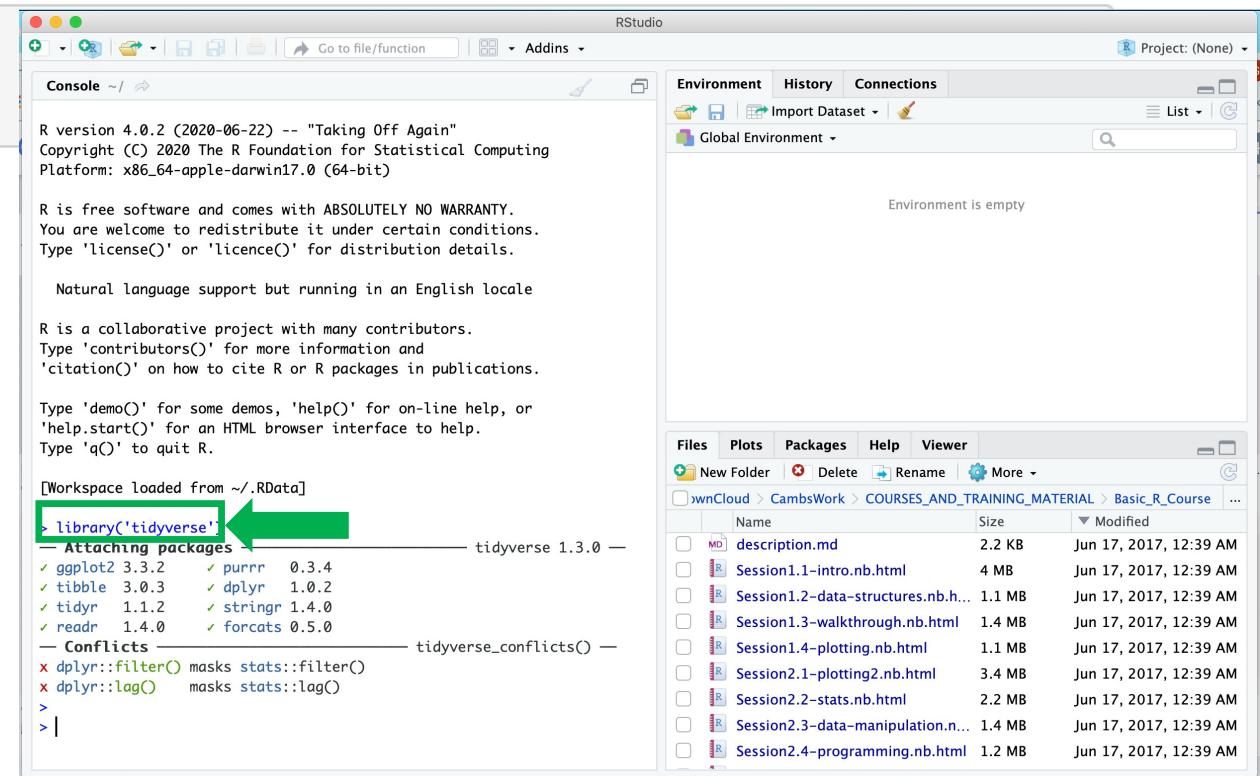


```
#install the tidyverse package  
install.packages("tidyverse")
```

# 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)

# What's a tibble

- Object of class `tbl`
- A reimagining of the `data.frame`
- Looks and acts like a `data.frame`
- ..but its even better
- Retains useful features
- drops frustrating features
- `tidyverse` works with tibbles

```
Console ~/ownCloud/CambsWork/COURSES_AND_TRAINING_MATERIAL/Basic_R_Course/TEST.R/ ┌─────────┐ ┌─────────┐
>
> 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

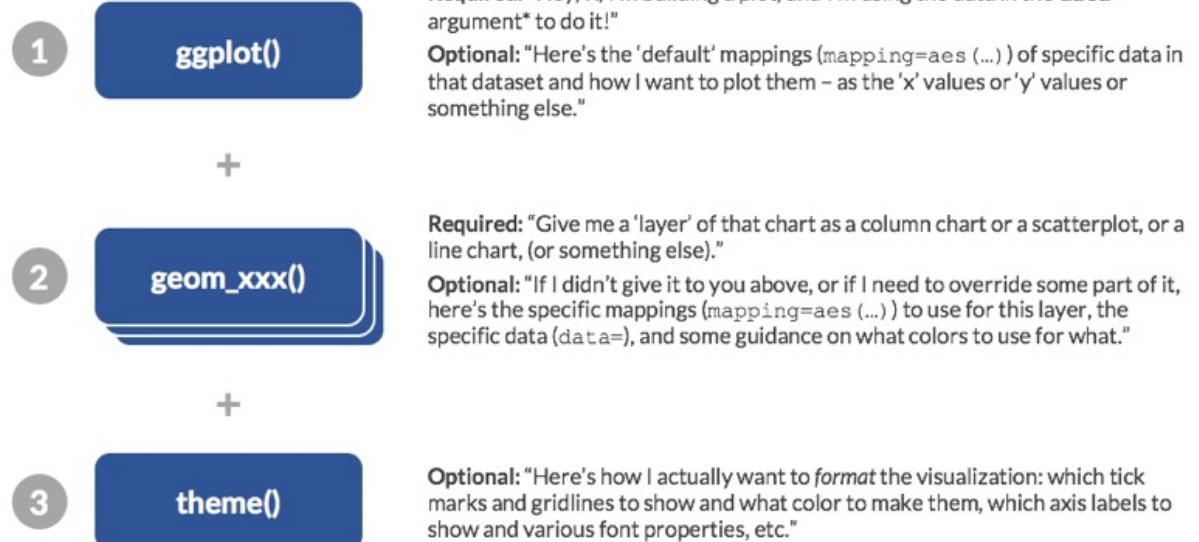
# Visualizing data in R

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.



\* 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

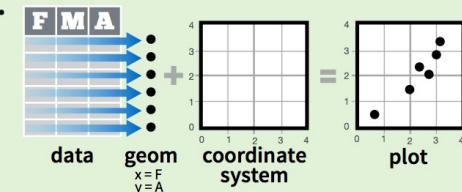
**Scales** - How your variables map onto the coordinate system (e.g. a log scale)

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

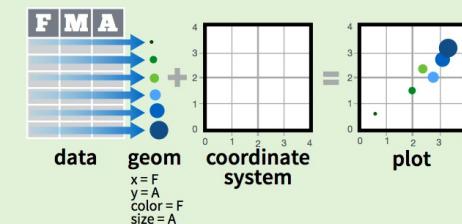
**Thinking about what you want to produce via the components above  
will get you to your desired plot quicker**

## 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

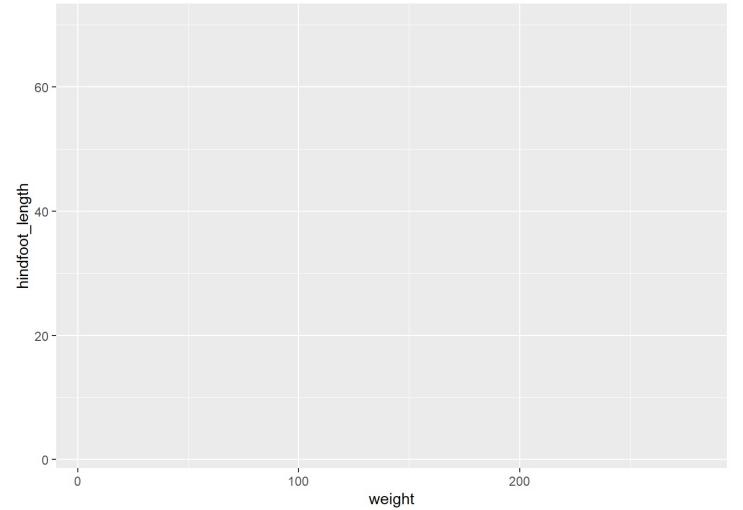
Think of GEOMS as Layers!

Geoms - Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.		
<b>One Variable</b> <ul style="list-style-type: none"> <li><b>Continuous</b> <ul style="list-style-type: none"> <li>a + <code>geom_area(stat = "bin")</code> x, y, alpha, color, fill, linetype, size</li> <li>a + <code>geom_density(kernel = "gaussian")</code> x, y, alpha, color, fill, linetype, size, weight</li> <li>a + <code>geom_dotplot()</code> x, y, alpha, color, fill</li> <li>a + <code>geom_freqpoly()</code> x, y, alpha, color, linetype, size b + <code>geom_freqpoly(aes(y = ..density..))</code></li> <li>a + <code>geom_histogram(binwidth = 5)</code> x, y, alpha, color, fill, linetype, size, weight b + <code>geom_histogram(aes(y = ..density..))</code></li> </ul> </li> </ul>	<b>Two Variables</b> <ul style="list-style-type: none"> <li><b>Continuous X, Continuous Y</b> <ul style="list-style-type: none"> <li>f + <code>geom_blank()</code></li> <li>f + <code>geom_jitter()</code> x, y, alpha, color, fill, shape, size</li> <li>f + <code>geom_point()</code> x, y, alpha, color, fill, shape, size</li> <li>f + <code>geom_quantile()</code> x, y, alpha, color, linetype, size, weight</li> <li>f + <code>geom_rug(sides = "bl")</code> alpha, color, linetype, size</li> <li>f + <code>geom_smooth(model = lm)</code> x, y, alpha, color, fill, linetype, size, weight</li> <li>C f + <code>geom_text(aes(label = cty))</code> x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust</li> </ul> </li> </ul>	<b>Continuous Bivariate Distribution</b> <ul style="list-style-type: none"> <li>i + <code>geom_bin2d(binwidth = c(5, 0.5))</code> xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size, weight</li> <li>i + <code>geom_density2d()</code> x, y, alpha, colour, linetype, size</li> <li>i + <code>geom_hex()</code> x, y, alpha, colour, fill, size</li> </ul>
<b>Graphical Primitives</b> <ul style="list-style-type: none"> <li>c + <code>geom_bar(stat = "identity")</code> x, y, alpha, color, fill, linetype, size</li> <li>c + <code>geom_polygon(aes(group = group))</code> x, y, alpha, color, fill, linetype, size</li> <li>d + <code>geom_path(lineend = "butt", linejoin = "round", linemitre = 1)</code> x, y, alpha, color, linetype, size</li> <li>d + <code>geom_ribbon(aes(ymin = unemploy - 900, ymax = unemploy + 900))</code> x, ymax, ymin, alpha, color, fill, linetype, size</li> <li>e + <code>geom_segment(aes(xend = long + delta_long, yend = lat + delta_lat))</code> x, xend, y, yend, alpha, color, linetype, size</li> <li>e + <code>geom_rect(aes(xmin = long, ymin = lat, xmax = long + delta_long, ymax = lat + delta_lat))</code> xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size</li> </ul>	<b>Discrete X, Continuous Y</b> <ul style="list-style-type: none"> <li>g + <code>geom_boxplot()</code> lower, middle, upper, x, ymax, ymin, alpha, color, fill, linetype, shape, size, weight</li> <li>g + <code>geom_dotplot(binaxis = "y", stackdir = "center")</code> x, y, alpha, color, fill</li> <li>g + <code>geom_violin(scale = "area")</code> x, y, alpha, color, fill, linetype, size, weight</li> </ul>	<b>Continuous Function</b> <ul style="list-style-type: none"> <li>j + <code>geom_area()</code> x, y, alpha, color, fill, linetype, size</li> <li>j + <code>geom_line()</code> x, y, alpha, color, linetype, size</li> <li>j + <code>geom_step(direction = "hv")</code> x, y, alpha, color, linetype, size</li> </ul>
	<b>Discrete X, Discrete Y</b> <ul style="list-style-type: none"> <li>h + <code>geom_jitter()</code> x, y, alpha, color, fill, shape, size</li> </ul>	<b>Visualizing error</b> <ul style="list-style-type: none"> <li>k + <code>geom_crossbar(fatten = 2)</code> x, y, ymax, ymin, alpha, color, fill, linetype, size</li> <li>k + <code>geom_errorbar()</code> x, ymax, ymin, alpha, color, linetype, size, width (also <code>geom_errorbarh()</code>)</li> <li>k + <code>geom_linerange()</code> x, ymin, ymax, alpha, color, linetype, size</li> <li>k + <code>geom_pointrange()</code> x, y, ymin, ymax, alpha, color, fill, linetype, shape, size</li> </ul>
		<b>Maps</b> <ul style="list-style-type: none"> <li>l + <code>geom_map(aes(map_id = state), map = map) + expand_limits(x = map\$long, y = map\$lat)</code> map_id, alpha, color, fill, linetype, size</li> </ul>
		<b>Three Variables</b> <ul style="list-style-type: none"> <li>m + <code>geom_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)</code> x, y, alpha, fill</li> <li>m + <code>geom_tile(aes(fill = z))</code> x, y, alpha, color, fill, linetype, size</li> </ul>

# 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

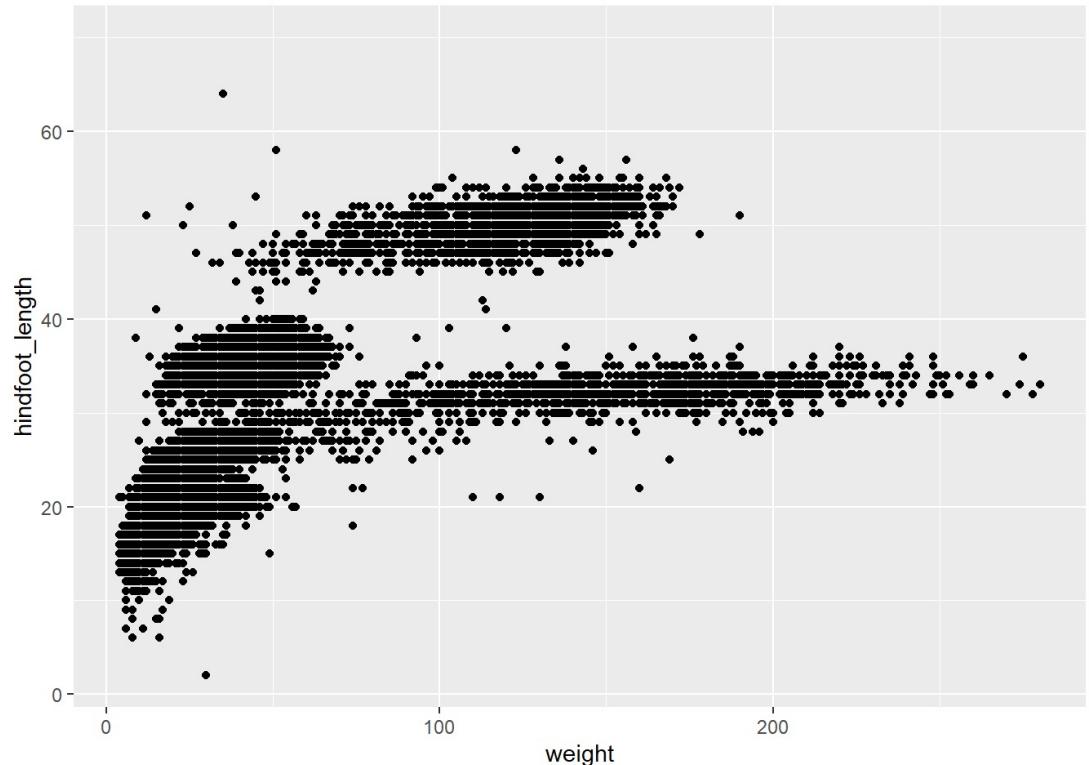
```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point()
```



# Visualizing data in R

➤ Add Layers

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



# Visualizing data in R

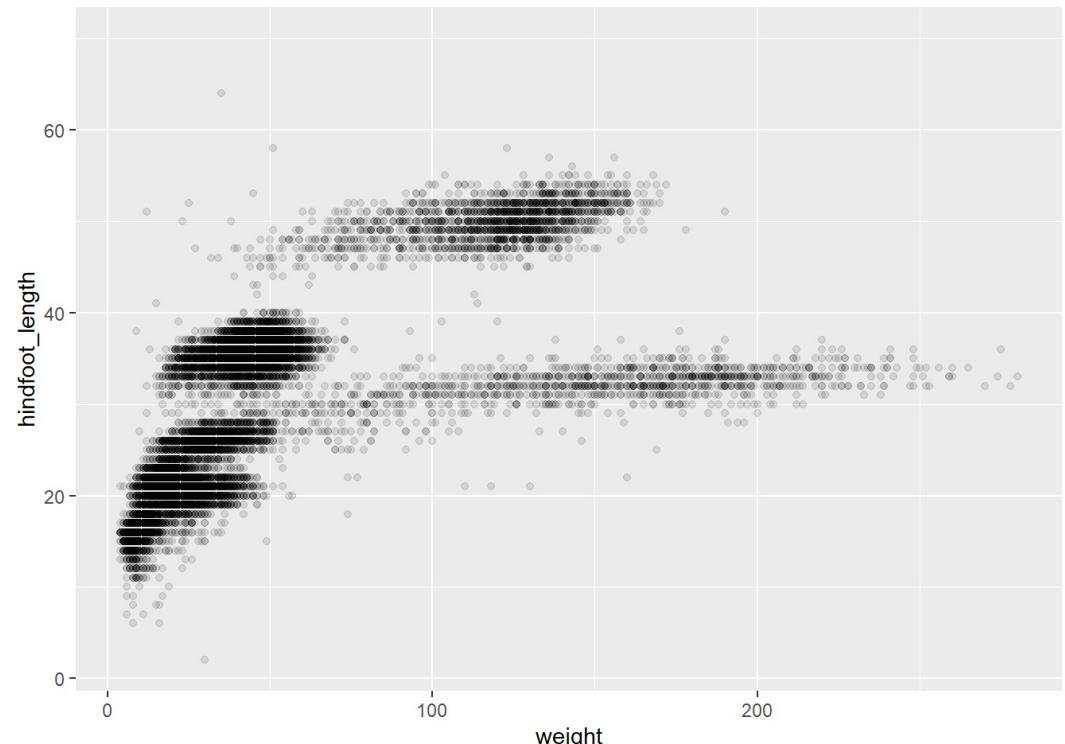
- Customize: Add Transparency

```
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

- Add color

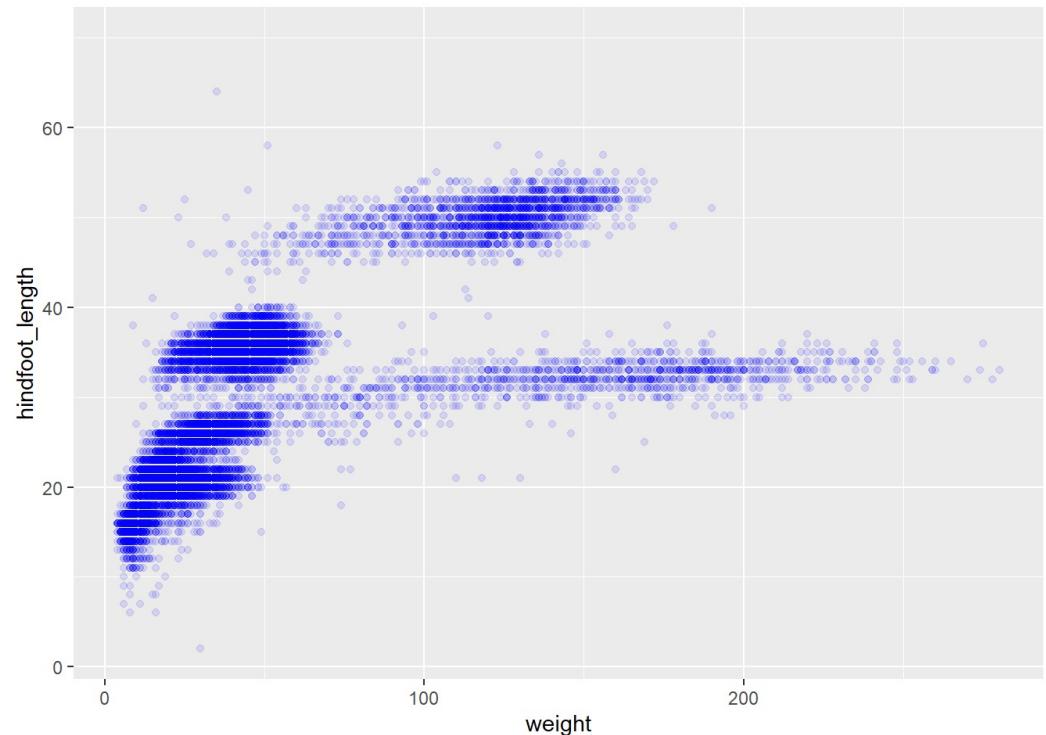
```
ggplot (data = surveys , mapping = aes ( x = weight, y = hindfoot_length)) + geom_point (alpha = 0.1, color = 'blue')
```



# 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

**Let's 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()
```

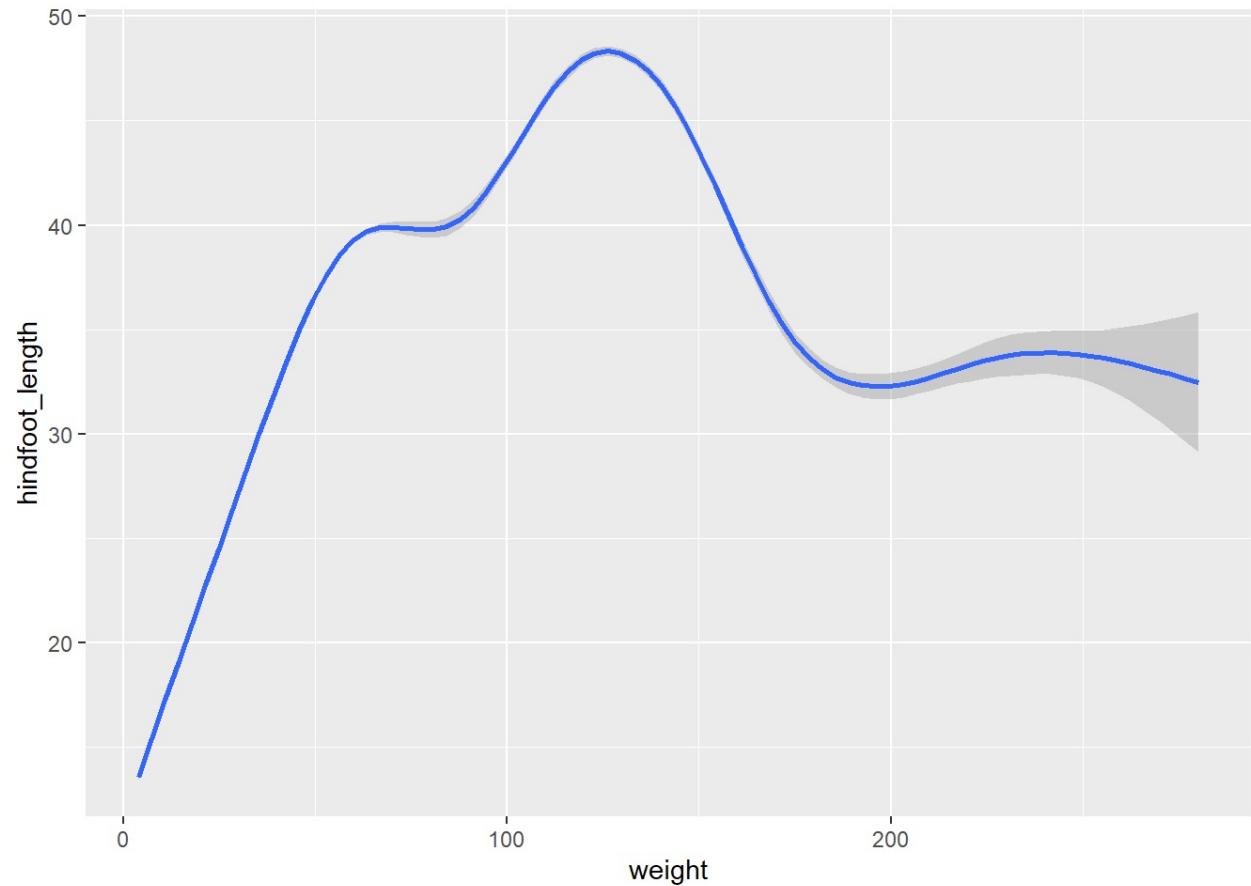
- Now draw a geom\_smooth plot

```
surveys_plot +  
  geom_smooth()
```

# Visualizing data in R

## A geom-smooth plot

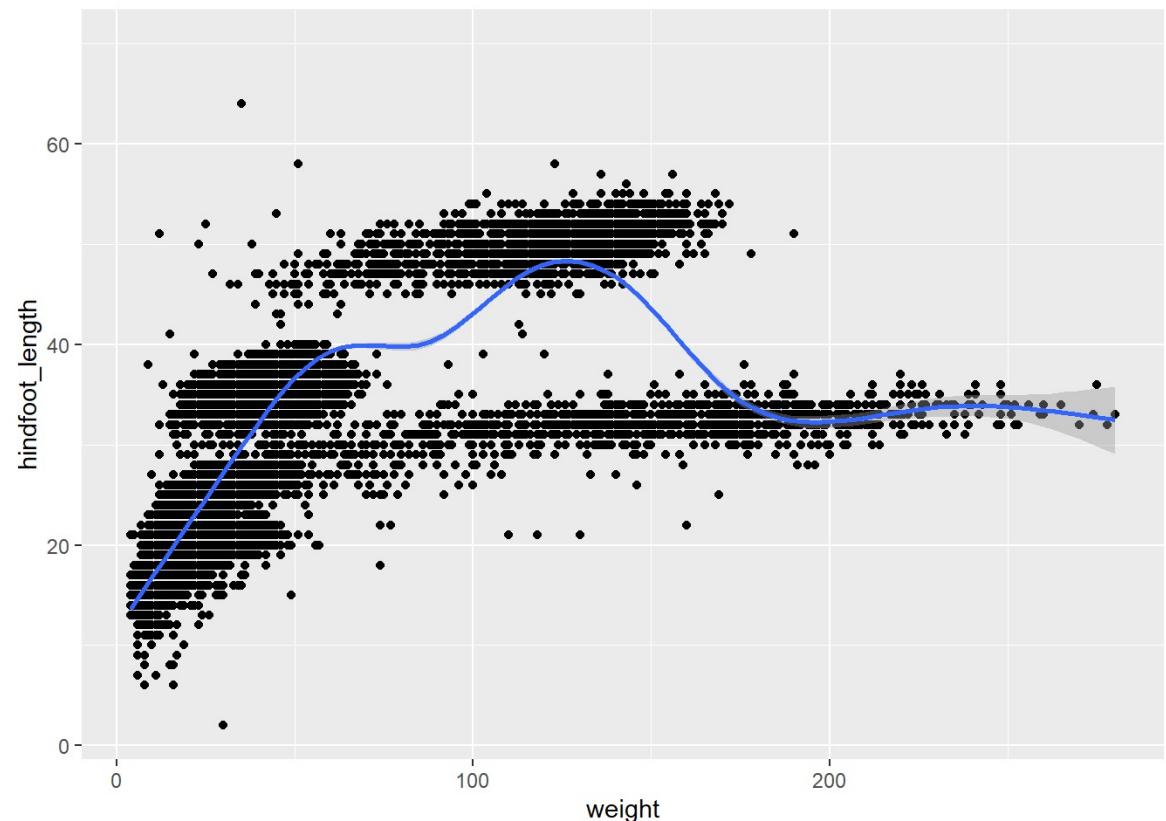
```
surveys_plot +  
  geom_smooth()
```



# 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. Uses this dataset to plot **weight** on x-axis and **hindfoot\_length** on y-axis in a **geom\_density2d** plot

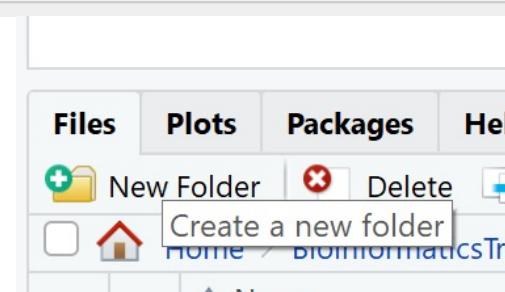
# Saving a plot to a file

```
#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 the folder ‘img\_output’ beforehand!  
Create from within RStudio



# DATA TRANSFORMATION

dplyr : go wrangling

