

Introduction to use of R in Data Analysis: Data cleaning, Data exploration and Data visualization

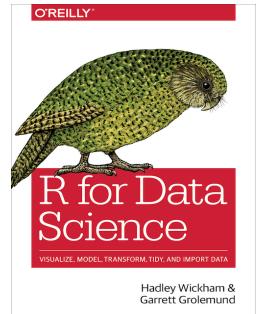
Umar Ahmad

September 30, 2020



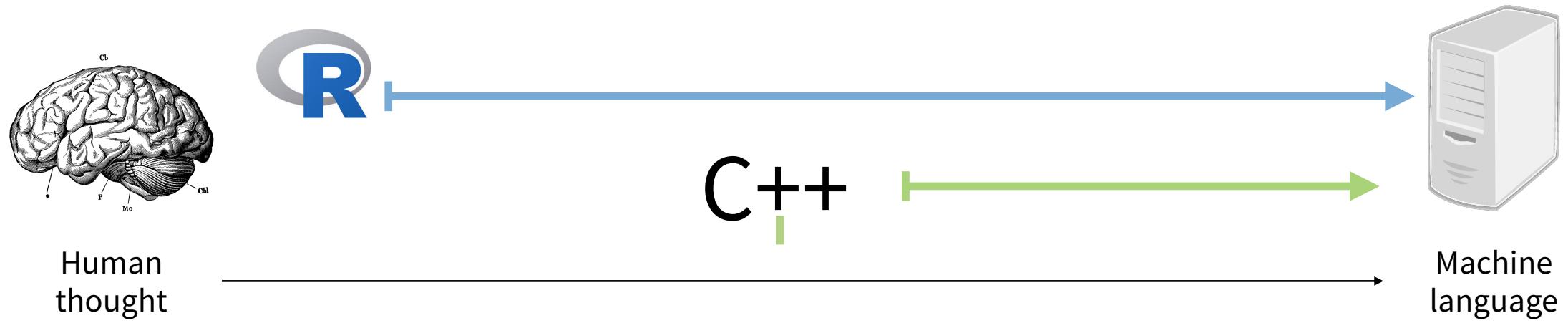
Data Cleaning

Umar Ahmad



**This webinar is will be based on
R for Data Science <http://r4ds.had.co.nz/>**

R - A computer language for scientists



RStudio

R: Engine



RStudio: Dashboard



Courtesy [Modern Dive](#)

RStudio

The screenshot displays the RStudio interface with the following components:

- Script Editor (Top Left):** Shows an R script titled "VisualisingData.R" containing code for importing data from CSV files using the `readr` package.
- Environment Browser (Top Right):** Shows the global environment with objects like `cancer`, `df`, `df2`, `f`, `gg`, `mtcars`, `mydata`, `mydata2`, `myplot`, and `ToothGrowth`.
- File Browser (Bottom Right):** Shows the project structure under "Data-visualization".
- Console (Bottom Left):** Displays the output of running the R script, including column specifications, a warning message about duplicated column names, and the head of the `cancer` dataset.

```
15
16 # IMPORTING WITH READR #####
17
18 # CSV
19
20 library(readr)
21 cancer <- read_csv("data/cancer.csv")
22 head(cancer)
23
24 cancer <- read.csv(file.choose()) # difference ways to import data
25
26 # You can also create sample data which would be used further to demonstrate
27 # data exploration techniques. The program below creates random observations
28 # with replacement.
29
30 cancers = data.frame(Q1 = sample(1:6, 100, replace = TRUE),
31 Q2 = sample(1:6, 100, replace = TRUE),
32 Q3 = sample(1:6, 100, replace = TRUE))
33
34 # IMPORTING WITH READR
```

```
~/Data-visualization/
Parsed with column specification:
cols(
  .default = col_double()
)
See spec(...) for full column specifications.
Warning message:
Duplicated column names deduplicated: 'Q14' => 'Q14_1' [23], 'Q24' => 'Q24_1' [33], 'Q34' => 'Q34_1' [44]
> head(cancer)
# A tibble: 6 x 61
   Age Gender status Education `job role` `How long` Stage type treatment Q1   Q2   Q3
   <dbl> <dbl>
1    45     1     5      4      2     24     2     1      2     3     1     5
2    48     1     5      5      4     21     2     1      2     1     2     2
3    41     1     5      5      4     36     2     1      1     3     4     3
4    58     1     5      5      1     3      2     2      2     5     2     6
5    50     1     5      4      1     36     2     2      2     3     3     5
6    53     1     5      5      1     48     4     2      2     5     5     6
# ... with 49 more variables: Q4 <dbl>, Q5 <dbl>, Q6 <dbl>, Q7 <dbl>, Q8 <dbl>, Q9 <dbl>,
# Q14 <dbl>, Q11 <dbl>, Q12 <dbl>, Q13 <dbl>, Q14_1 <dbl>, Q15 <dbl>, Q16 <dbl>, Q17 <dbl>,
# Q18 <dbl>, Q19 <dbl>, Q24 <dbl>, Q21 <dbl>, Q22 <dbl>, Q23 <dbl>, Q24_1 <dbl>, Q25 <dbl>,
# WSCTotal_Score <dbl>, Q26 <dbl>, Q27 <dbl>, Q28 <dbl>, Q29 <dbl>, Q34 <dbl>, Q31 <dbl>,
# Q32 <dbl>, Q33 <dbl>, Q34_1 <dbl>, Q35 <dbl>, Q36 <dbl>, Q37 <dbl>, Q38 <dbl>, Q39 <dbl>,
# Q44 <dbl>, Q41 <dbl>, Q42 <dbl>, WRITotal_Score <dbl>, BQ1 <dbl>, BQ2 <dbl>, BQ3 <dbl>,
# BQ4 <dbl>, BQ5 <dbl>, BQ6 <dbl>, BQ7 <dbl>, GAD7Total_Score <dbl>
```

```
> head(cancer)
```

```
# A tibble: 6 x 61
```

	Age	Gender	status	Education	`job role`	`How long`	Stage	type	treatment	Q1	Q2	Q3
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	45	1	5	4	2	24	2	1	2	3	1	5
2	48	1	5	5	4	21	2	1	2	1	2	2
3	41	1	5	5	4	36	2	1	1	3	4	3
4	58	1	5	5	1	3	2	2	2	5	2	6
5	50	1	5	4	1	36	2	2	2	3	3	5
6	53	1	5	5	1	48	4	2	2	5	5	6

```
# ... with 49 more variables: Q4 <dbl>, Q5 <dbl>, Q6 <dbl>, Q7 <dbl>, Q8 <dbl>, Q9 <dbl>,
#   Q14 <dbl>, Q11 <dbl>, Q12 <dbl>, Q13 <dbl>, Q14_1 <dbl>, Q15 <dbl>, Q16 <dbl>, Q17 <dbl>,
#   Q18 <dbl>, Q19 <dbl>, Q24 <dbl>, Q21 <dbl>, Q22 <dbl>, Q23 <dbl>, Q24_1 <dbl>, Q25 <dbl>,
#   WSCTotal_Score <dbl>, Q26 <dbl>, Q27 <dbl>, Q28 <dbl>, Q29 <dbl>, Q34 <dbl>, Q31 <dbl>,
#   Q32 <dbl>, Q33 <dbl>, Q34_1 <dbl>, Q35 <dbl>, Q36 <dbl>, Q37 <dbl>, Q38 <dbl>, Q39 <dbl>,
#   Q44 <dbl>, Q41 <dbl>, Q42 <dbl>, WRITotal_Score <dbl>, BQ1 <dbl>, BQ2 <dbl>, BQ3 <dbl>,
#   BQ4 <dbl>, BQ5 <dbl>, BQ6 <dbl>, BQ7 <dbl>, GAD7Total_Score <dbl>
```

Why do we clean your data?

- Accurate - Is the data true
- Uniform - Are all units within a column and across datasets measured the same way
- Consistent - Is data consistent across variables
- Valid - Does the data conform to constraints
- Complete - Missingness will happen but there are ways to mitigate this

Other good data cleaning rules

- Only one piece of information per column
- No unclear values in cells
- Descriptive variable names with no spaces
- No characters (ex: \$) in columns unless it is a string variable
- No duplicate entries

Typical data cleaning steps

1. Read in file/s and explore data
2. Drop columns (Ex: De-identify data - remove name)
3. Rename columns (Ex: Descriptive names, remove spaces)
4. Filter data (Ex: remove those with missing IDs)
5. Remove duplicates
6. Transform/create cols (Ex: string->numeric, or remove \$ or %)
7. Recode variables (Ex: NA->0, reverse code likert scale)
8. Transform data (Ex: from wide to long or long to wide)
9. Merge data and/or append data
10. Add variable labels
11. Make codebook
12. Export data

Associated packages and functions

Step	Package::Function
Read file	<code>readxl::read_excel; readr::read_csv</code>
Read files from folder	<code>list.files, lapply (base)</code>
Explore data	<code>dplyr::glimpse; names, str, summary, table (base)</code>
Explore cont.	<code>skimr::skim; janitor::tabyl</code>
Select cols	<code>dplyr:: select; starts_with, contains, ends_with</code>
Rename cols	<code>purr::set_names; setNames (base R); dplyr::rename dplyr::filter</code>

Packages

Packages add functionality that is not present in base R. They're where much of the power of R is found.

R: A new phone



R Packages: Apps you can download



Courtesy [Modern Dive](#)

Common data cleaning (cont.)

Step	Package::Function
Remove duplicate rows	dplyr::distinct
Create/Transform cols	dplyr::mutate, stringr::str_remove, str_extract; tidyverse::extract
Split column	tidyverse::separate; stringr::str_split
Concatenate 2 cols	paste0 (base); tidyverse::glue
Change col class	lubridate::mdy; as.numeric, as.string, as.factor (base)
Recode cols	dplyr::recode, na_if; tidyverse::replace_na; ifelse (base)
Add value labels	labelled::labelled

Tidyverse (<https://www.tidyverse.org>)

Tidyverse

Packages Blog Learn Help Contribute

R packages for data science

The tidyverse is an opinionated **collection of R packages** designed for data science. All packages share an underlying design philosophy, grammar, and data structures.

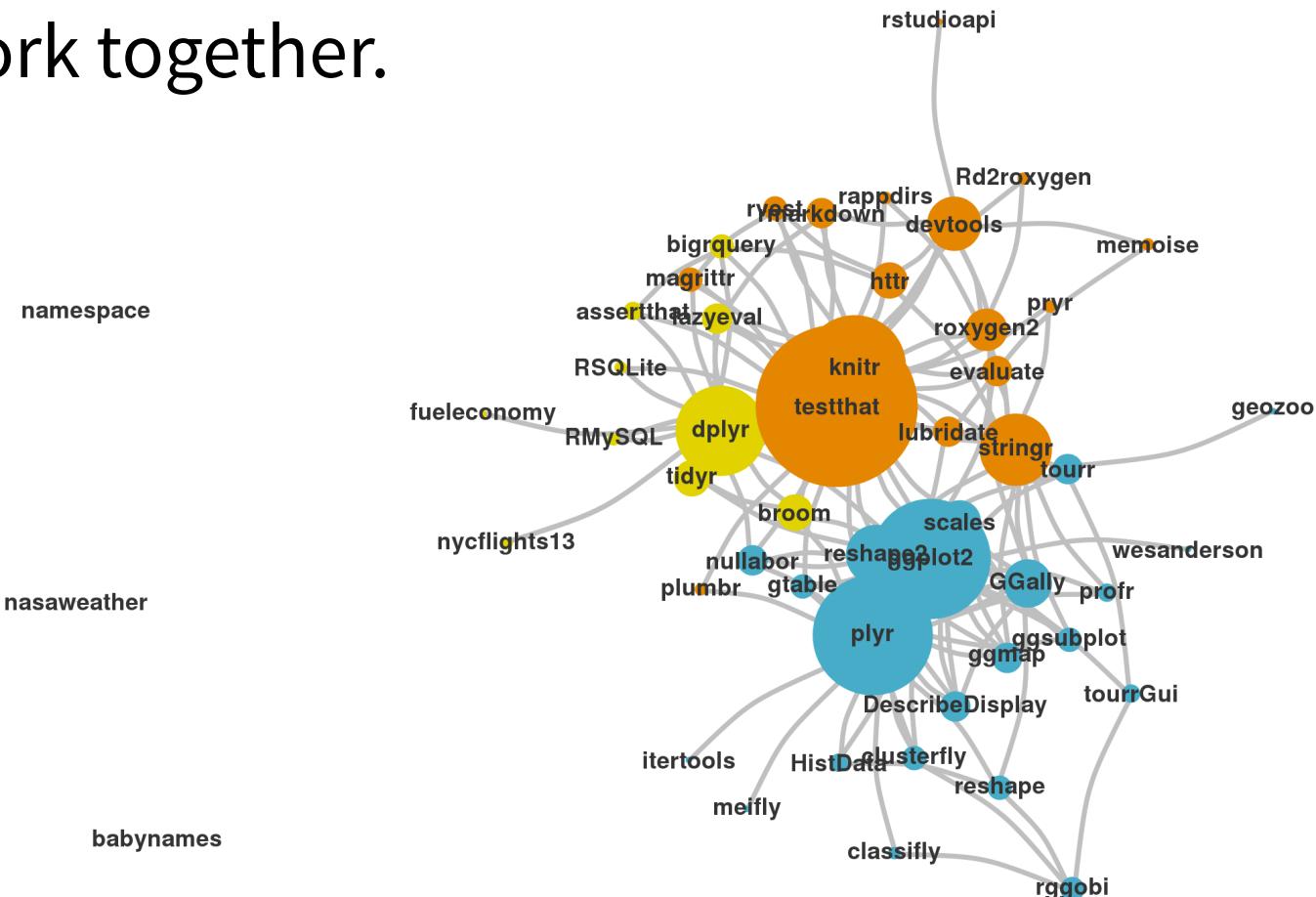
Install the complete tidyverse with:

```
install.packages("tidyverse")
```

Learn the tidyverse

The Tidyverse

A collection of modern R packages that share common philosophies, embed best practices, and are designed to work together.



tidyverse

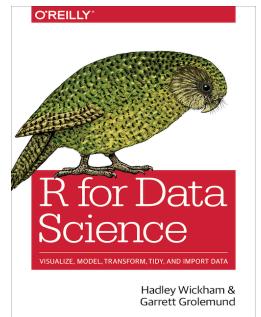
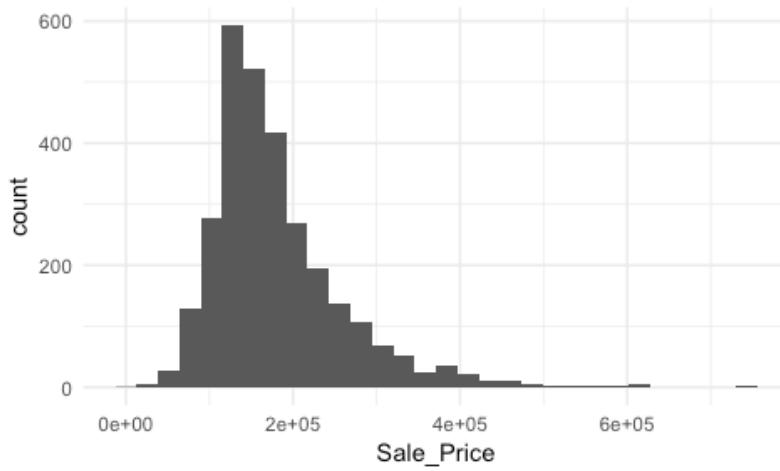


An R package that serves as a short cut for installing and loading the components of the tidyverse.

```
library("tidyverse")
```

Data Exploration

Umar Ahmad



This webinar is will be based on
R for Data Science <http://r4ds.had.co.nz/>

Examine the Data

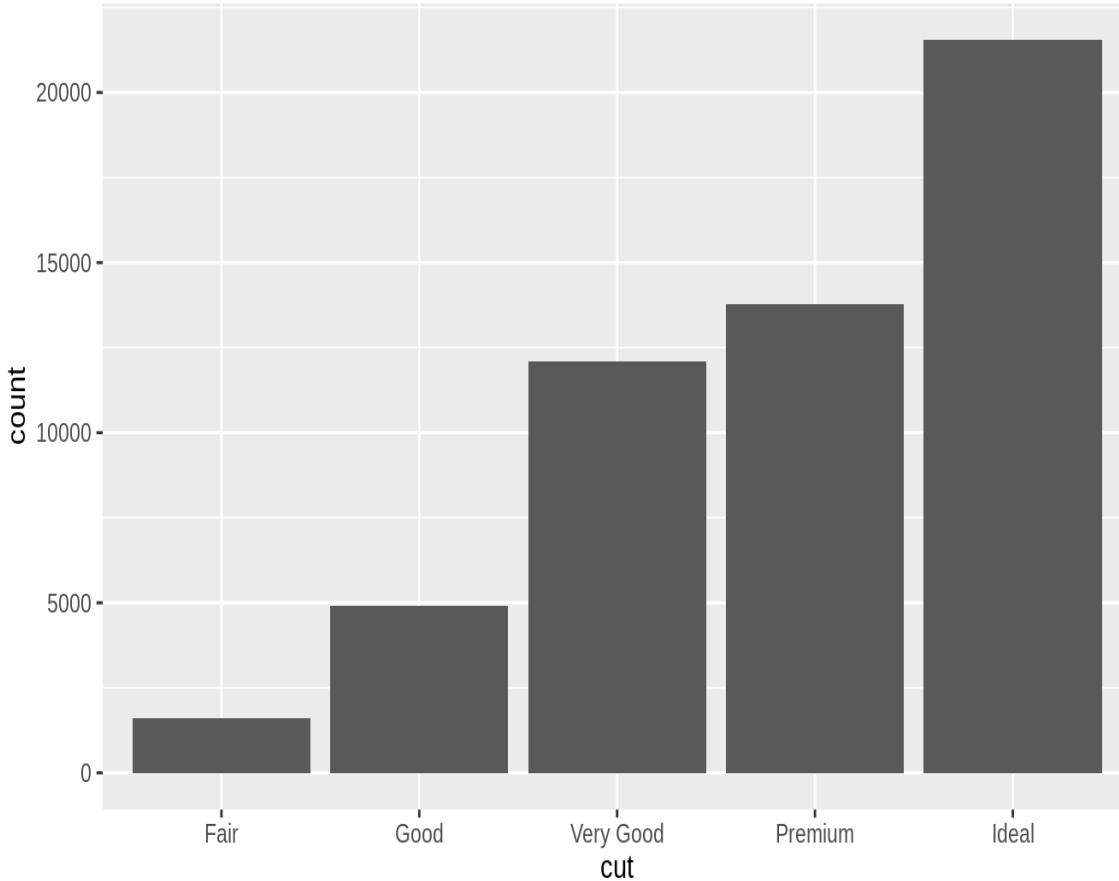
If you type the name of your data frame (i.e. cancer), R will output the following:

```
# A tibble: 6 x 61
  Age Gender status Education `job role` `How long` Stage type treatment   Q1   Q2   Q3
  <dbl> <dbl>
1 45     1     5     4     2     24    2     1     2     3     1     5
2 48     1     5     5     4     21    2     1     2     1     2     2
3 41     1     5     5     4     36    2     1     1     3     4     3
4 58     1     5     5     1     3     2     2     2     5     2     6
5 50     1     5     4     1     36    2     2     2     3     3     5
6 53     1     5     5     1     48    4     2     2     5     5     6
# ... with 49 more variables: Q4 <dbl>, Q5 <dbl>, Q6 <dbl>, Q7 <dbl>, Q8 <dbl>, Q9 <dbl>,
# Q14 <dbl>, Q11 <dbl>, Q12 <dbl>, Q13 <dbl>, Q14_1 <dbl>, Q15 <dbl>, Q16 <dbl>, Q17 <dbl>,
# Q18 <dbl>, Q19 <dbl>, Q24 <dbl>, Q21 <dbl>, Q22 <dbl>, Q23 <dbl>, Q24_1 <dbl>, Q25 <dbl>,
# WSCTotal_Score <dbl>, Q26 <dbl>, Q27 <dbl>, Q28 <dbl>, Q29 <dbl>, Q34 <dbl>, Q31 <dbl>,
# Q32 <dbl>, Q33 <dbl>, Q34_1 <dbl>, Q35 <dbl>, Q36 <dbl>, Q37 <dbl>, Q38 <dbl>, Q39 <dbl>,
# Q44 <dbl>, Q41 <dbl>, Q42 <dbl>, WRITotal_Score <dbl>, BQ1 <dbl>, BQ2 <dbl>, BQ3 <dbl>,
# BQ4 <dbl>, BQ5 <dbl>, BQ6 <dbl>, BQ7 <dbl>, GAD7Total_Score <dbl>
```

```
+> head(cancer)
> head(cancer)
# A tibble: 6 x 61
  Age Gender status Education `job role` `How long` Stage type treatment Q1 Q2 Q3
  <dbl> <dbl>
1 45     1     5     4     2     24    2     1     2     3     1     5
2 48     1     5     5     4     21    2     1     2     1     2     2
3 41     1     5     5     4     36    2     1     1     3     4     3
4 58     1     5     5     1     3     2     2     2     5     2     6
5 50     1     5     4     1     36    2     2     2     3     3     5
6 53     1     5     5     1     48    4     2     2     5     5     6
# ... with 49 more variables: Q4 <dbl>, Q5 <dbl>, Q6 <dbl>, Q7 <dbl>, Q8 <dbl>, Q9 <dbl>,
#   Q14 <dbl>, Q11 <dbl>, Q12 <dbl>, Q13 <dbl>, Q14_1 <dbl>, Q15 <dbl>, Q16 <dbl>, Q17 <dbl>,
#   Q18 <dbl>, Q19 <dbl>, Q24 <dbl>, Q21 <dbl>, Q22 <dbl>, Q23 <dbl>, Q24_1 <dbl>, Q25 <dbl>,
#   WSCTotal_Score <dbl>, Q26 <dbl>, Q27 <dbl>, Q28 <dbl>, Q29 <dbl>, Q34 <dbl>, Q31 <dbl>,
#   Q32 <dbl>, Q33 <dbl>, Q34_1 <dbl>, Q35 <dbl>, Q36 <dbl>, Q37 <dbl>, Q38 <dbl>, Q39 <dbl>,
#   Q44 <dbl>, Q41 <dbl>, Q42 <dbl>, WRITotal_Score <dbl>, BQ1 <dbl>, BQ2 <dbl>, BQ3 <dbl>,
#   BQ4 <dbl>, BQ5 <dbl>, BQ6 <dbl>, BQ7 <dbl>, GAD7Total_Score <dbl>
> summary(cancer[4])
Education
Min. :0.0
1st Qu.:4.0
Median :5.0
Mean   :4.4
3rd Qu.:5.0
Max.   :6.0
> |
```

Visualising distribution...

```
ggplot(data = diamonds) +  
  geom_bar(mapping = aes(x = cut))
```

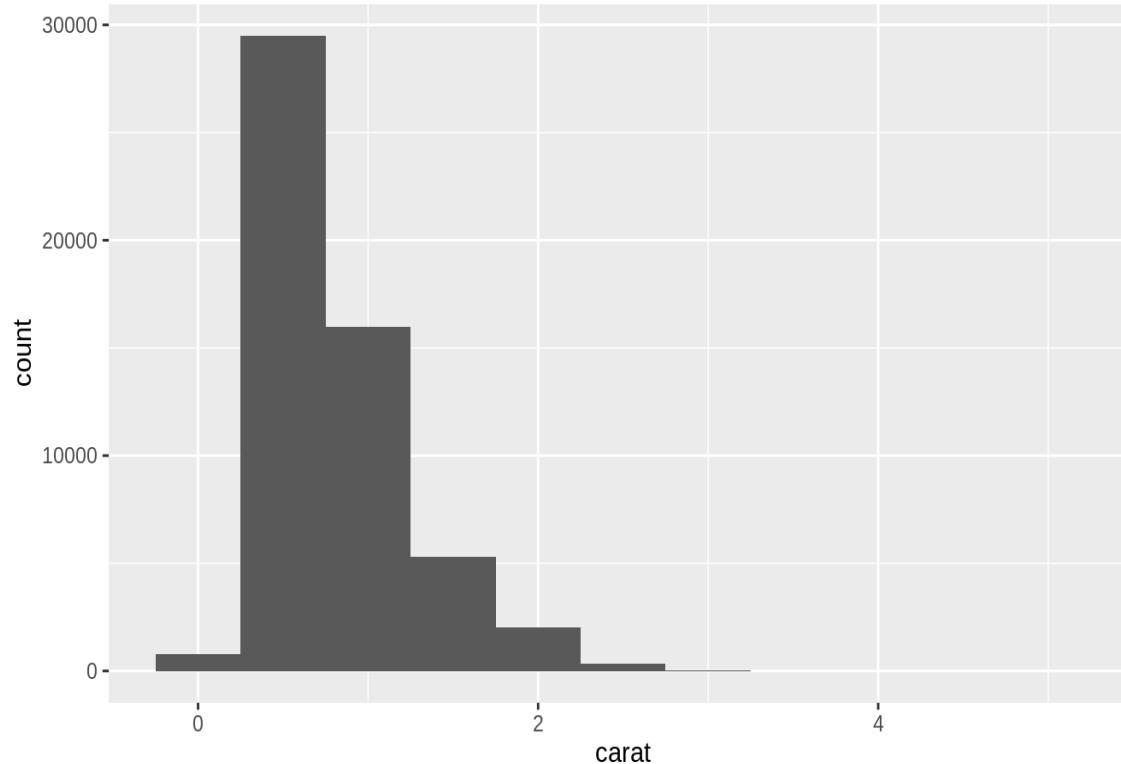


- The height of the bars displays how many observations with each x value
- You can compute these values manually with `dplyr::count()`:

```
iamonds %>%  
  count(cut)  
#> # A tibble: 5 x 2  
#>   cut     n  
#>   <ord> <int>  
#> 1 Fair    1610  
#> 2 Good    4906  
#> 3 Very Good 12082  
#> 4 Premium  13791  
#> 5 Ideal    21551
```

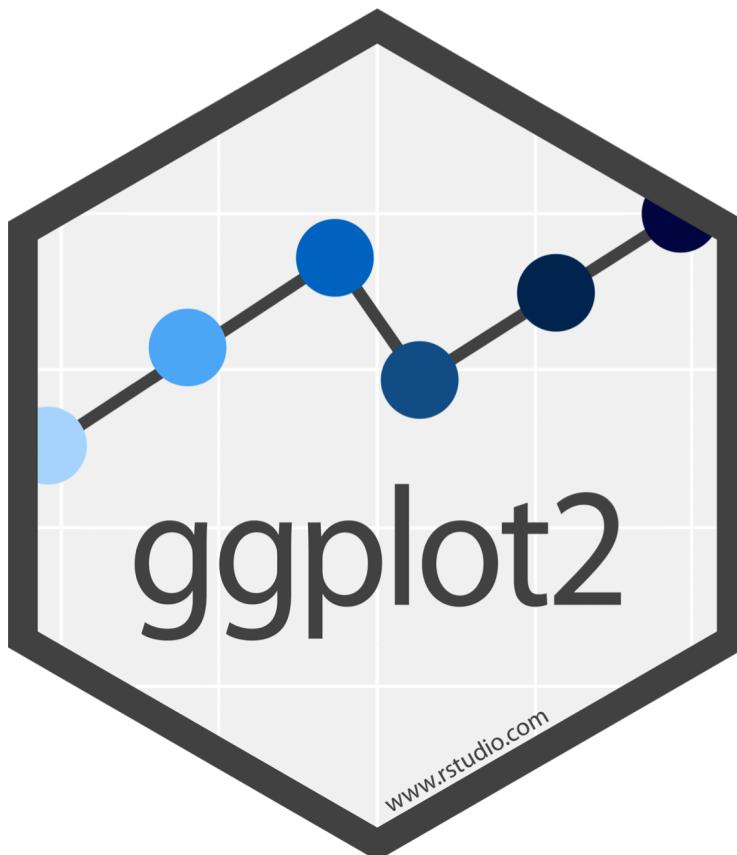
Visualising distribution...

```
ggplot(data = diamonds) +  
  geom_histogram(mapping = aes(x = carat),  
  binwidth = 0.5)
```

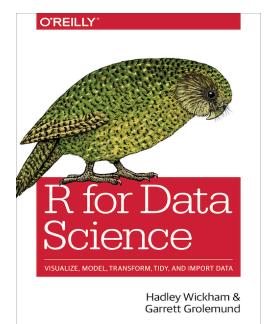


```
diamonds %>%  
  count(cut_width(carat, 0.5))  
#> # A tibble: 11 x 2  
#>   `cut_width(carat, 0.5)`   n  
#>   <fct>                 <int>  
#> 1 [-0.25,0.25]           785  
#> 2 (0.25,0.75]           29498  
#> 3 (0.75,1.25]          15977  
#> 4 (1.25,1.75]           5313  
#> 5 (1.75,2.25]           2002  
#> 6 (2.25,2.75]           322  
#> # ... with 5 more rows
```

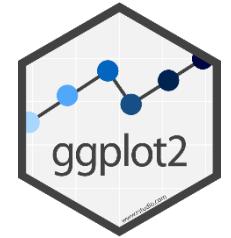
Data visualization with



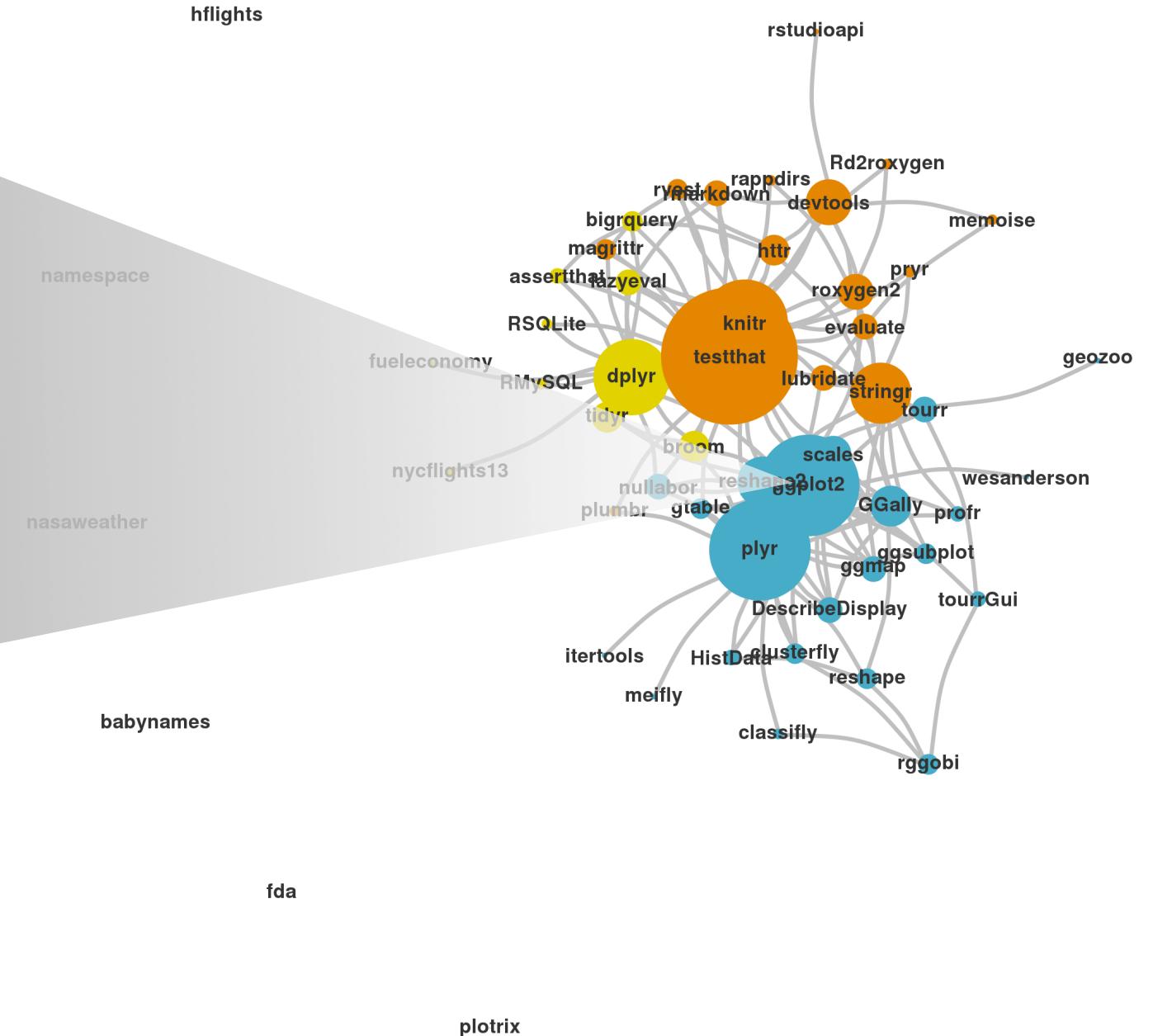
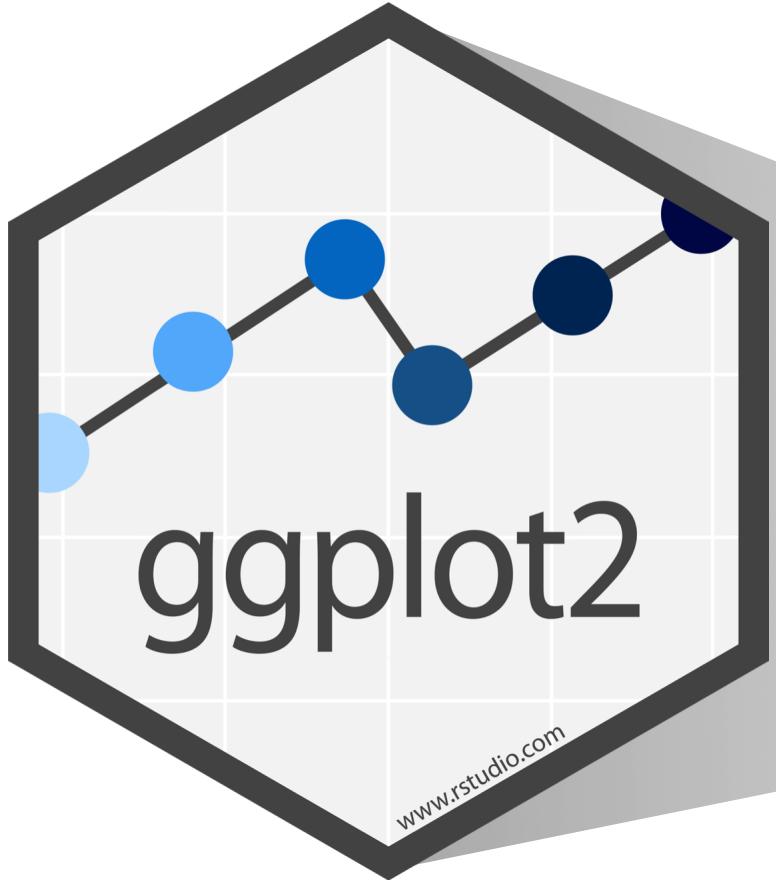
This webinar is will be based on
R for Data Science <http://r4ds.had.co.nz/>



Plotting figures and graphs with ggplot

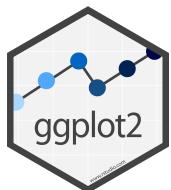


- ggplot is the plotting library for tidyverse
 - Powerful
 - Flexible
- Follows the same conventions as the rest of tidyverse
 - Data stored in tibbles
 - Data is arranged in 'tidy' format
 - Tibble is the first argument to each function



A template

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



```
ggplot(data = bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross))
```

data

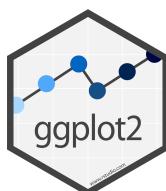
+ before new line

type of layer

aes()

x variable

y variable



Geometries and Aesthetics

- Geometries are types of plot

`geom_point()` Point geometry, (x/y plots, stripcharts etc)

`geom_line()` Line graphs

`geom_boxplot()` Box plots

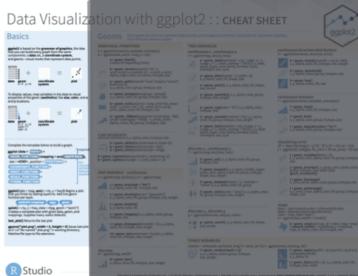
`geom_col()` Barplots

`geom_histogram()` Histogram plots

- Aesthetics are graphical parameters which can be adjusted in a given geometry

geom_ functions

Each requires a mapping argument.



Geoms Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemployed))
b <- ggplot(seals, aes(x = long, y = lat))
```

TWO VARIABLES

continuous x, continuous y

```
e <- ggplot(mpg, aes(cty, hwy))
```

continuous bivariate distribution

```
h <- ggplot(diamonds, aes(carat, price))
```

continuous function

```
i <- ggplot(economics, aes(date, unemployed))
```

visualizing error

```
df <- data.frame(grp = c("A", "B"), fit = 4.5, se = 1.2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

maps

```
data <- data.frame(murder ~ USArrests$Murder,
                    state ~ tolower(rownames(USArrests)))
map <- map_data("state")
k <- ggplot(data, aes(fill = murder))
```

ggplot2

LINE SEGMENTS common aesthetics: x, y, alpha, color, linetype, size

```
b + geom_abline(aes(intercept=0, slope=1))
b + geom_hline(aes(intercept = lat))
b + geom_vline(aes(xintercept = long))

b + geom_segment(aes(yend = lat + 1, xend = long + 1))
b + geom_spoke(aes(angle = 1:1155, radius = 1))
```

ONE VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
```

discrete x, continuous y

```
f <- ggplot(mpg, aes(class, hwy))
```

discrete x, discrete y

```
g <- ggplot(diamonds, aes(cut, color))
```

THREE VARIABLES

```
sealsSz <- with(seals, sqrt(delta_long^2 + delta_lat^2))
```

ggplot2

Aesthetics

```
ggplot(bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, color = clean_test))
```

aesthetic
property

Variable to
map it to

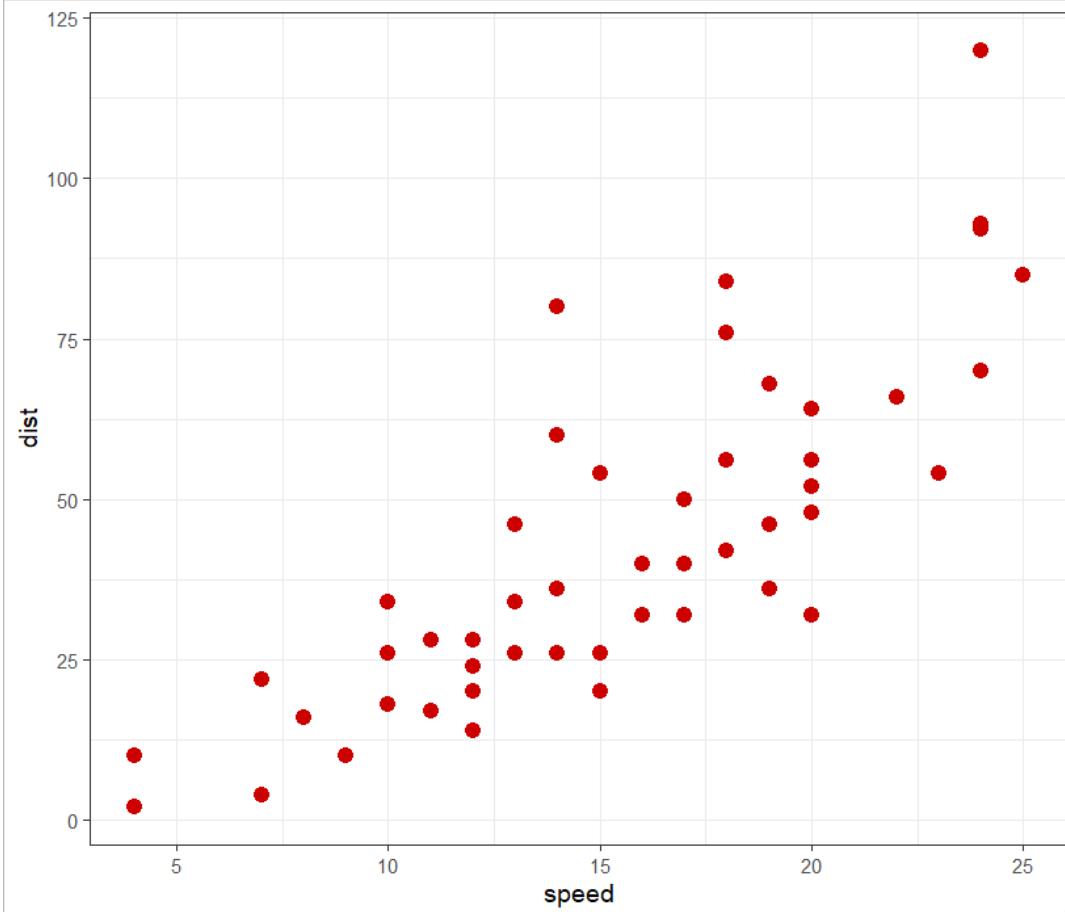
```
ggplot(bechdel) +  
  geom_point(mapping = aes(x = budget, y = domgross, size = clean_test))
```



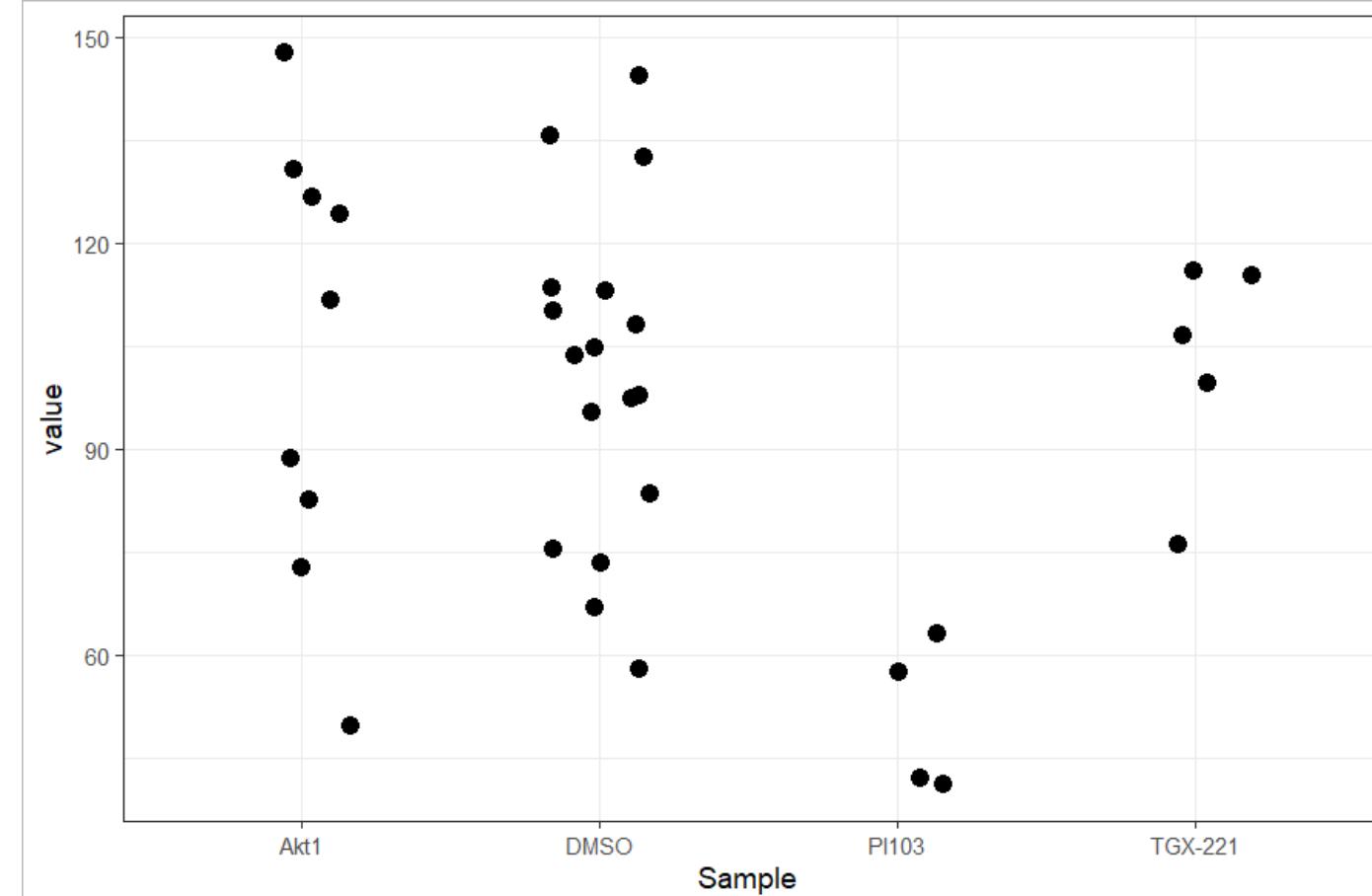
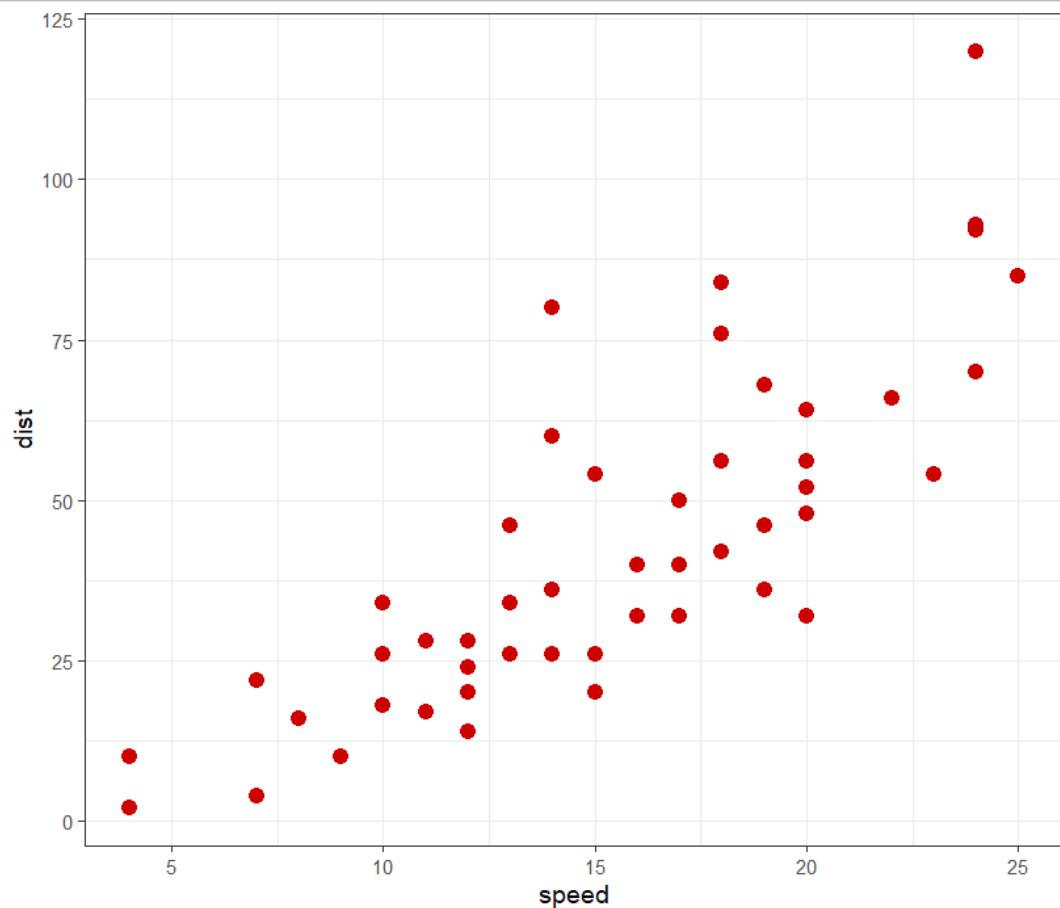
Code structure of a ggplot graph

- Start with a call to `ggplot()`
 - Pass the tibble of data
 - Say which columns you want to use
 - Generates a value which you can store or print
- Say which graphical representation you want to use
 - Points, lines, barplots etc
 - "Add" results to the value from `ggplot`
- Customise labels, colours annotations etc.
- Print the value – draws the plot

Aesthetics for `geom_point()`



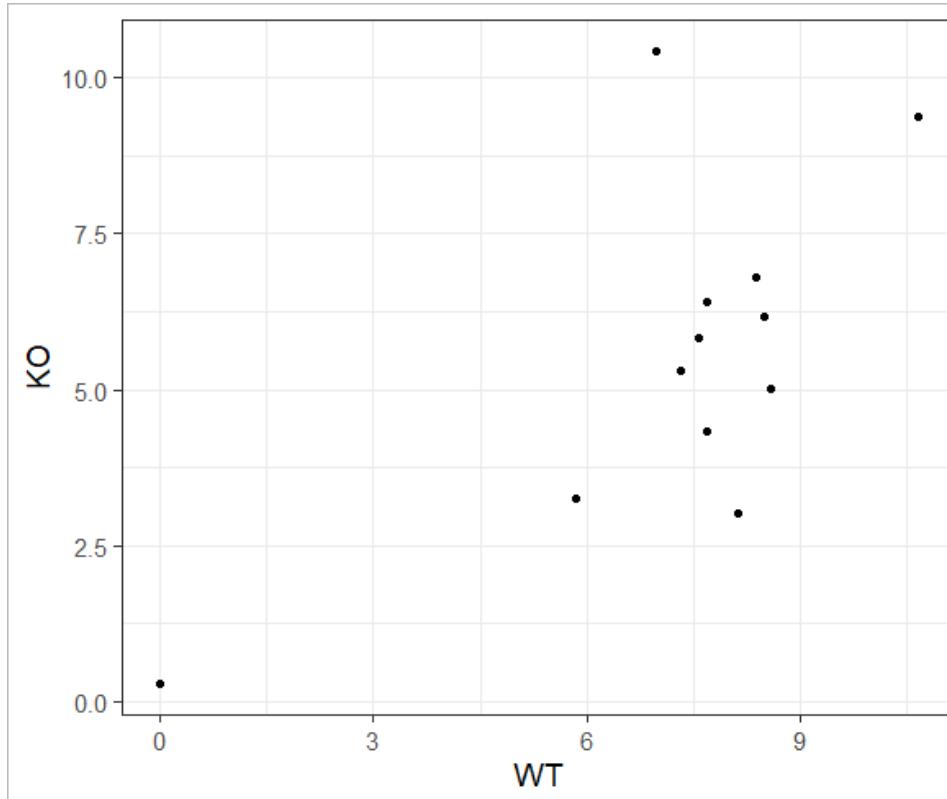
Mappings can be quantitative or categorical



Our first plot...

```
ggplot(expression, aes(x=WT, y=KO)) + geom_point()
```

```
> expression
# A tibble: 12 x 4
  Gene        WT      KO pValue
  <chr>     <dbl>    <dbl>   <dbl>
1 Mial       5.83    3.24    0.1
2 Snrpa      8.59    5.02    0.001
3 Itpkc      8.49    6.16    0.04
4 Adck4      7.69    6.41    0.2
5 Numb1      8.37    6.81    0.1
6 Ltbp4      6.96    10.4    0.001
7 Shkbp1     7.57    5.83    0.1
8 Spnb4      10.7    9.38    0.2
9 Blvrb      7.32    5.29    0.05
10 Pgaml     0        0.285   0.5
11 Sertad3    8.13    3.02    0.0001
12 Sertad1    7.69    4.34    0.01
```

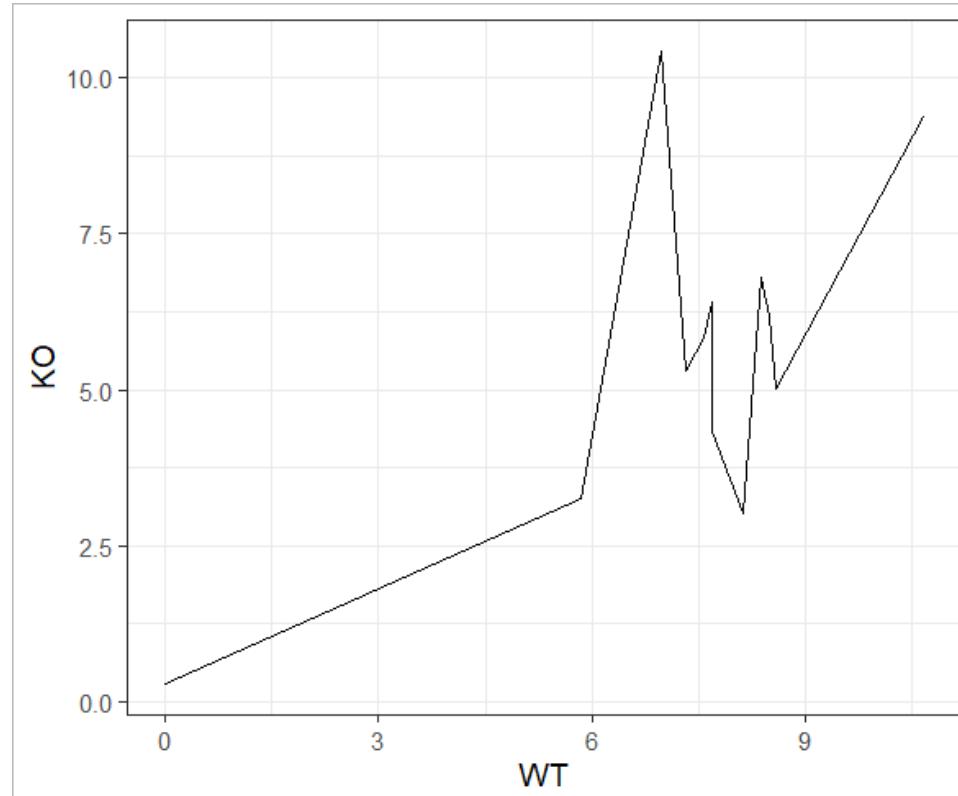


- Identify the tibble with the data you want to plot
- Decide on the geometry (plot type) you want to use
- Decide which columns will modify which aesthetic
- Call `ggplot(aes(...))`
- Add a `geom_xxx` function call

Our second plot...

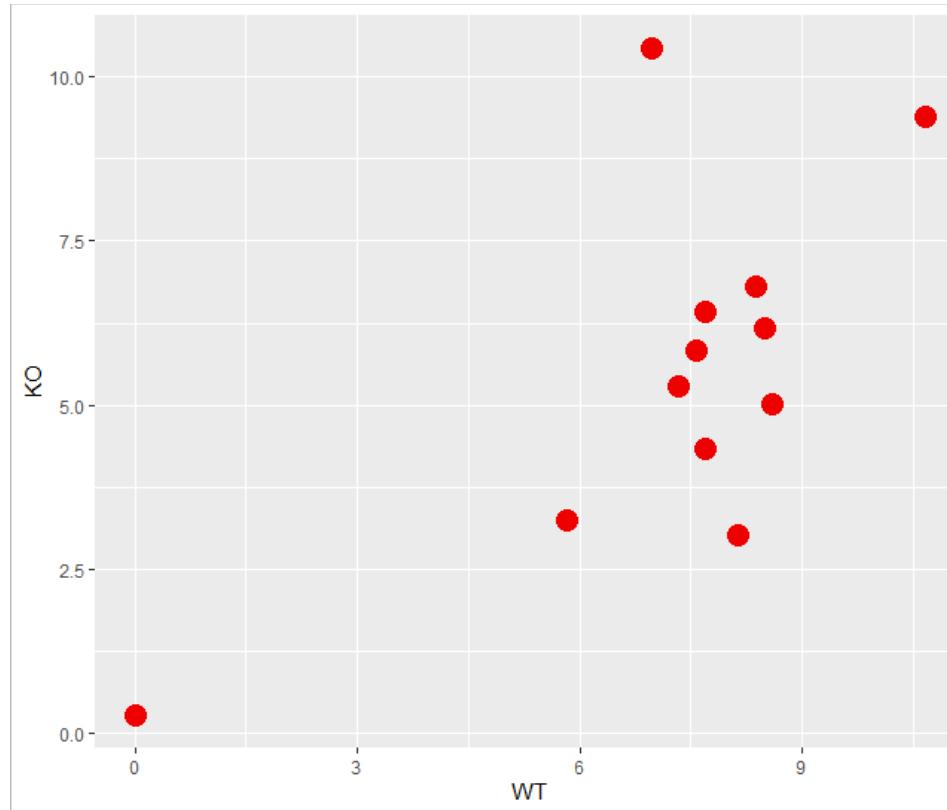
```
ggplot(expression, aes(x=WT, y=KO)) + geom_line()
```

```
> expression
# A tibble: 12 x 4
  Gene        WT      KO pValue
  <chr>     <dbl>   <dbl>  <dbl>
1 Mial       5.83   3.24  0.1
2 Snrpa      8.59   5.02  0.001
3 Itpkc      8.49   6.16  0.04
4 Adck4      7.69   6.41  0.2
5 Numb1      8.37   6.81  0.1
6 Ltbp4      6.96  10.4   0.001
7 Shkbp1     7.57   5.83  0.1
8 Spnb4      10.7   9.38  0.2
9 Blvrb      7.32   5.29  0.05
10 Pgaml      0     0.285  0.5
11 Sertad3    8.13   3.02  0.0001
12 Sertad1    7.69   4.34  0.01
```



Our third plot...

```
expression %>%
  ggplot (aes(x=WT, y=KO)) +
  geom_point(colour="red2", size=5)
```



Other data plot types (geometries)

- Barplots
 - geom_bar
 - geom_col
- Stripcharts
 - geom_jitter
- Distribution Summaries
 - geom_histogram
 - geom_density
 - geom_violin
 - geom_boxplot

Drawing a barplot (`geom_col()` or `geom_bar()`)

- Two different functions – depends on the nature of the data
- If your data has values which represents the height of the bars use `geom_col`
- If your data has individual values and you want the plot to either count them or calculate a summary (usually the mean) then use `geom_bar`

Drawing a bar height barplot (`geom_col()`)

Aesthetics

`geom_bar()` understands the following aesthetics
(required aesthetics are in bold):

- `x`
- `y`
- `alpha`
- `colour`
- `fill`
- `group`
- `linetype`
- `size`

- Plot the expression values for the WT samples for all genes

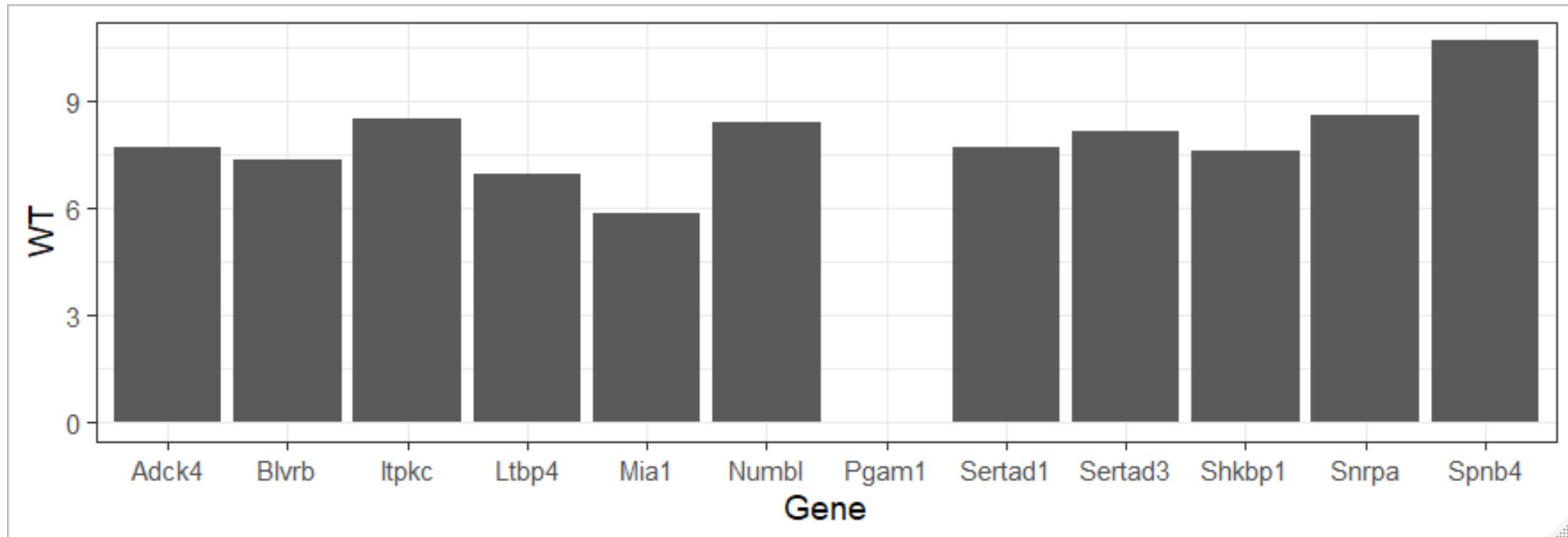
- What is your X?
- What is your Y?

> expression

```
# A tibble: 12 x 4
  Gene      WT      KO pValue
  <chr>    <dbl>   <dbl>  <dbl>
  1 Mial     5.83    3.24   0.1
  2 Snrpa    8.59    5.02   0.001
```

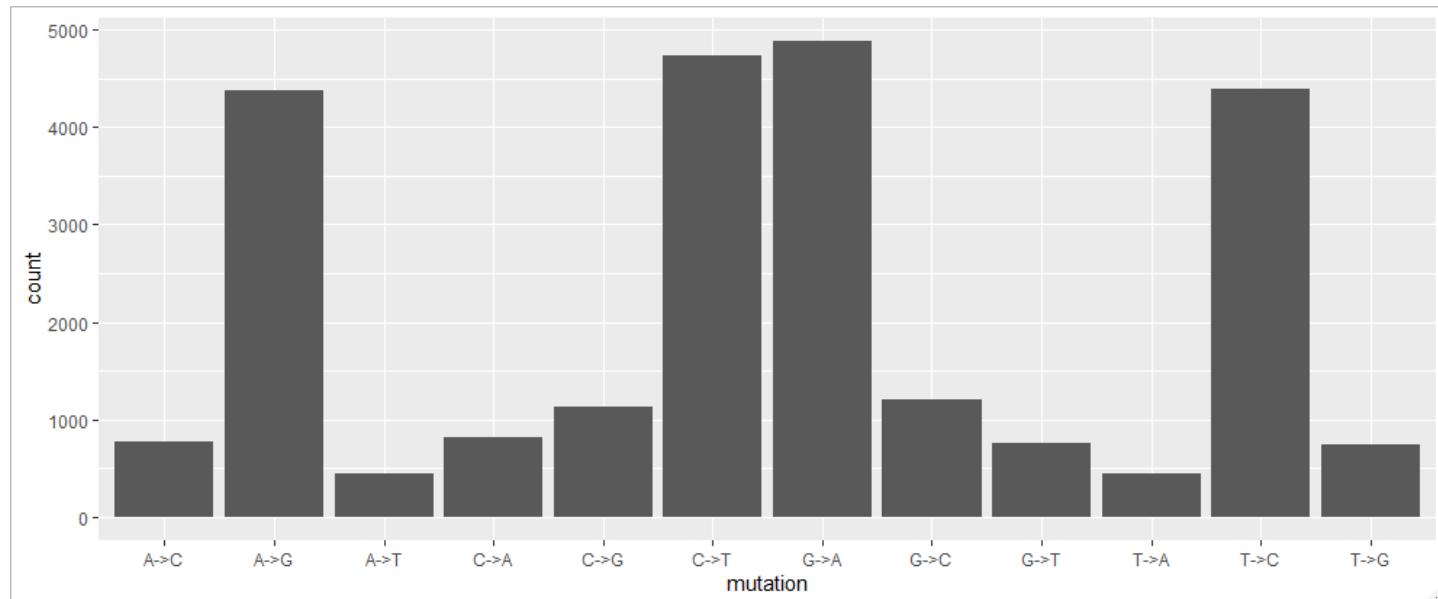
A bar height barplot

```
ggplot(expression, aes(x=Gene, y=WT)) +  
  geom_col()
```



A summarised barplot (geom_bar) - counts

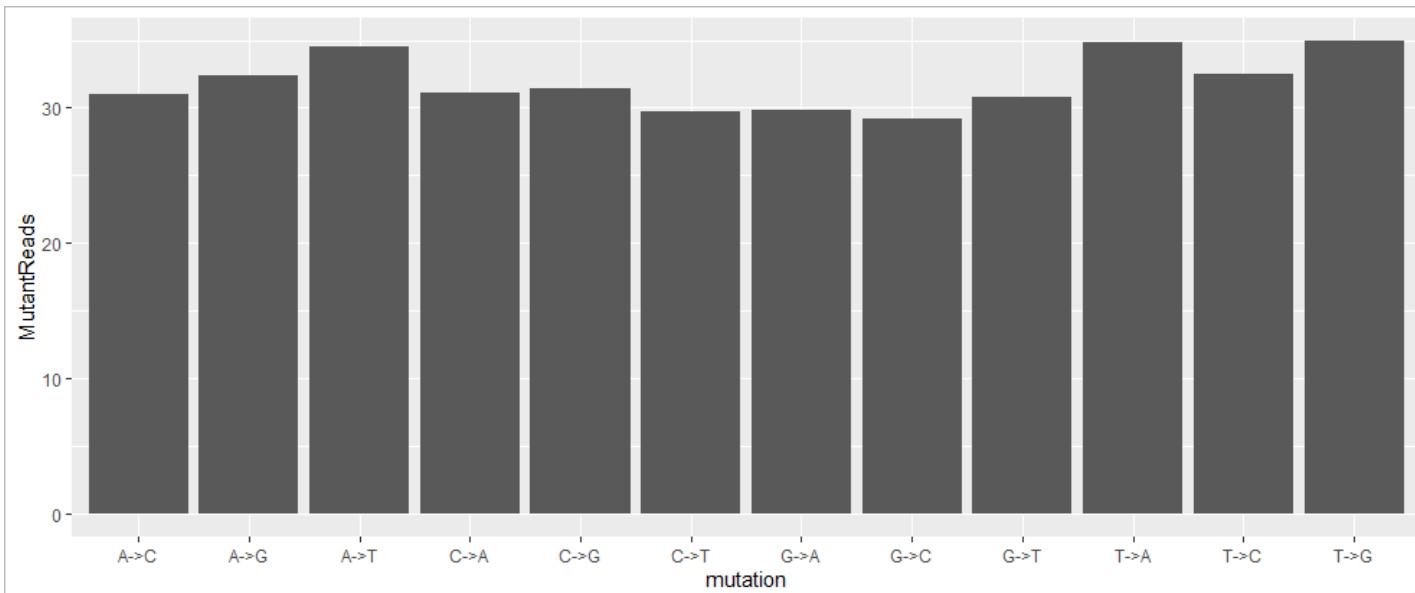
```
mutation.plotting.data %>%
  ggplot(aes(x=mutation)) +
  geom_bar()
```



```
> mutation.plotting.data
# A tibble: 24,686 x 9
  CHR      POS dbSNP      mutation  QUAL GENE    ENST      MutantReads  COVERAGE
  <chr>   <dbl> <chr>      <chr>   <dbl> <chr>    <chr>      <dbl>       <dbl>
1 1       69270 .        A->G      16  OR4F5  ENST00000335137     3          4
2 1       69511 rs75062661 A->G      200 OR4F5  ENST00000335137    24         27
3 1       69761 .        A->T      200 OR4F5  ENST00000335137     8          8
4 1       69897 rs75758884 T->C      59  OR4F5  ENST00000335137     3          3
5 1       877831 rs6672356 T->C     200 SAMD11 ENST00000342066    10         11
6 1       881627 rs2272757 G->A     200 NOC2L   ENST00000327044    52         56
7 1       887801 rs3828047 A->G     200 NOC2L   ENST00000327044    47         48
8 1       888639 rs3748596 T->C     200 NOC2L   ENST00000327044    23         24
9 1       888659 rs3748597 T->C     200 NOC2L   ENST00000327044    17         21
10 1      889158 rs13303056 G->C     200 NOC2L   ENST00000327044   25         28
```

A summarised barplot (geom_bar) - means

```
mutation.plotting.data %>%
  ggplot(aes(x=mutation, y=MutantReads))+
  geom_bar(stat="summary", fun="mean")
```

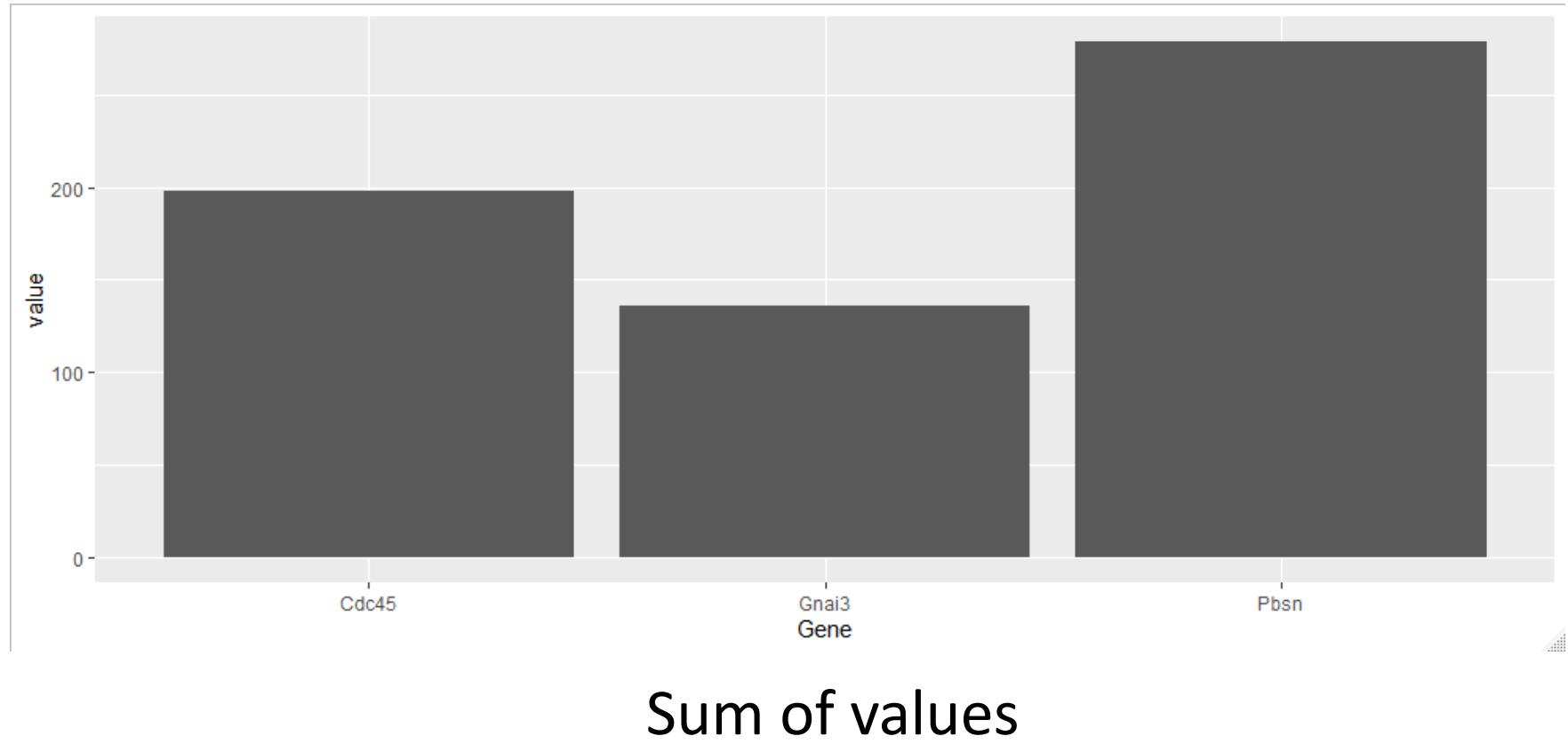


```
> mutation.plotting.data
# A tibble: 24,686 x 9
  CHR      POS dbSNP      mutation  QUAL GENE    ENST      MutantReads  COVERAGE
  <chr>   <dbl> <chr>     <chr>   <dbl> <chr> <chr>      <dbl>       <dbl>
1 1       69270 .        A->G      16  OR4F5  ENST00000335137 3          4
2 1       69511 rs75062661 A->G      200 OR4F5  ENST00000335137 24         27
3 1       69761 .        A->T      200 OR4F5  ENST00000335137 8          8
4 1       69897 rs75758884 T->C      59  OR4F5  ENST00000335137 3          3
5 1       877831 rs6672356 T->C     200 SAMD11 ENST00000342066 10         11
6 1       881627 rs2272757 G->A     200 NOC2L   ENST00000327044 52         56
7 1       887801 rs3828047 A->G     200 NOC2L   ENST00000327044 47         48
8 1       888639 rs3748596 T->C     200 NOC2L   ENST00000327044 23         24
9 1       888659 rs3748597 T->C     200 NOC2L   ENST00000327044 17         21
10 1      889158 rs13303056 G->C    200 NOC2L   ENST00000327044 25         28
```

Stacked and Grouped Barplots

```
bar.group %>%
  ggplot(aes(x=Gene, y=value)) +
  geom_col()
```

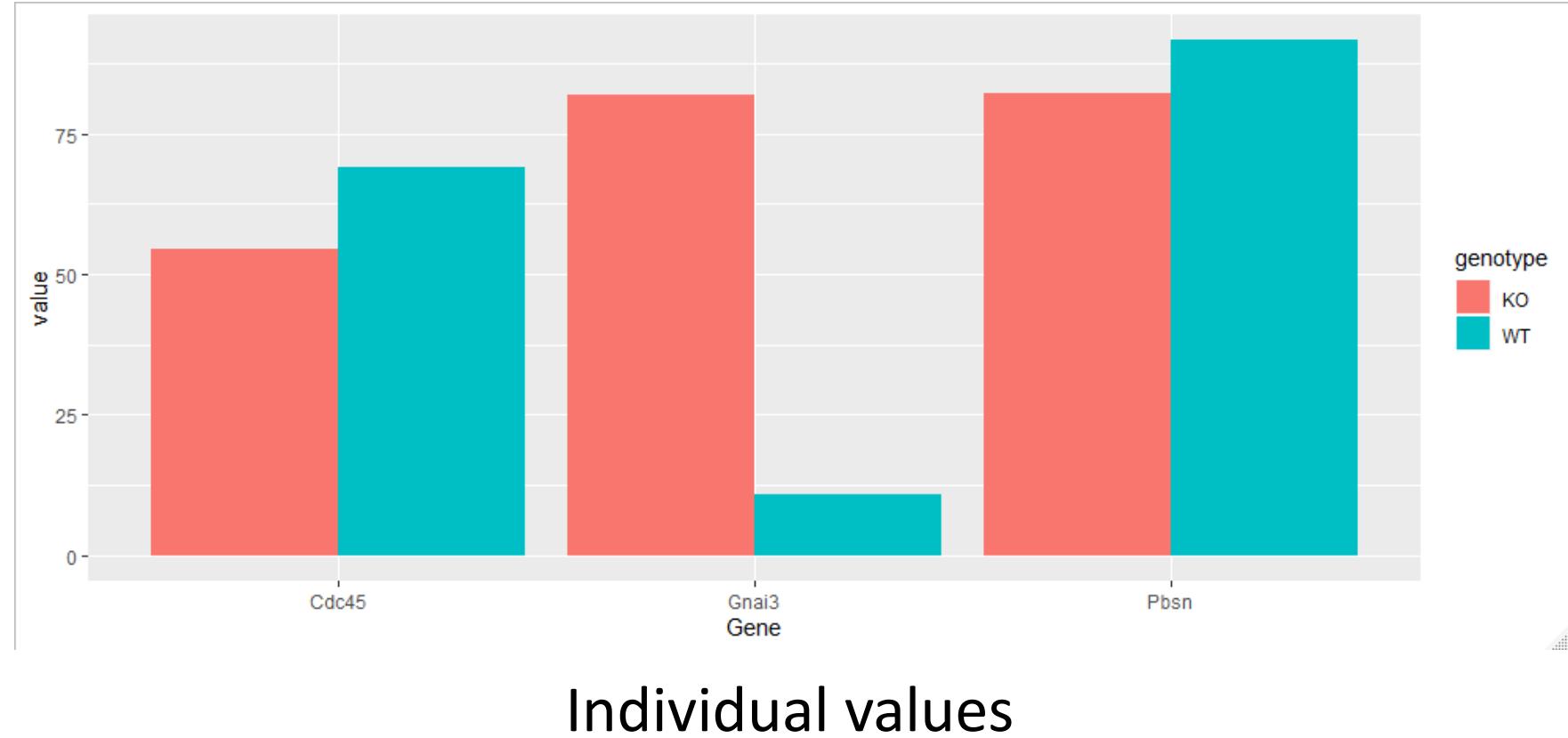
```
> bar.group
# A tibble: 12 x 3
  Gene   genotype value
  <chr> <chr>    <dbl>
1 Gnai3  WT       9.39
2 Pbsn   WT      91.7 
3 Cdc45  WT      69.2 
4 Gnai3  WT      10.9 
5 Pbsn   WT      59.6 
6 Cdc45  WT      36.1 
7 Gnai3  KO      33.5 
8 Pbsn   KO      45.3 
9 Cdc45  KO      54.4 
10 Gnai3 KO     81.9 
11 Pbsn   KO     82.3 
12 Cdc45 KO     38.1
```



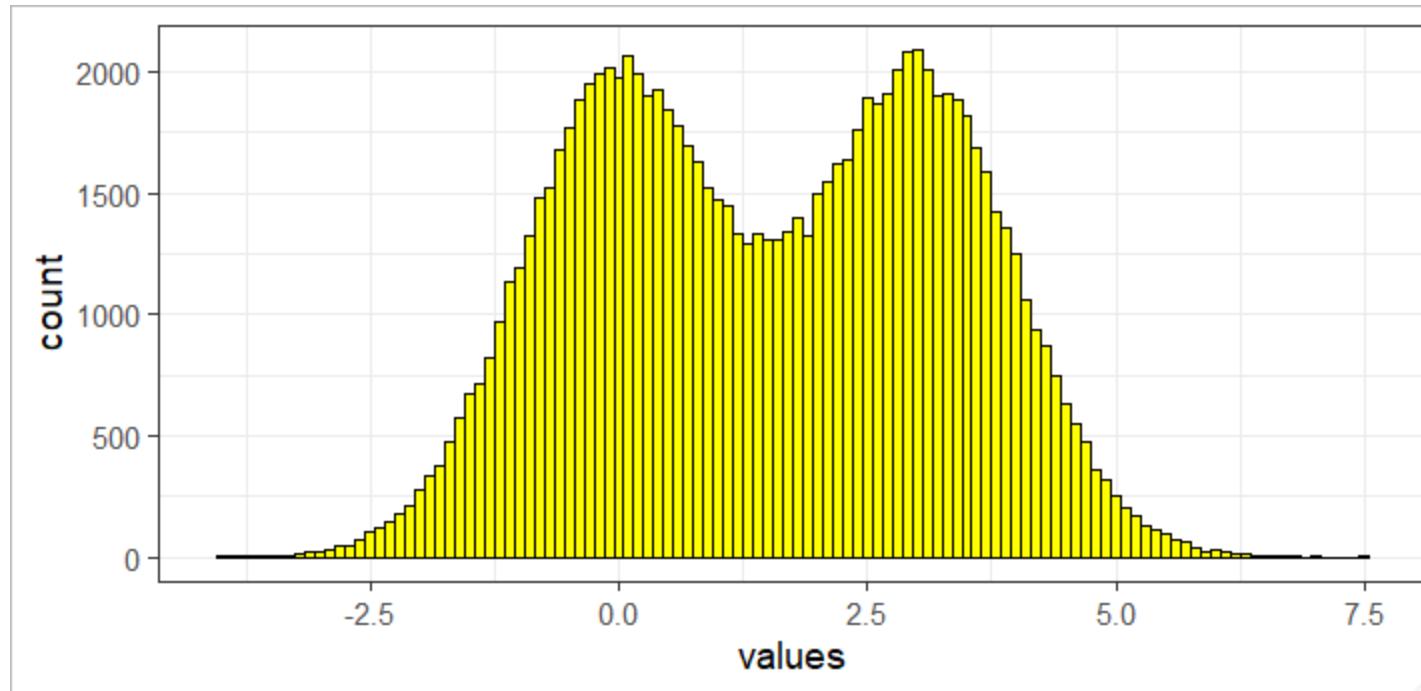
Stacked and Grouped Barplots

```
> bar.group  
# A tibble: 12 x 3  
  Gene   genotype value  
  <chr> <chr>    <dbl>  
1 Gnai3  WT      9.39  
2 Pbsn   WT      91.7  
3 Cdc45  WT      69.2  
4 Gnai3  WT      10.9  
5 Pbsn   WT      59.6  
6 Cdc45  WT      36.1  
7 Gnai3  KO      33.5  
8 Pbsn   KO      45.3  
9 Cdc45  KO      54.4  
10 Gnai3 KO      81.9  
11 Pbsn   KO      82.3  
12 Cdc45 KO      38.1
```

```
bar.group %>%  
  ggplot(aes(x=Gene, y=value, fill=genotype)) +  
  geom_col(position="dodge")
```



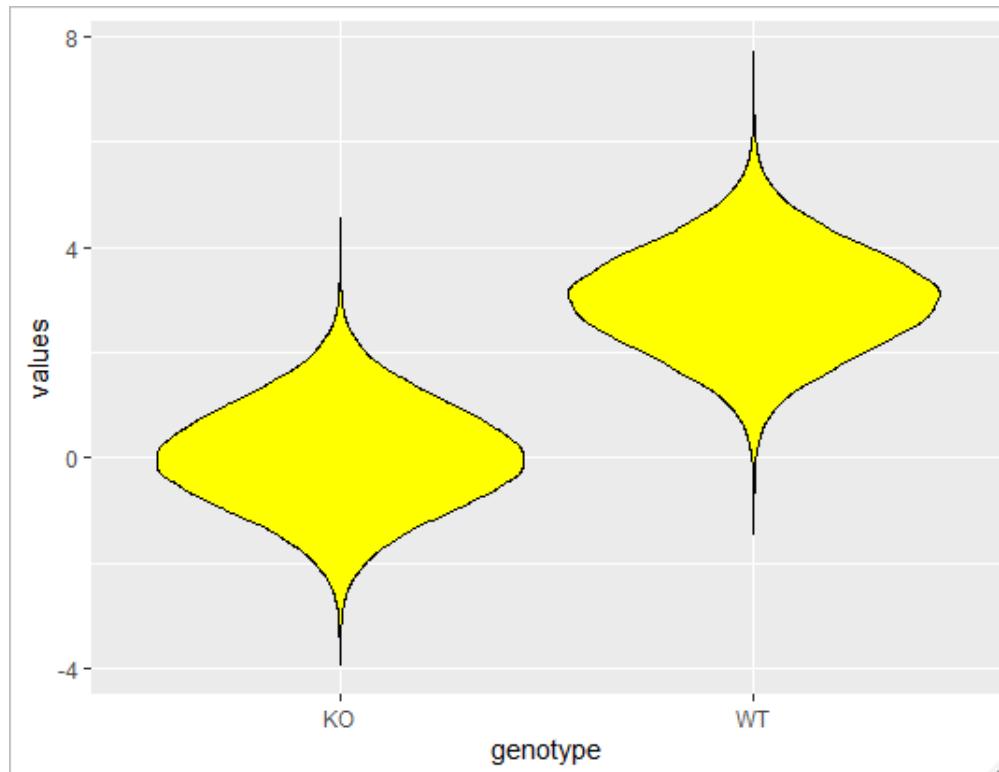
Plotting distributions - histograms



```
many.values %>%
  ggplot(aes(x=values)) +
  geom_histogram(binwidth = 0.1, fill="yellow", colour="black")
```

```
> many.values
# A tibble: 100,000 x 2
  values genotype
  <dbl> <chr>
1 1.90  KO
2 2.39  WT
3 4.32  KO
4 2.94  KO
5 0.728 WT
6 -0.280 WT
7 0.337 WT
8 -1.31  WT
9 1.55  WT
10 1.86  KO
```

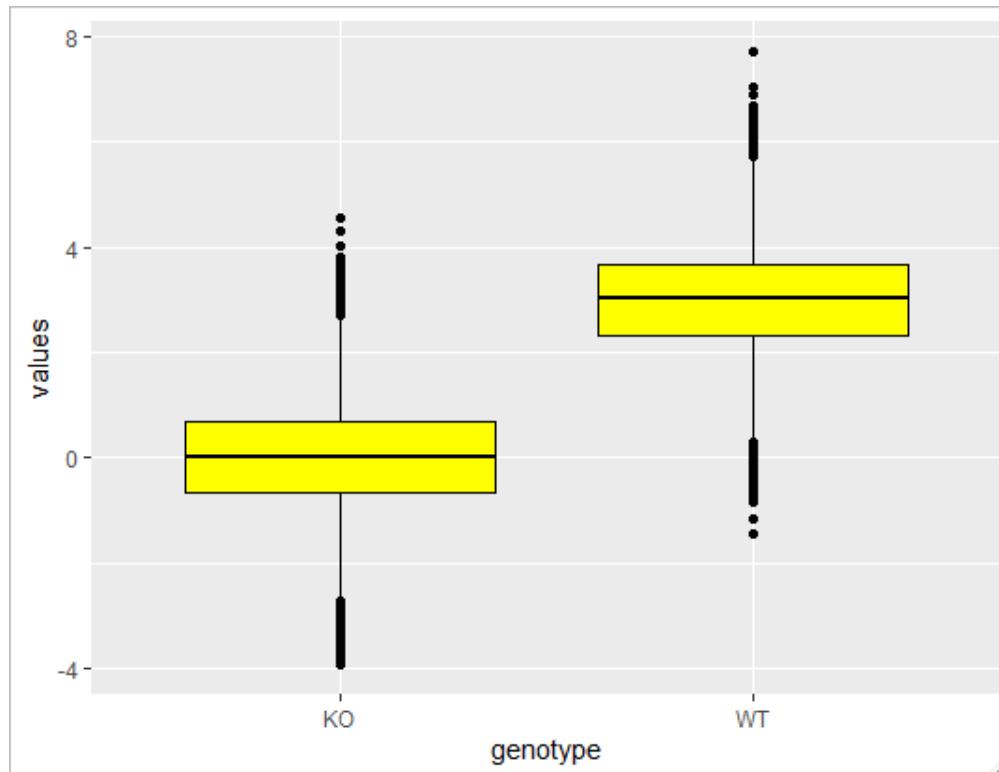
Plotting distributions – violin plots



```
many.values %>%
  ggplot(aes(x=genotype, y=values)) +
  geom_violin(colour="black", fill="yellow")
```

```
> many.values
# A tibble: 100,000 × 2
  values genotype
  <dbl> <chr>
1 1.90  KO
2 2.39  WT
3 4.32  KO
4 2.94  KO
5 0.728 WT
6 -0.280 WT
7 0.337 WT
8 -1.31  WT
9 1.55  WT
10 1.86  KO
```

Plotting distributions – boxplots



```
many.values %>%
  ggplot(aes(x=genotype, y=values)) +
  geom_boxplot(colour="black", fill="yellow")
```

```
> many.values
# A tibble: 100,000 x 2
  values genotype
  <dbl> <chr>
1 1.90  KO
2 2.39  WT
3 4.32  KO
4 2.94  KO
5 0.728 WT
6 -0.280 WT
7 0.337 WT
8 -1.31  WT
9 1.55  WT
10 1.86  KO
```

Annotation, Scaling and Colours

Titles and axis labels

- Can set everything with `labs()`
 - `title="Main title"`
 - `x="X axis"`
 - `y="Y axis"`
- Can use functions to set them individually
 - `ggtitle()`
 - `xlab()`
 - `ylab()`

Changing scaling

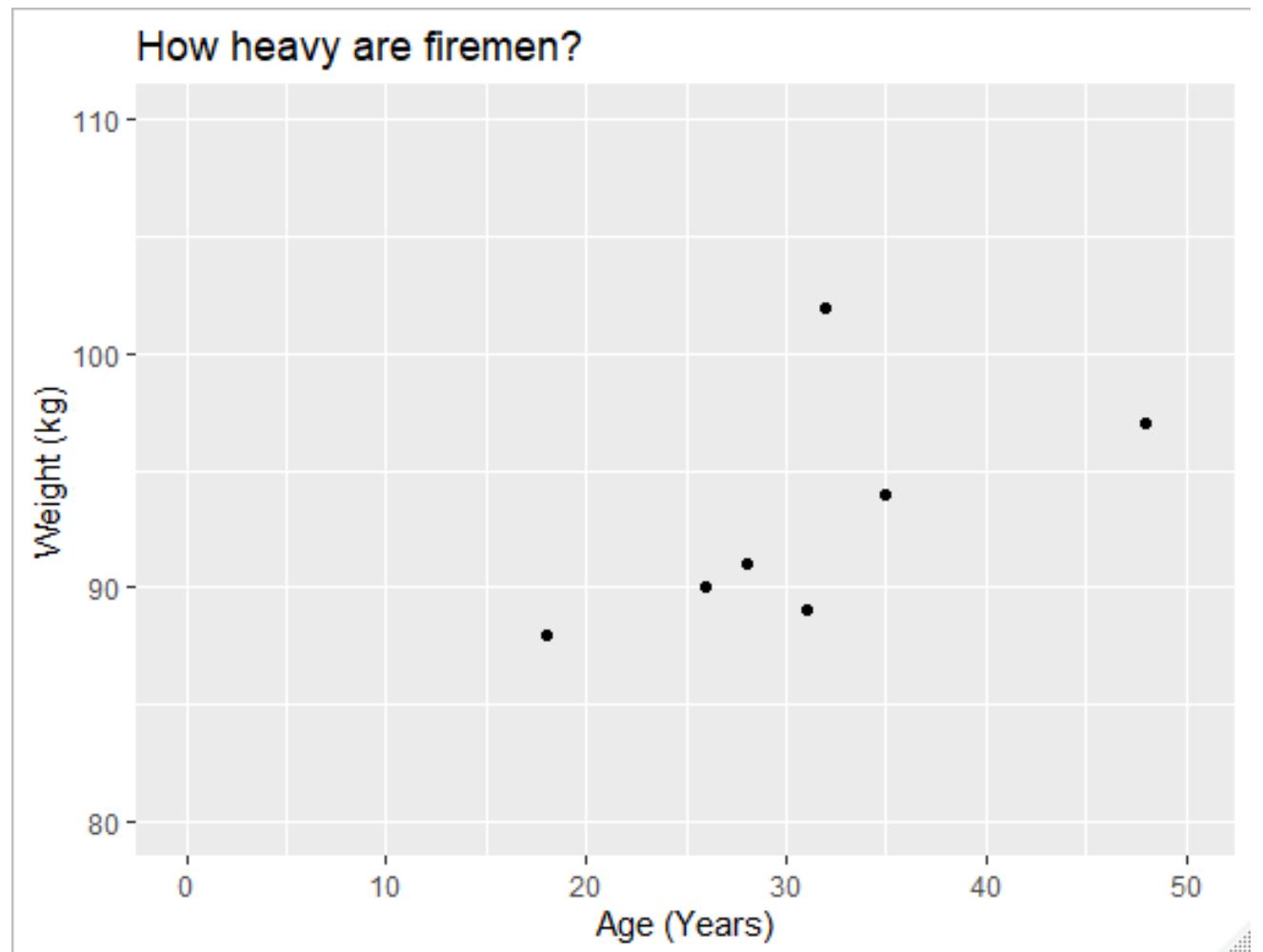
- Alter the data before plotting
 - `mutate(value=log(value))`
- Alter the data whilst plotting
 - `ggplot(aes(log(value)))`
- Alter the scale of the plot
 - Add an option to adjust the scaling of the axis

Axis scaling options

- Transforming scales
 - `scale_x_log10()`
 - `scale_x_sqrt()`
 - `scale_x_reverse()`
- Equivalent `_y_` versions also exist
- Switching axes
 - `coord_flip()`
- Adjusting ranges
 - `scale_x_continuous()`
 - `limits=c(-5, 5)`
 - `breaks=seq(from=-5, by=2, to=5)`
 - `minor_breaks`
 - `labels`
 - `coord_cartesian()`
 - `xlim=c(-5, 5)`
 - `ylim=c(10, 20)`

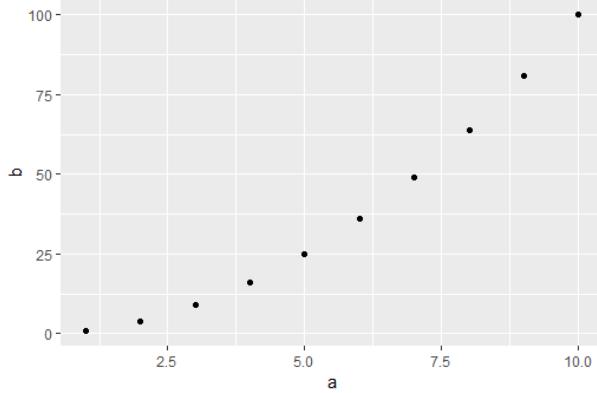
Annotation and scaling example

```
trumpton %>%  
  ggplot(aes(x=Age, y=weight)) +  
  geom_point() +  
  
  xlab("Age (Years)") +  
  ylab("Weight (kg)") +  
  ggtitle("How heavy are firemen?") +  
  
  coord_cartesian(  
    xlim=c(0,50),  
    ylim=c(80,110)  
)
```



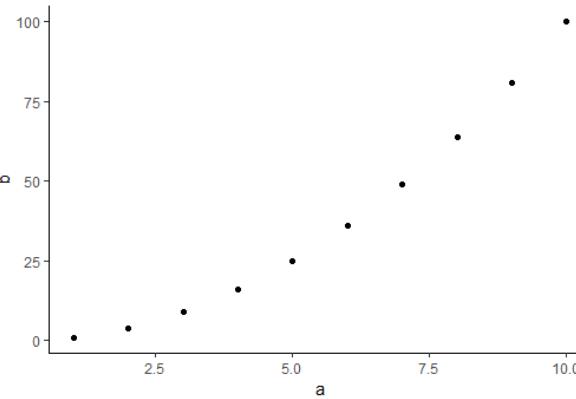
ggPlot Themes

Theme Gray/Grey

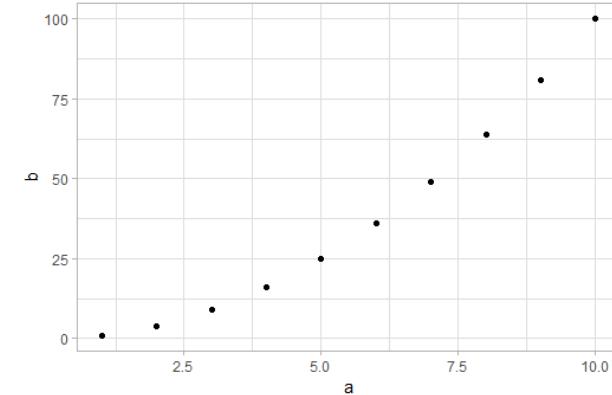


- `theme_grey()`
- `theme_bw()`
- `theme_dark()`
- `theme_light()`
- `theme_minimal()`
- `theme_classic()`
- `theme_linedraw()`

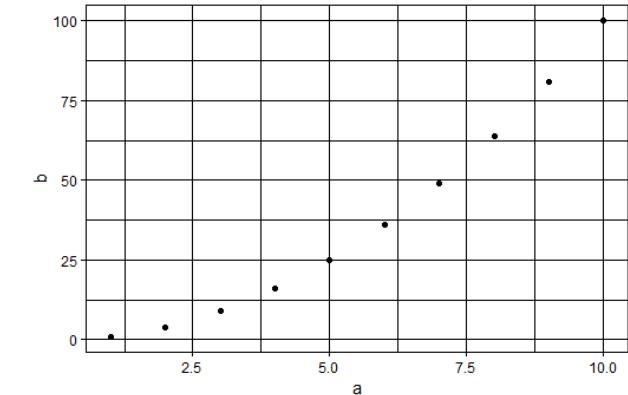
Theme classic



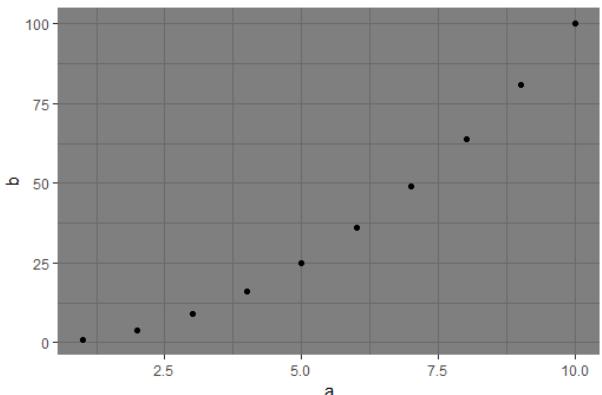
Theme light



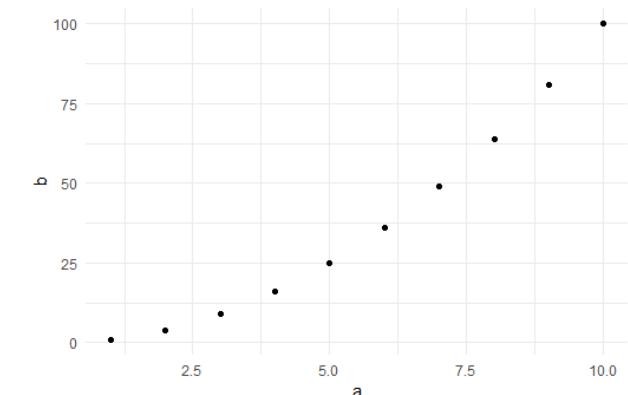
Theme linedraw



Theme dark



Theme minimal



Setting and Customising themes

- Globally

```
theme_set(theme_bw(base_size=14))
```

```
theme_update(plot.title = element_text(hjust = 0.5))
```

- In a single plot

```
+theme_dark()
```

```
+theme(plot.title = element_text(hjust=0.5))
```

What can you customise?

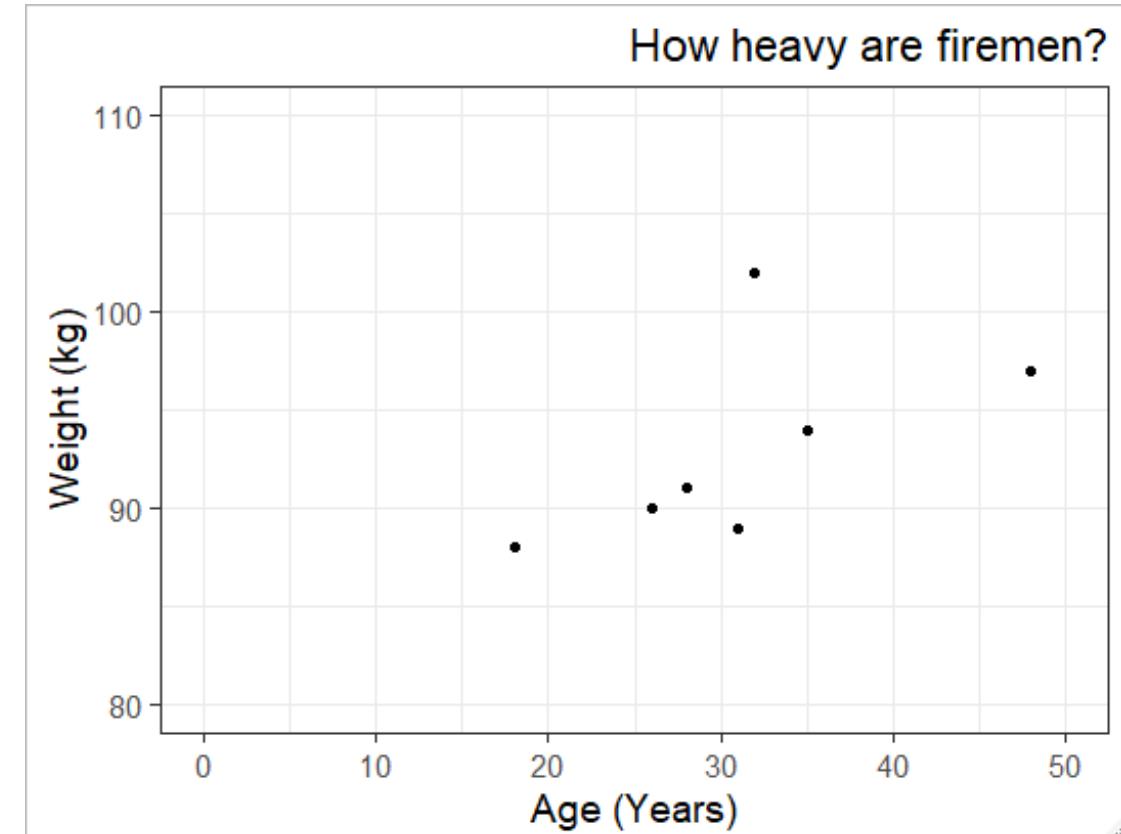
```
theme(line, rect, text, title, aspect.ratio, axis.title, axis.title.x,  
axis.title.x.top, axis.title.x.bottom, axis.title.y, axis.title.y.left,  
axis.title.y.right, axis.text, axis.text.x, axis.text.x.top,  
axis.text.x.bottom, axis.text.y, axis.text.y.left, axis.text.y.right,  
axis.ticks, axis.ticks.x, axis.ticks.x.top, axis.ticks.x.bottom,  
axis.ticks.y, axis.ticks.y.left, axis.ticks.y.right, axis.ticks.length,  
axis.line, axis.line.x, axis.line.x.top, axis.line.x.bottom, axis.line.y,  
axis.line.y.left, axis.line.y.right, legend.background, legend.margin,  
legend.spacing, legend.spacing.x, legend.spacing.y, legend.key,  
legend.key.size, legend.key.height, legend.key.width, legend.text,  
legend.text.align, legend.title, legend.title.align, legend.position,  
legend.direction, legend.justification, legend.box, legend.box.just,  
legend.box.margin, legend.box.background, legend.box.spacing,  
panel.background, panel.border, panel.spacing, panel.spacing.x,  
panel.spacing.y, panel.grid, panel.grid.major, panel.grid.minor,  
panel.grid.major.x, panel.grid.major.y, panel.grid.minor.x,  
panel.grid.minor.y, panel.on top, plot.background, plot.title,  
plot.subtitle, plot.caption, plot.tag, plot.tag.position, plot.margin,  
strip.background, strip.background.x, strip.background.y,  
strip.placement, strip.text, strip.text.x, strip.text.y,  
strip.switch.pad.grid, strip.switch.pad.wrap
```

Theme setting example

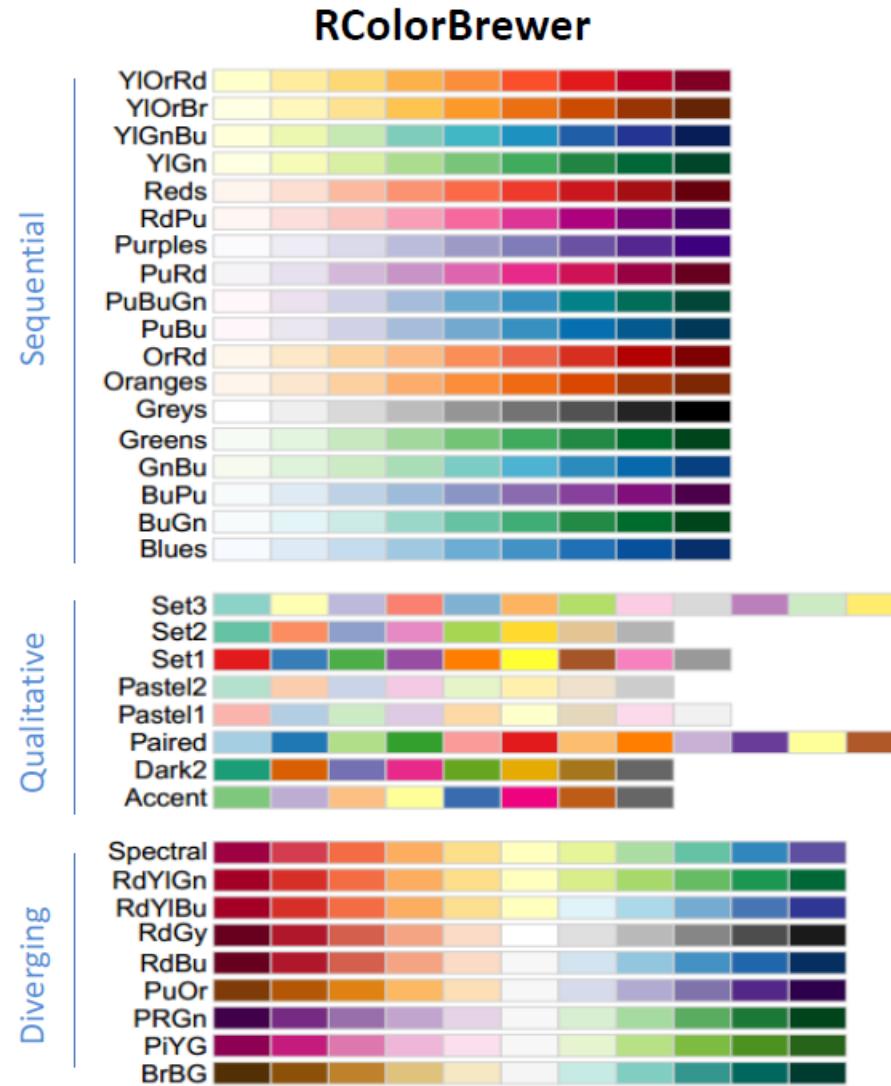
```
theme_set(theme_bw(base_size = 14))  
theme_update(plot.title = element_text(hjust=1))
```

OR

```
my.plot +  
  theme_bw(base_size = 14) +  
  theme(plot.title = element_text(hjust=1))
```



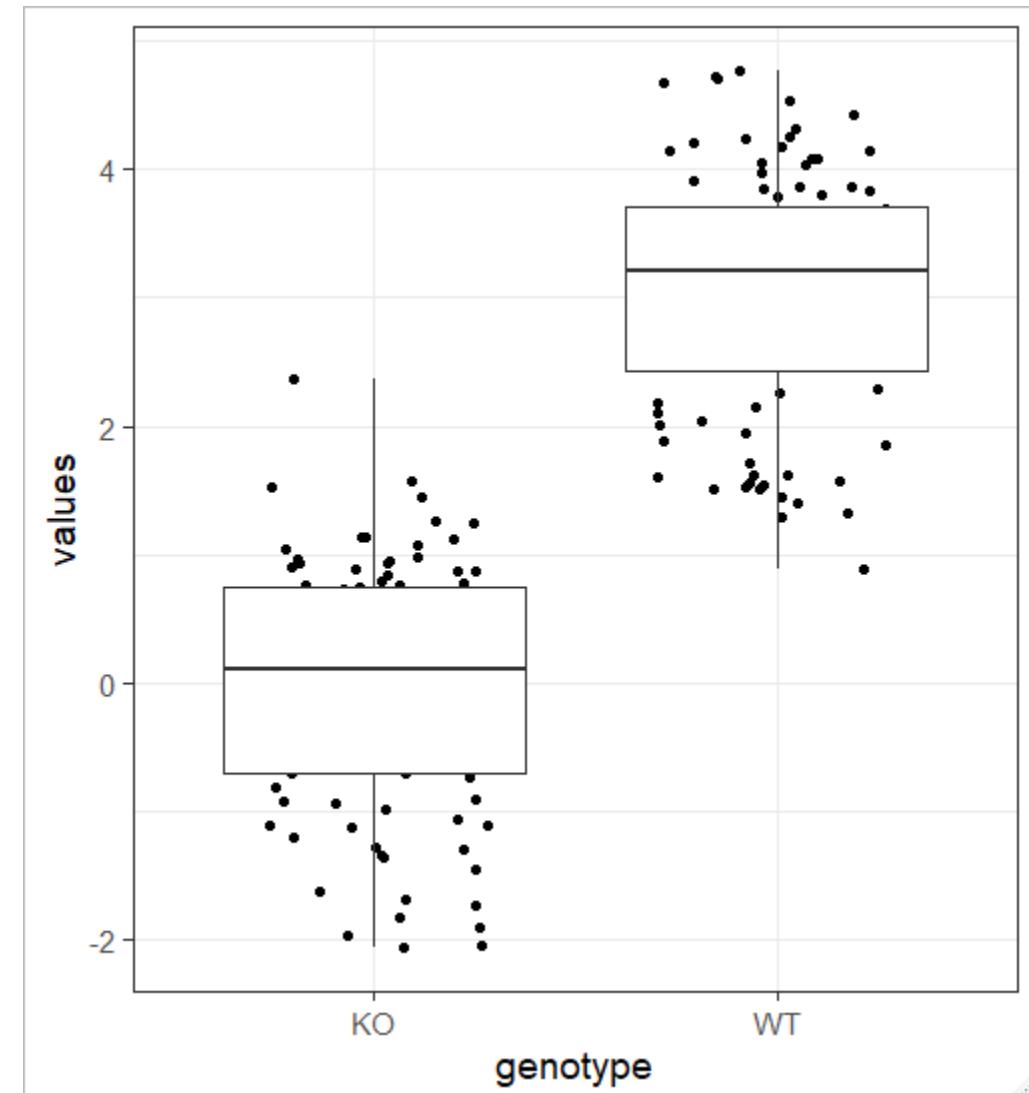
ColorBrewer Scales



`scale_color_brewer` for qualitative
`scale_color_distiller` for quantitative

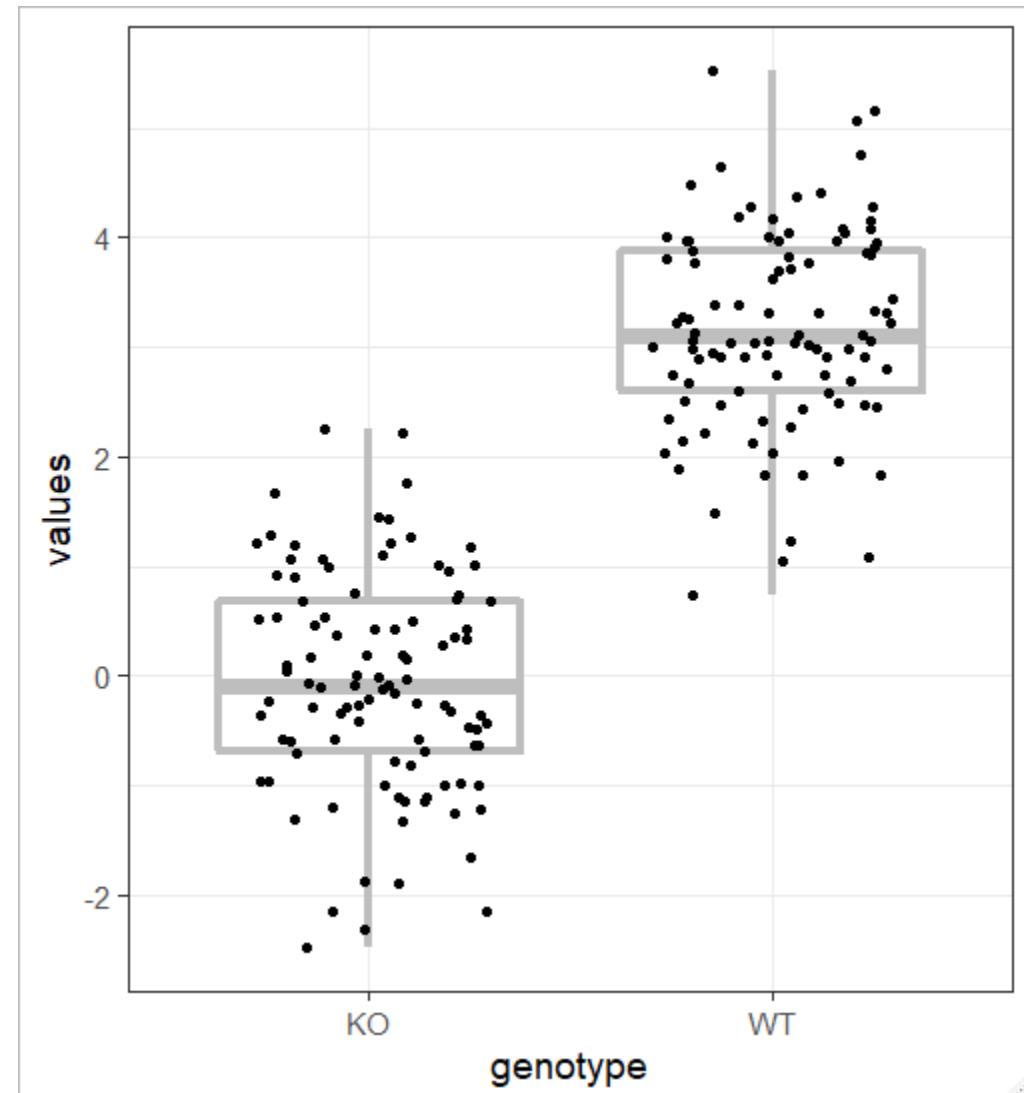
Overlaying raw data and summaries

```
many.values %>%
  group_by(genotype) %>%
  sample_n(100) %>%
  ggplot(aes(x=genotype, y=values)) +
  geom_jitter(height=0, width = 0.3) +
  geom_boxplot()
```



Overlaying raw data and summaries

```
many.values %>%
  group_by(genotype) %>%
  sample_n(100) %>%
  ggplot(aes(x=genotype, y=values)) +
  geom_boxplot(size=1.5, colour="grey") +
  geom_jitter(height=0, width = 0.3)
```



Saving plots

ggsave() saves the last plot.

Uses size on screen:

```
ggsave("my-plot.pdf")  
ggsave("my-plot.png")
```

Specify size in inches

```
ggsave("my-plot.pdf", width = 6, height = 6)
```



To make a graph

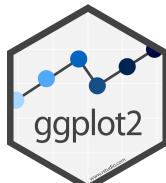


1. Pick a **data** set

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

2. Choose a **geom**
to display cases

3. **Map** aesthetic
properties to
variables

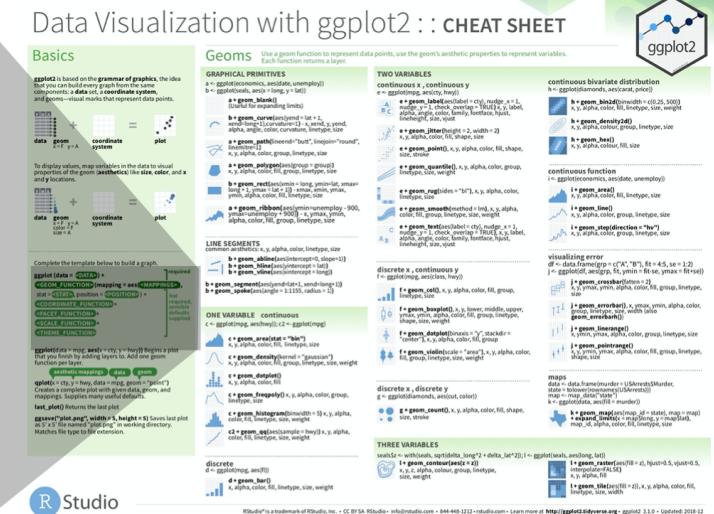


A ggplot2 template

Make any plot by filling in the parameters of this template

```
ggplot (data = <DATA>) +  
<GEOM_FUNCTION>(mapping = aes(<MAPPINGS>),  
  stat = <STAT>, position = <POSITION>) +  
<COORDINATE_FUNCTION> +  
<FACET_FUNCTION> +  
<SCALE_FUNCTION> +  
<THEME_FUNCTION>
```

↑ required
↓ Not required, sensible defaults supplied



Overview

ggplot2 is a system for declaratively creating graphics, based on [The Grammar of Graphics](#). You provide the data, tell ggplot2 how to map variables to aesthetics, what graphical primitives to use, and it takes care of the details.

Installation

```
# The easiest way to get ggplot2 is to install the whole tidyverse:  
install.packages("tidyverse")  
  
# Alternatively, install just ggplot2:  
install.packages("ggplot2")  
  
# Or the development version from GitHub:  
# install.packages("devtools")  
devtools::install_github("tidyverse/ggplot2")
```

Cheatsheet



Links

Download from CRAN at

<https://cloud.r-project.org/>
package=ggplot2

Browse source code at

<https://github.com/tidyverse/ggplot2/>

Report a bug at

<https://github.com/tidyverse/ggplot2/>
issues

Learn more at

<https://r4ds.had.co.nz/data-visualisation.html>

Extensions at

[https://exts.ggplot2.tidyverse.org/
gallery/](https://exts.ggplot2.tidyverse.org/gallery/)

License

[Full license](#)

[GPL-2](#) | file [LICENSE](#)

Community

[Contributing guide](#)

[Code of conduct](#)

Citation

Books

- Elements of Data Analytic Style, Jeff Leek
- R Programming for Data Science, Roger Peng
- The Art of Data Science, Roger Peng
- R Cookbook. Both a website, and a book, Winston Chang
- R for Data Science. Both a website and a book. Hadley Wickham and Garrett Grolemund.

