

MTHM501:Working with Data

Mark Kelson, m.j.kelson@exeter.ac.uk



OVERVIEW

- ▶ Online
 - ▶ 2 practicals per week - each 2 hours
 - ▶ Self-directed learning outside of that
 - ▶ 10 hours of lectures
 - ▶ 20 hours of practicals
 - ▶ 36 hours of background reading
 - ▶ 84 hours guided independent study - assessments
 - ▶ Student hour: Wednesday 9-10 (changed since the video!)- Book in advance
 - ▶ Assessment:
 - ▶ 4 formative courseworks
 - ▶ 1 project (10 pages)
 - ▶ 100% weighting for the project

AIMS

- ▶ to bring everyone up to speed on the fundamentals of analysis and statistical thinking
 - ▶ to establish core competencies in RStudio
 - ▶ to instil confidence in your ability to handle data
 - ▶ to encourage rigour in your analytical approaches
 - ▶ to exercise your self-directed learning skills

TOPICS

- ▶ Data- types, wrangling, coding
 - ▶ RStudio and Markdown
 - ▶ Visualisation
 - ▶ Maps
 - ▶ Missing data
 - ▶ Uncertainty in data
 - ▶ P-values, communication

VISUALISATION

- ▶ Data visualisation is the presentation of data in a graphical format.
 - ▶ It can provide a valuable insight into your data and help in identifying patterns.
 - ▶ Numerous methods are available to visualise your data
 - ▶ bar charts
 - ▶ pie charts
 - ▶ scatterplots
 - ▶ histograms
 - ▶ box plots
 - ▶ line plots
 - ▶ maps
 - ▶ ... and many more!

VISUALISATION

Deviation	Correlation	Ranking	Distribution	Change over Time	Magnitude	Part-to-whole	Spatial	Flow
 Dot plot Dot plot visualises deviation from a central value.	 Scatter PT Scatter plot showing income & life expectancy.	 Horizontal bar Horizontal bar chart showing the rank of countries.	 Histogram Histogram showing the distribution of data.	 Line Line chart showing the change over time.	 Bar Bar chart showing the magnitude of values.	 Treemap PT Treemap showing part-to-whole relationships.	 Choropleth Choropleth map showing the spatial distribution of data.	 Treemap PT Treemap showing flow of data.
 Dot plot Dot plot showing how something changes over time.	 Scatter PT Scatter plot showing the standard deviation of data.	 Horizontal bar Horizontal bar chart showing the ordered values of data.	 Histogram Histogram showing the distribution of data.	 Line Line chart showing the change in a series.	 Bar Bar chart showing column values.	 Treemap Treemap showing hierarchical structure.	 Choropleth Choropleth map showing the spatial distribution of data.	 Treemap Treemap showing flow of data.
 Dot plot Dot plot showing how something changes over time.	 Live Live chart showing the latest data.	 Horizontal bar Horizontal bar chart showing the ordered values of data.	 Histogram Histogram showing the distribution of data.	 Line Line chart showing the change in a series.	 Bar Bar chart showing column values.	 Treemap Treemap showing hierarchical structure.	 Choropleth Choropleth map showing the spatial distribution of data.	 Treemap Treemap showing flow of data.
 Dot plot Dot plot showing how something changes over time.	 Live Live chart showing the latest data.	 Horizontal bar Horizontal bar chart showing the ordered values of data.	 Histogram Histogram showing the distribution of data.	 Line Line chart showing the change in a series.	 Bar Bar chart showing column values.	 Treemap Treemap showing hierarchical structure.	 Choropleth Choropleth map showing the spatial distribution of data.	 Treemap Treemap showing flow of data.
 Dot plot Dot plot showing how something changes over time.	 Live Live chart showing the latest data.	 Horizontal bar Horizontal bar chart showing the ordered values of data.	 Histogram Histogram showing the distribution of data.	 Line Line chart showing the change in a series.	 Bar Bar chart showing column values.	 Treemap Treemap showing hierarchical structure.	 Choropleth Choropleth map showing the spatial distribution of data.	 Treemap Treemap showing flow of data.
 Dot plot Dot plot showing how something changes over time.	 Live Live chart showing the latest data.	 Horizontal bar Horizontal bar chart showing the ordered values of data.	 Histogram Histogram showing the distribution of data.	 Line Line chart showing the change in a series.	 Bar Bar chart showing column values.	 Treemap Treemap showing hierarchical structure.	 Choropleth Choropleth map showing the spatial distribution of data.	 Treemap Treemap showing flow of data.
 Dot plot Dot plot showing how something changes over time.	 Live Live chart showing the latest data.	 Horizontal bar Horizontal bar chart showing the ordered values of data.	 Histogram Histogram showing the distribution of data.	 Line Line chart showing the change in a series.	 Bar Bar chart showing column values.	 Treemap Treemap showing hierarchical structure.	 Choropleth Choropleth map showing the spatial distribution of data.	 Treemap Treemap showing flow of data.
 Dot plot Dot plot showing how something changes over time.	 Live Live chart showing the latest data.	 Horizontal bar Horizontal bar chart showing the ordered values of data.	 Histogram Histogram showing the distribution of data.	 Line Line chart showing the change in a series.	 Bar Bar chart showing column values.	 Treemap Treemap showing hierarchical structure.	 Choropleth Choropleth map showing the spatial distribution of data.	 Treemap Treemap showing flow of data.

Visual vocabulary

Designing with data

There are so many ways to visualise data... how do we know which one to pick? Use the categories across the top to decide which parts of the data are most important in your story, then look at the different types of chart within the category to form some initial ideas about what might work best. This is not meant to be exhaustive, there are many more ways to visualise data, creating informative and meaningful data visualisations.

© 2014 The Financial Times Ltd. All rights reserved. Reproduced by kind permission of the Marketing Academy.

ft.com/vocabulary

FT

VISUALISATION

Correlation

Show the relationship between two or more variables. Be mindful that, unless you tell them otherwise, many readers will assume the relationships you show them to be causal (e.g. one causes the other).

Example FT uses
inflation & unemployment, income & life expectancy

Scatterplot



The standard way to show the relationship between two continuous variables, each of which has its own axis.

Use + Column



A good way of showing the relationship between an amount (columns) and a rate (line).

Connected scatterplot



Usually used to show how the relationship between 2 variables has changed over time.

Bubble



Like a scatterplot, but adds additional detail by sizing the circles according to a third variable.

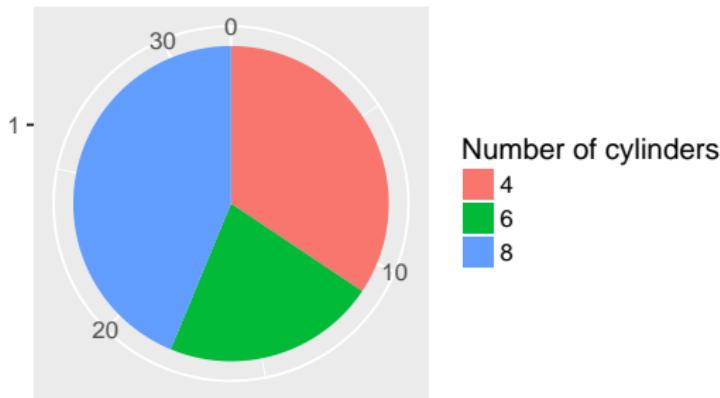
XY heatmap



A good way of showing the patterns between 2 categories of data, less good at showing fine differences in amounts.

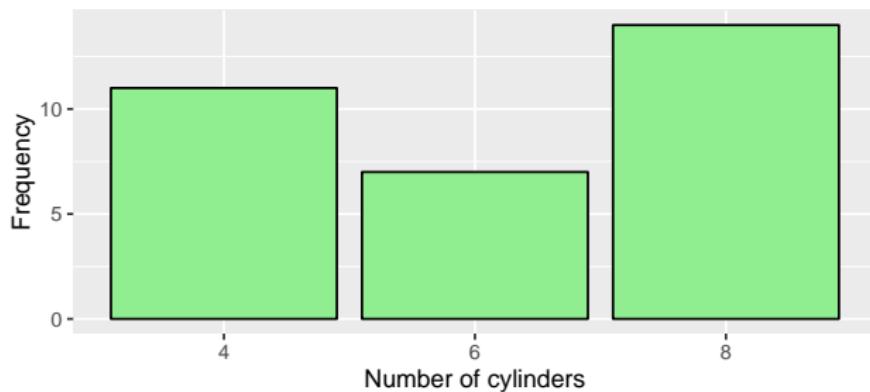
PIE CHARTS

- ▶ You can use pie charts to display data where **proportions are important**.
 - ▶ For example, the proportion of cars with 4, 6 and 8 cylinders tested in the `mtcars` dataset in R.
 - ▶ However, they can be difficult to read and interpret



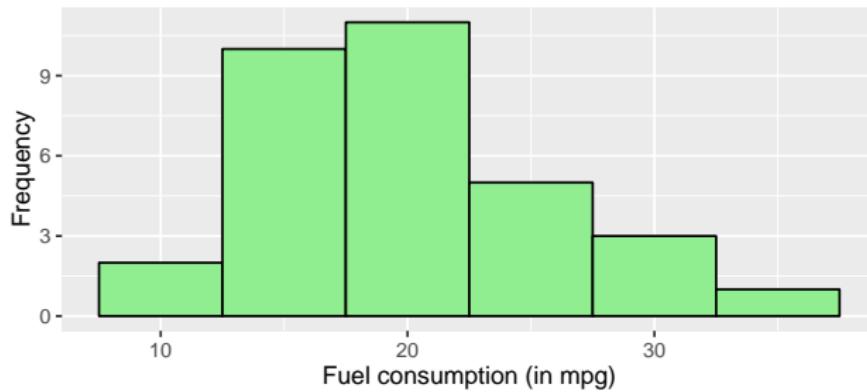
BAR CHARTS

- ▶ You can use bar charts to display **frequencies for qualitative variables**.
 - ▶ The value of a qualitative variable is represented by a bar.
 - ▶ For example, the number of cars with 4, 6 and 8 cylinders tested in the `mtcars` dataset in R.



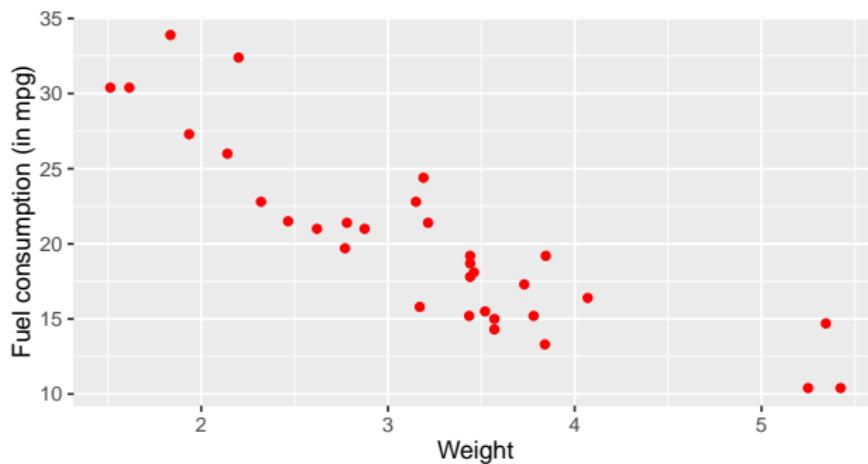
HISTOGRAMS

- ▶ You can use histograms to display the **distribution of a quantitative variable using relative frequencies**.
 - ▶ The area of each bar has a natural interpretation as a proportion of the total area of all the bars displayed
 - ▶ There is no space between the bars, and only one variable can be displayed on a single graph.
 - ▶ For example, histogram of fuel consumption (in miles per gallon) from the `mtcars` dataset in R.



SCATTER PLOTS

- ▶ You can use scatter plots to display **pairs of values of two quantitative variables**, often to check for correlation and association.
- ▶ For example, fuel consumption against weight of cars from the `mtcars` dataset in R.

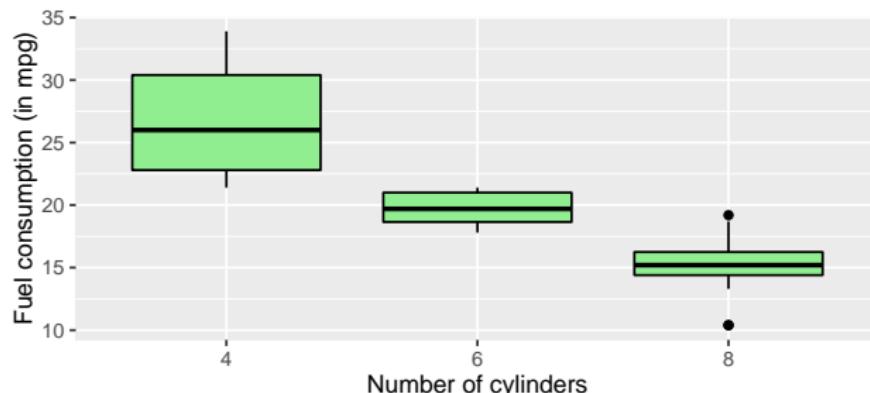


SCATTER PLOTS

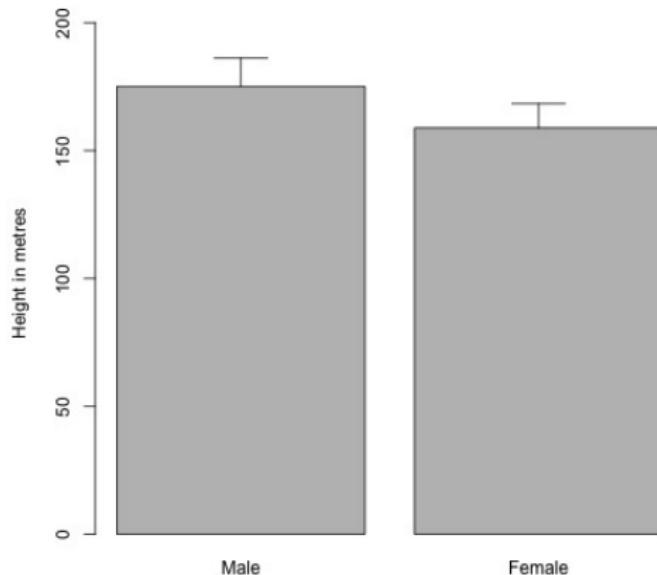
<https://twitter.com/DinaPomeranz/status/859198869061677056>

BOX PLOTS

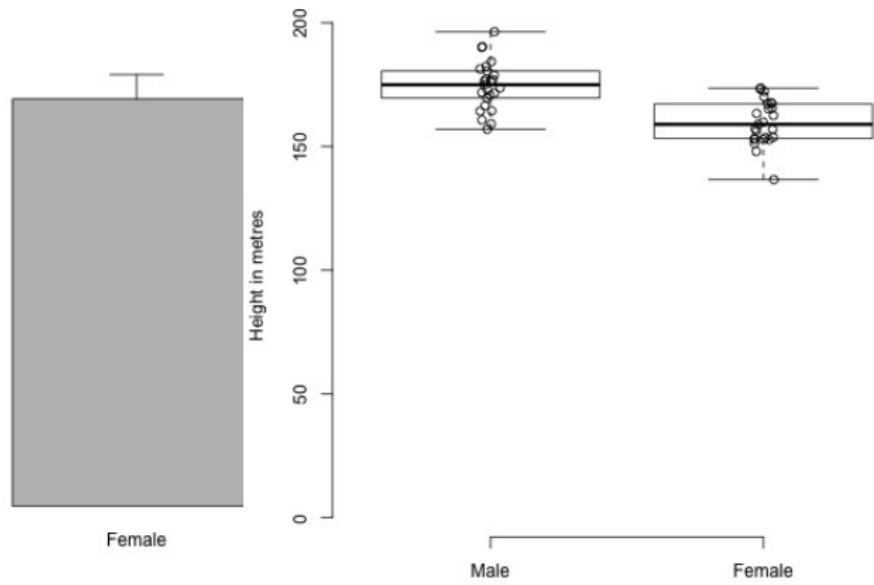
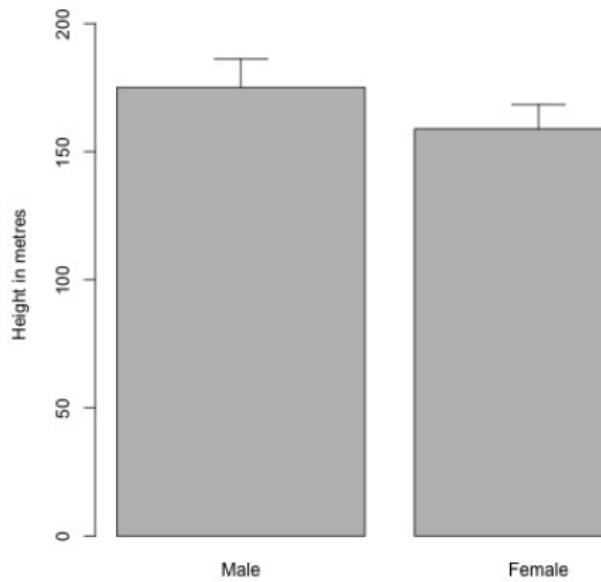
- ▶ You can use box plots to display the **median and variability between several sets of observations.**
- ▶ The central line is drawn at the median, and the box extends from the lower quartile to the upper quartile.
- ▶ For example, box plots of the fuel consumption (in miles per gallon) for cars with 4, 6 and 8 cylinders from the `mtcars` dataset in R.



DYNAMITE PLOTS

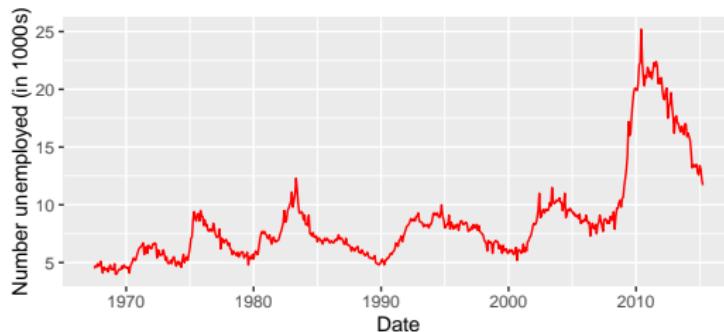


DYNAMITE PLOTS



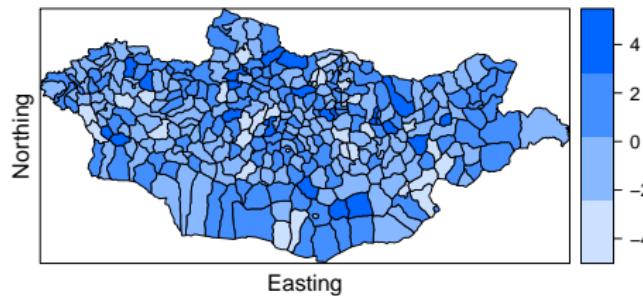
LINE PLOTS

- ▶ You can use line plots to show values of **one or more variables measured at different times**, connected by a curve.
- ▶ For example, the number of unemployed people in the US in thousands over time.



MAPS

- ▶ You can use maps to display information and variation over space.
- ▶ For example, here is a map of Mongolia.



CHOICE OF VISUALISATION

- ▶ The most appropriate type of visualisation depends on the type (qualitative/quantitative, explanatory/response) and number of variables being presented.
- ▶ Good visualisation consists of complex ideas communicated with clarity, precision, and efficiency.
- ▶ They give the viewer the greatest information in a small amount of space.

"There are no routine statistical questions, only questionable statistical routines."

D. R. Cox

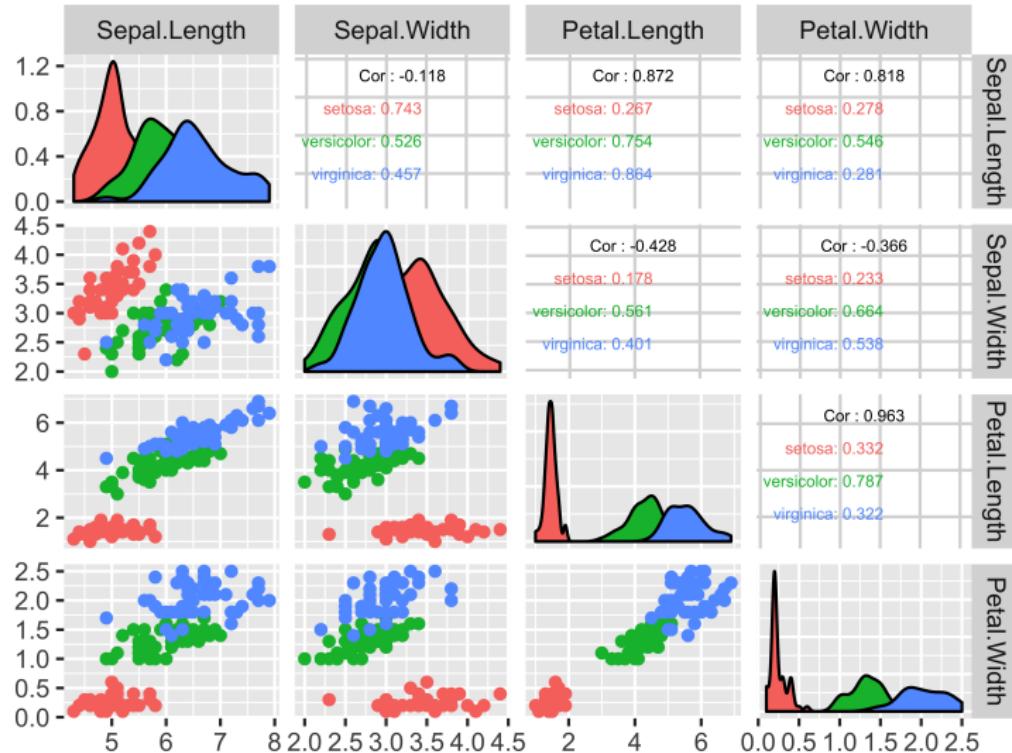
CHOICE OF VISUALISATION

- ▶ Visualisation can be done throughout an analysis
- ▶ Working
 - ▶ detect data errors and outliers
 - ▶ suggests models
 - ▶ may solve the problem alone.
- ▶ Presentation
 - ▶ effective communication (especially to non-technical audiences)
 - ▶ best and perhaps the only chance to get your message across.

CHOICE OF VISUALISATION

- For simple data sets, you can often present everything at once.

Anderson's Iris Data -- 3 species



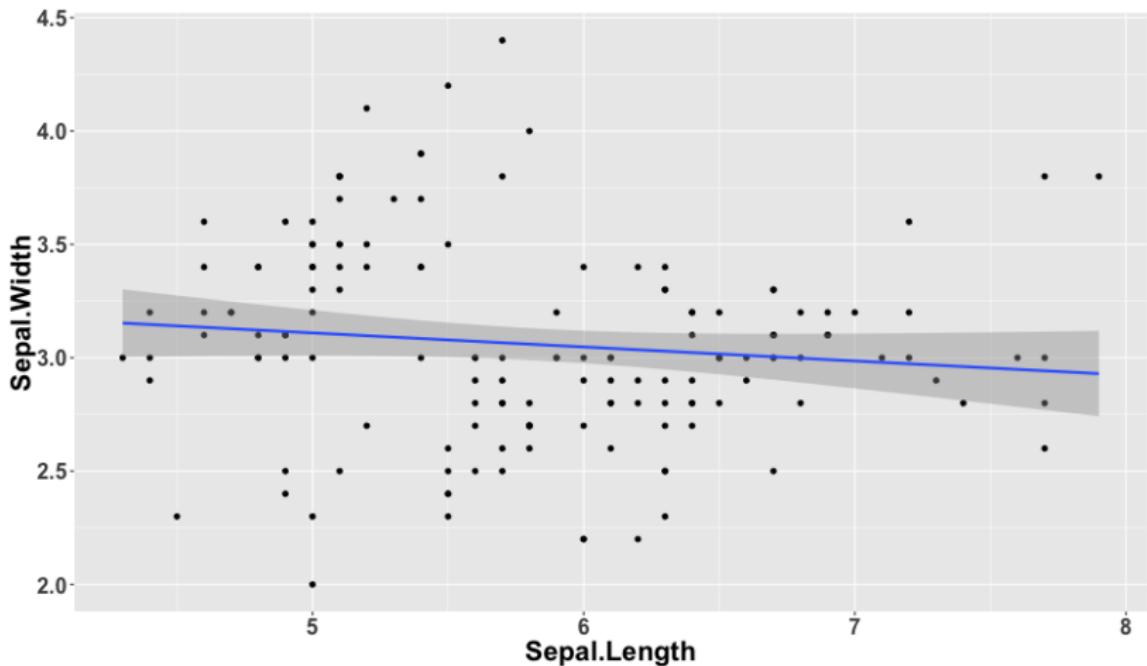
USING GRAPHICS TO PRESENT DATA

- ▶ However, it is much more difficult to view bigger datasets.
- ▶ It is important that you choose the right information to display.
- ▶ The general guidelines for visualisation are the similar to those for tables.
 - ▶ ensure that figures are self-explanatory
 - ▶ be consistent in the way that you display information
 - ▶ give clear, informative captions and titles
 - ▶ make sure your figures only contains information that adds value to your analysis and aids interpretation
 - ▶ no space is wasted.
- ▶ Always review as if you are a non-expert.

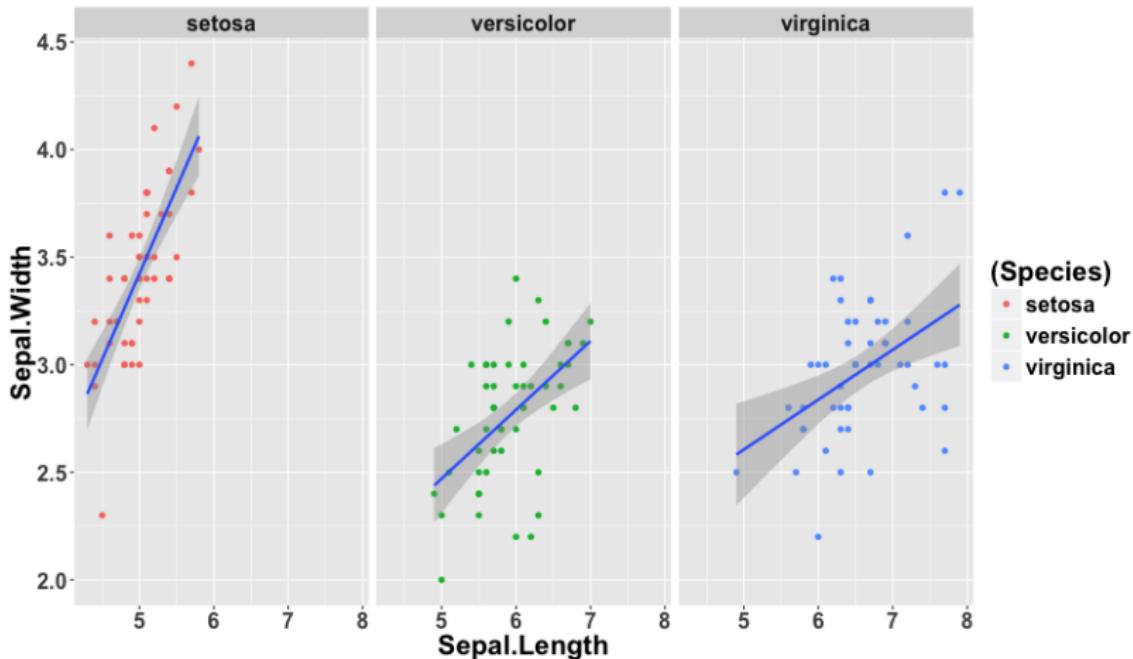
TUFTÉ'S RULES

http://www.sealthreinhold.com/school/tuftes-rules/rule_one.php

EXAMPLE: FISHER'S IRIS DATA



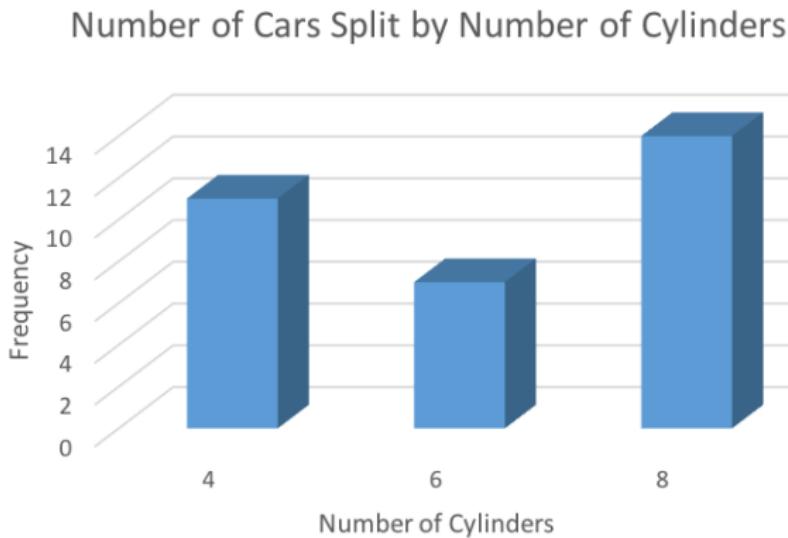
EXAMPLE: FISHER'S IRIS DATA



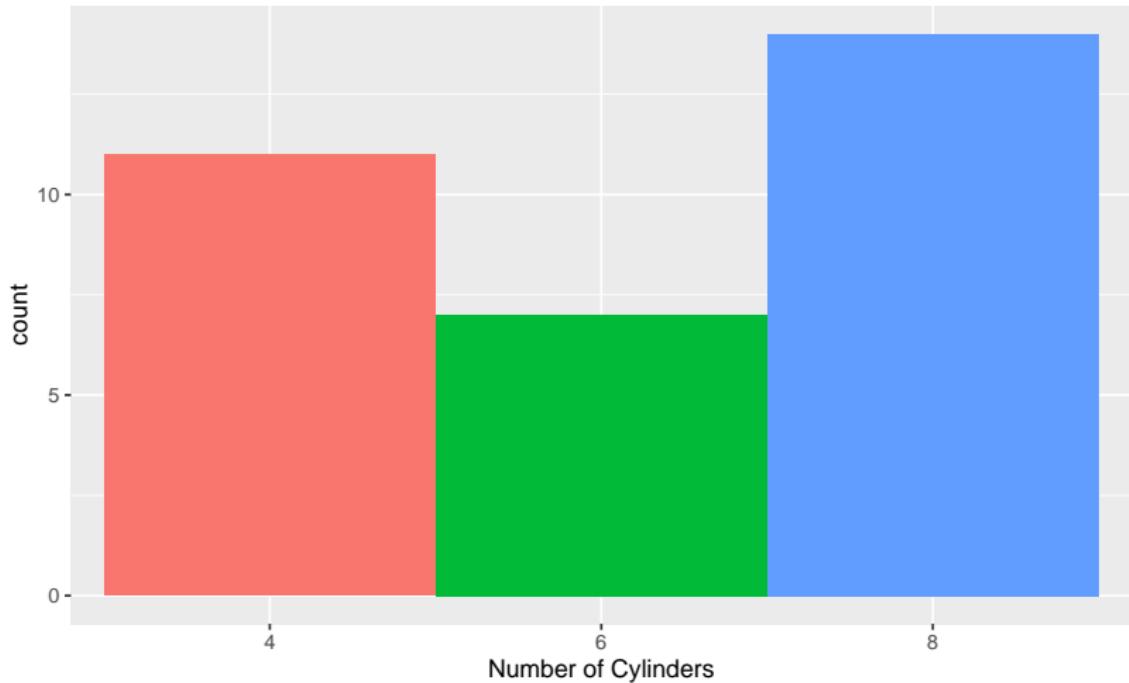
SIMPSON'S PARADOX

- ▶ Trends within groups of data can disappear or reverse when that data is aggregated
 - ▶ patterns from aggregated data do not carry over to individual-level data.
 - ▶ this means we need to explore our data very carefully to identify patterns
 - ▶ as a rule: always engage with a subject specific expert.

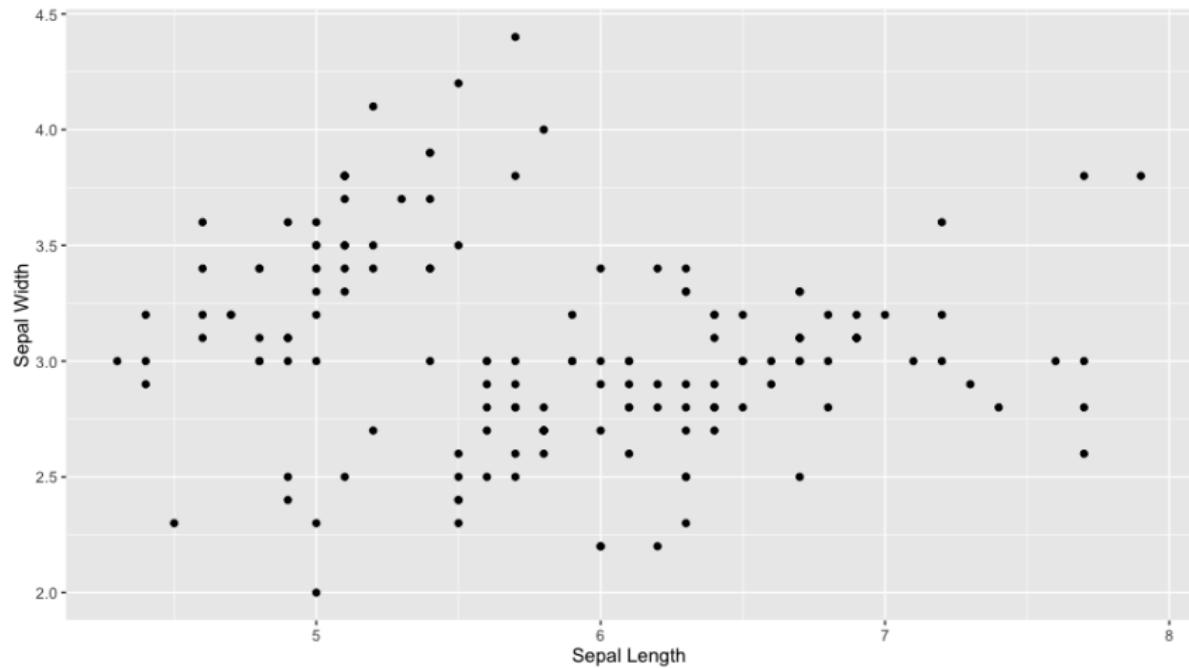
3D GRAPHICS CAN BE VERY MISLEADING



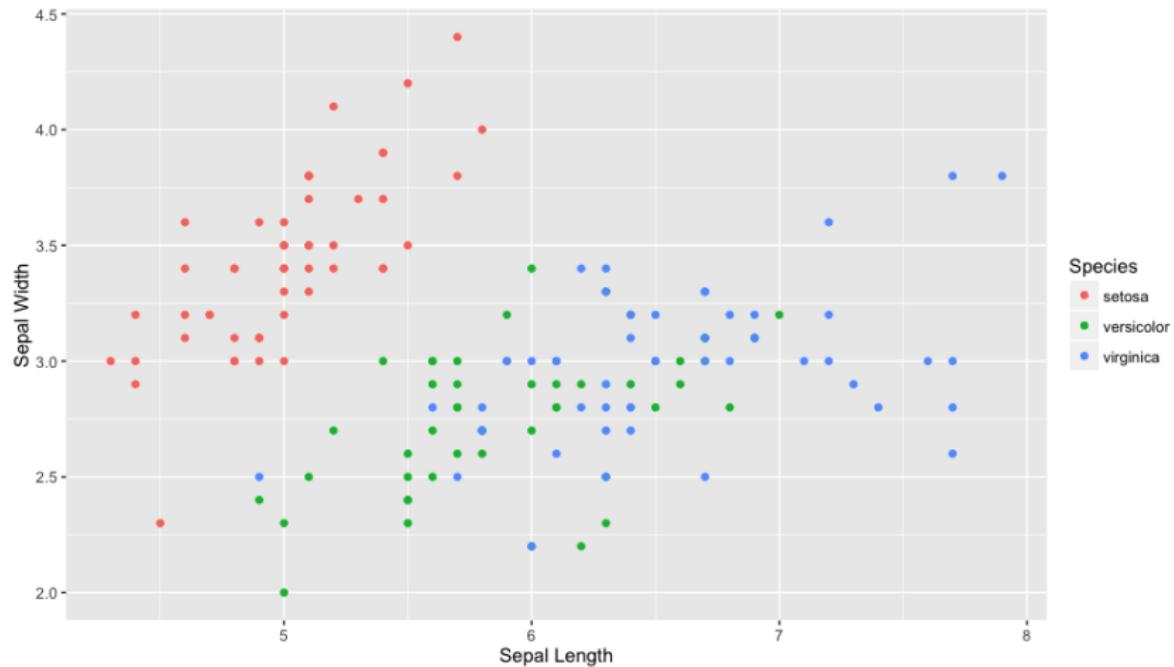
STICK TO 2D WHERE POSSIBLE



USE COLOUR



USE COLOUR

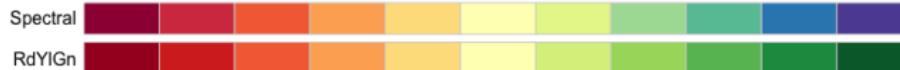


COLOUR SCHEMES

- ▶ Colour can be very helpful but there are practical issues.
- ▶ Colour scheme must be meaningful.
- ▶ Sequential colour schemes are good for ordered data, for example, population density.

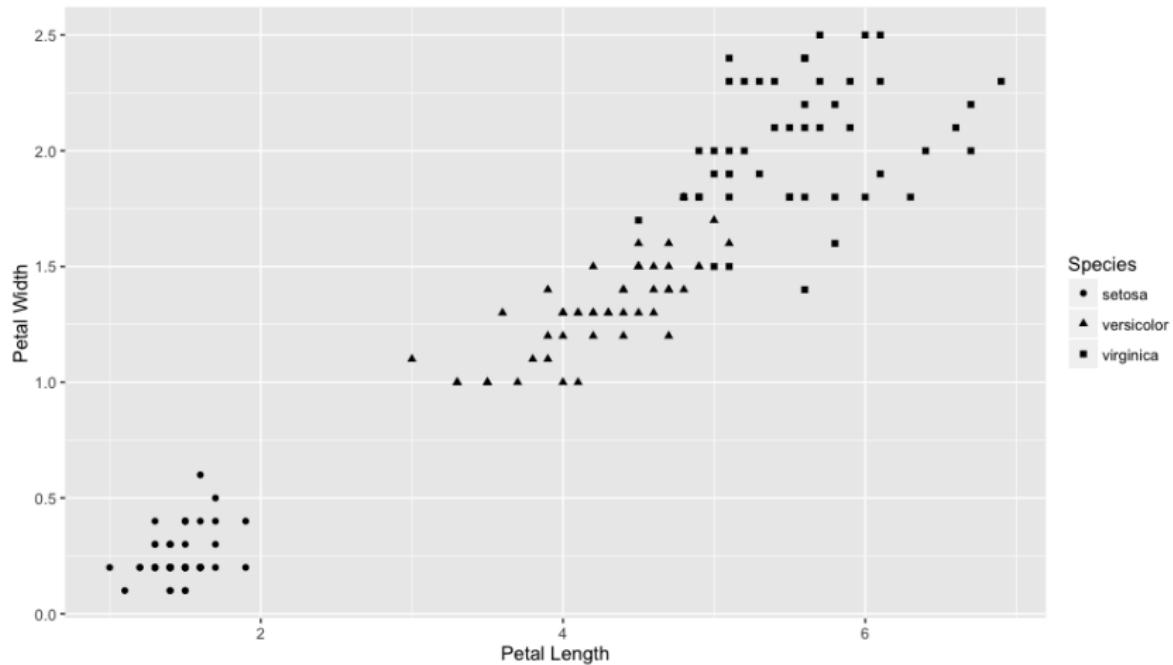


- ▶ Divergent colour schemes are good for ordered data where you want to focus on deviation from a mean level, deviance for average temperature

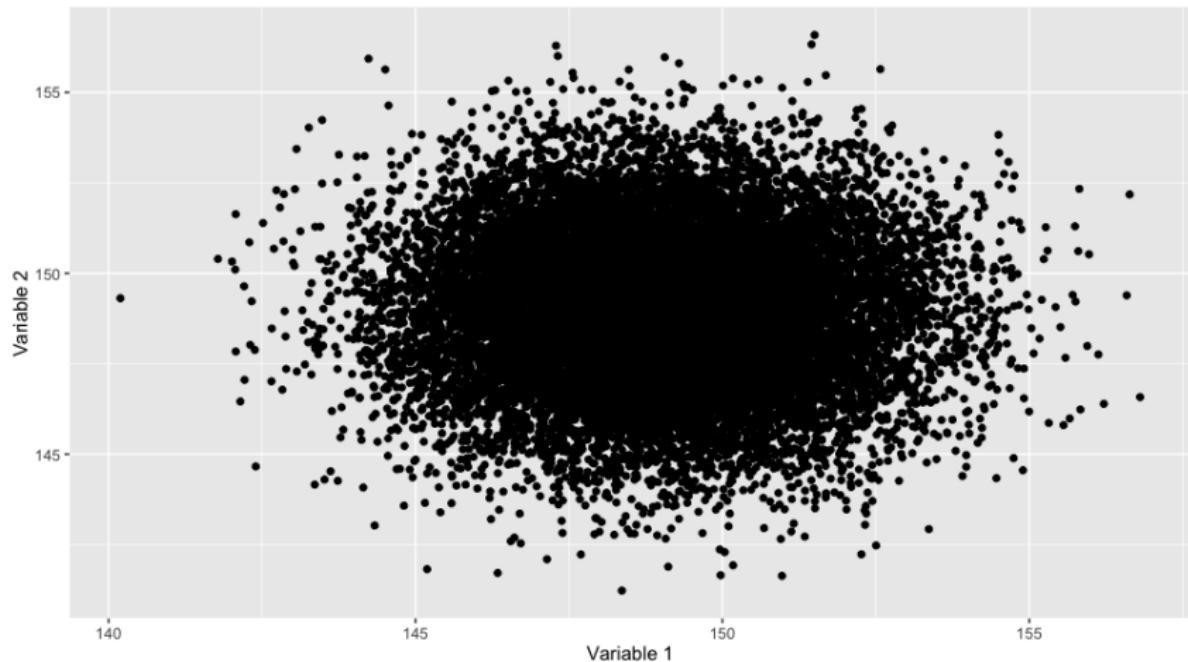


- ▶ A good choice of colour scheme are available from <http://colorbrewer2.org> and RColorBrewer R package.

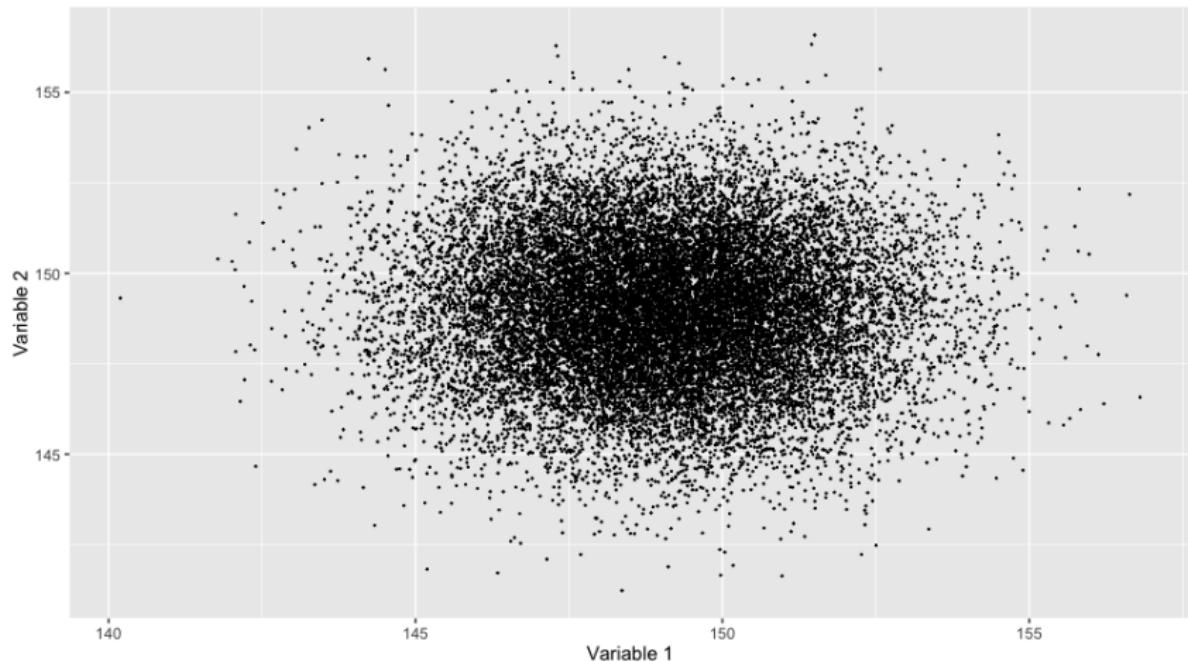
CHANGE THE STYLE



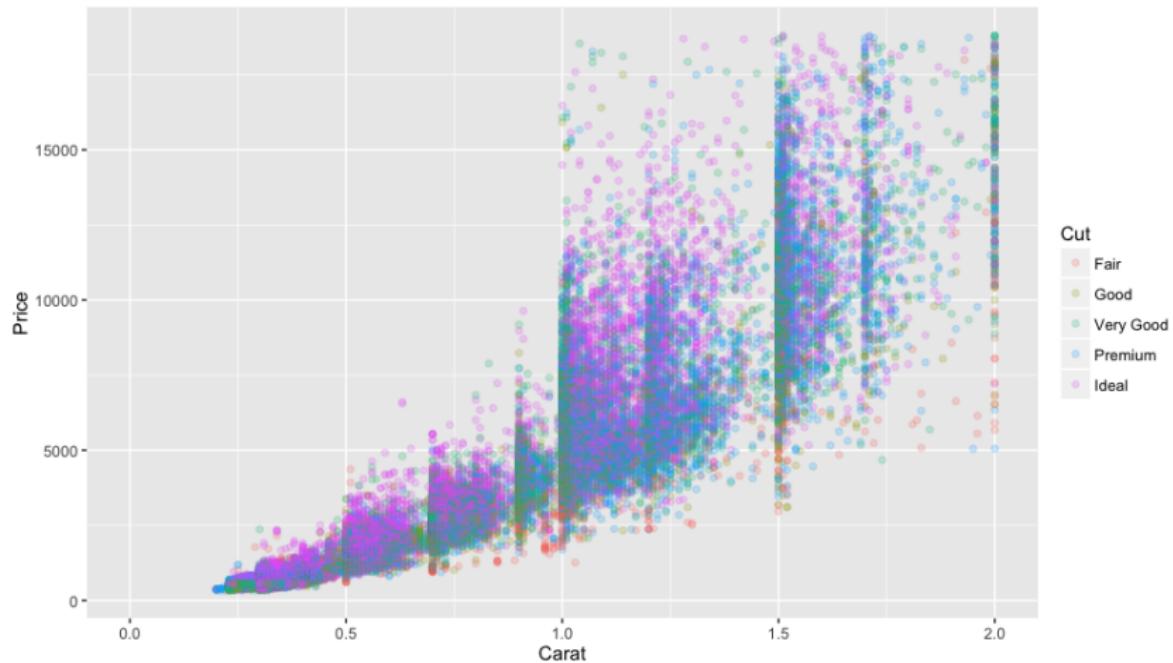
DISPLAYING BIG DATA CAN BE DIFFICULT



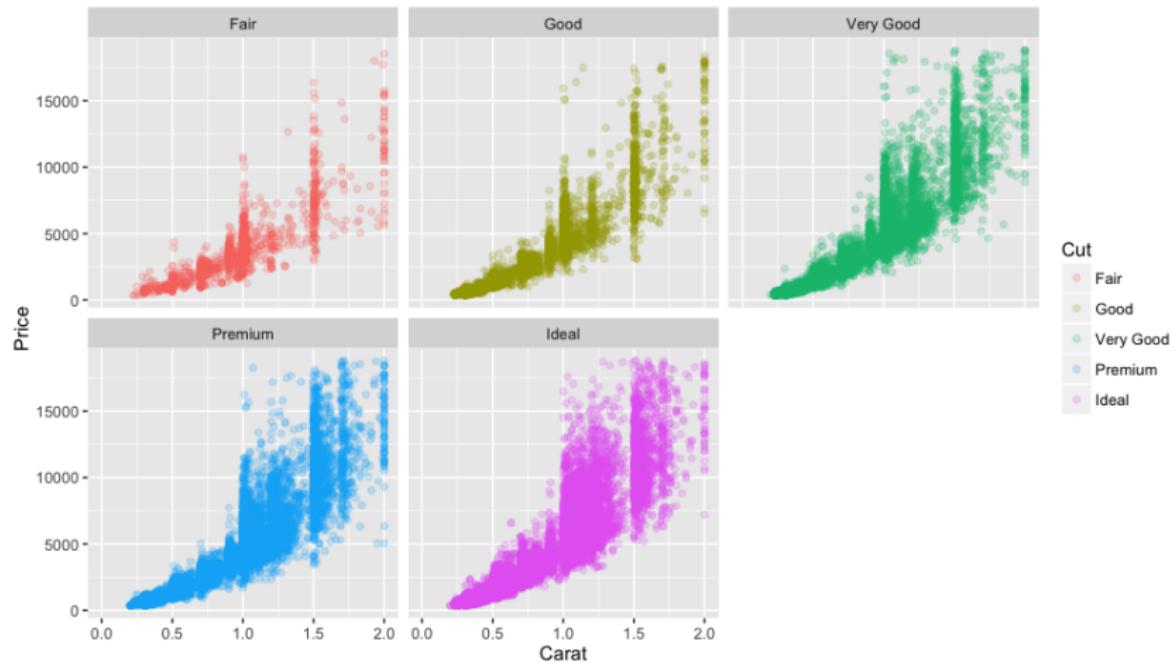
CHANGE THE SCALE



THERE ARE SOME LIMITATIONS TO THIS



USE FACETS



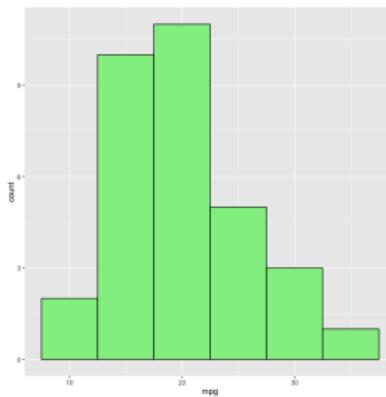
PROFESSIONAL GRAPHICS IN R

- ▶ The `ggplot2` package is a powerful graphics package in R.
 - ▶ You build a `ggplot` up piece by piece, combining the pieces with the “`+`” operator.
 - ▶ Graphics using `ggplot2` can be tailored to your analysis.

PROFESSIONAL GRAPHICS IN R

- ▶ For example we can create a histogram and store it in `p1`.

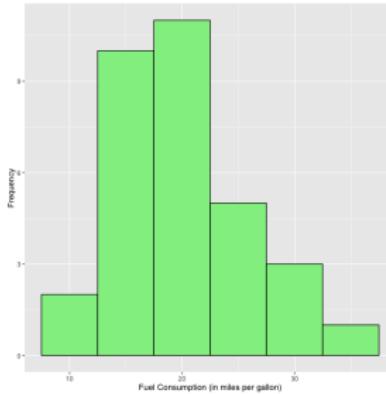
```
p1 <- ggplot(mtcars, aes(x=mpg)) + geom_histogram(binwidth=5,  
colour='black', fill='lightgreen')  
p1
```



PROFESSIONAL GRAPHICS IN R

- We can change the labels of the axes by adding to p1.

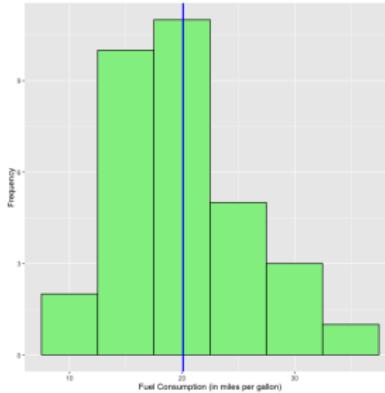
```
p1 <- p1 + labs(x = 'Fuel Consumption (in miles per gallon)',  
                 y='Frequency')  
p1
```



PROFESSIONAL GRAPHICS IN R

- We can add a line to p1 to indicate the location of the mean.

```
p1 <- p1 + vline(xintercept=mean(mtcars$mpg), col='blue', size=1)  
p1
```



OTHER RESOURCES

Data Visualization: A practical introduction, Kieran Healy

<https://socviz.co/>

Introduction to data visualisation

[https://rafalab.github.io/dsbook/
introduction-to-data-visualization.html](https://rafalab.github.io/dsbook/introduction-to-data-visualization.html)

Any Questions?