

# ETC5510: Introduction to Data Analysis

## Week 4, part B

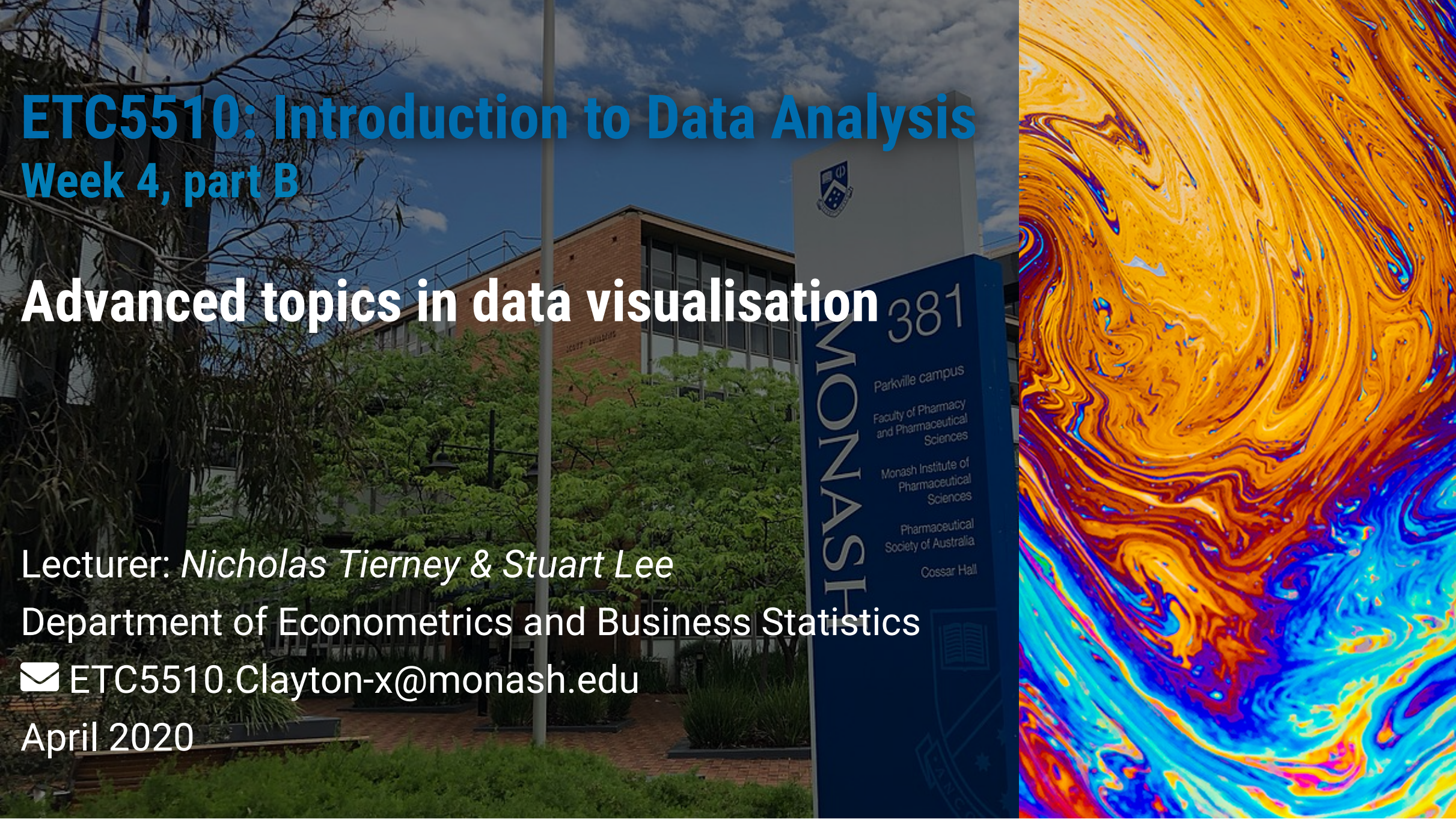
### Advanced topics in data visualisation

Lecturer: *Nicholas Tierney & Stuart Lee*

Department of Econometrics and Business Statistics

✉ [ETC5510.Clayton-x@monash.edu](mailto:ETC5510.Clayton-x@monash.edu)

April 2020





# While the song is playing...

Draw a mental model / concept map of last lectures content on joins.

# recap

- Joins
- venn diagrams
- feedback

Joins with a person and a coat, by Leight Tami

# Upcoming Due Dates

- Assignment 1: ...
- Other due dates?
- Stay tuned on ED for the upcoming dates

# Making effective data plots

1. Principles / science of data visualisation
2. Features of graphics

# Principles / science of data visualisation

- Palettes and colour blindness
- change blindness
- using proximity
- hierarchy of mappings

# Features of graphics

- Layering statistical summaries
- Themes
- adding interactivity



# Palettes and colour blindness

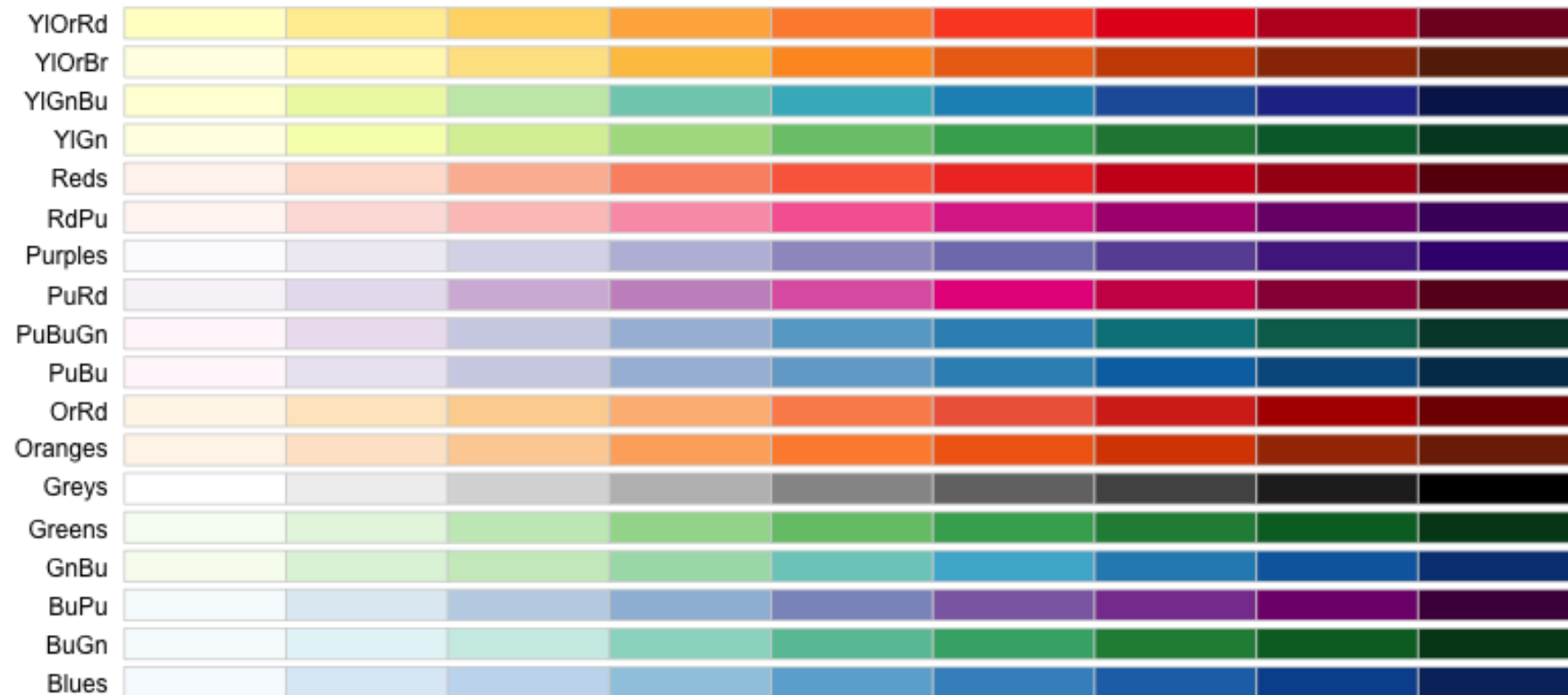
There are three main types of colour palette:

- Qualitative: categorical variables
- Sequential: low to high numeric values
- Diverging: negative to positive values

# Qualitative: categorical variables



# Sequential: low to high numeric values



# Diverging: negative to positive values



# Example: TB data

```
## # A tibble: 157,820 x 5
##   country      year count gender age
##   <chr>      <dbl> <dbl> <chr> <chr>
## 1 Afghanistan 1980     NA m     04
## 2 Afghanistan 1981     NA m     04
## 3 Afghanistan 1982     NA m     04
## 4 Afghanistan 1983     NA m     04
## 5 Afghanistan 1984     NA m     04
## 6 Afghanistan 1985     NA m     04
## 7 Afghanistan 1986     NA m     04
## 8 Afghanistan 1987     NA m     04
## 9 Afghanistan 1988     NA m     04
## 10 Afghanistan 1989     NA m     04
## # ... with 157,810 more rows
```

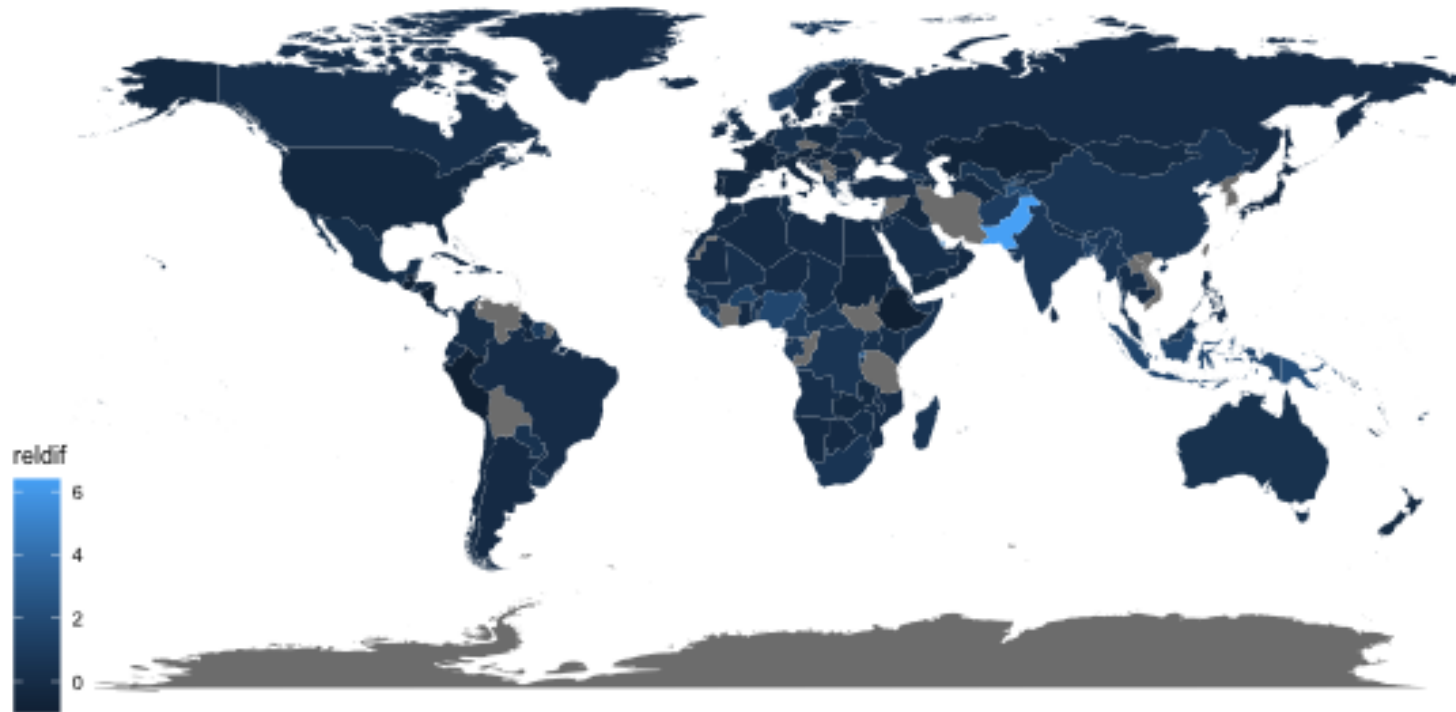
# Example: TB data: adding relative change

```
## # A tibble: 219 x 4
##   country      `2002` `2012` reldif
##   <chr>      <dbl> <dbl>   <dbl>
## 1 Afghanistan    6509  13907   1.14
## 2 Albania         225    185 -0.178
## 3 Algeria        8246   7510 -0.0893
## 4 American Samoa      1      0 -1
## 5 Andorra          2      2  0
## 6 Angola       17988  22106  0.229
## 7 Anguilla         0      0  0
## 8 Antigua and Barbuda  4      1 -0.75
## 9 Argentina       5383   4787 -0.111
## 10 Armenia        511    316 -0.382
## # ... with 209 more rows
```



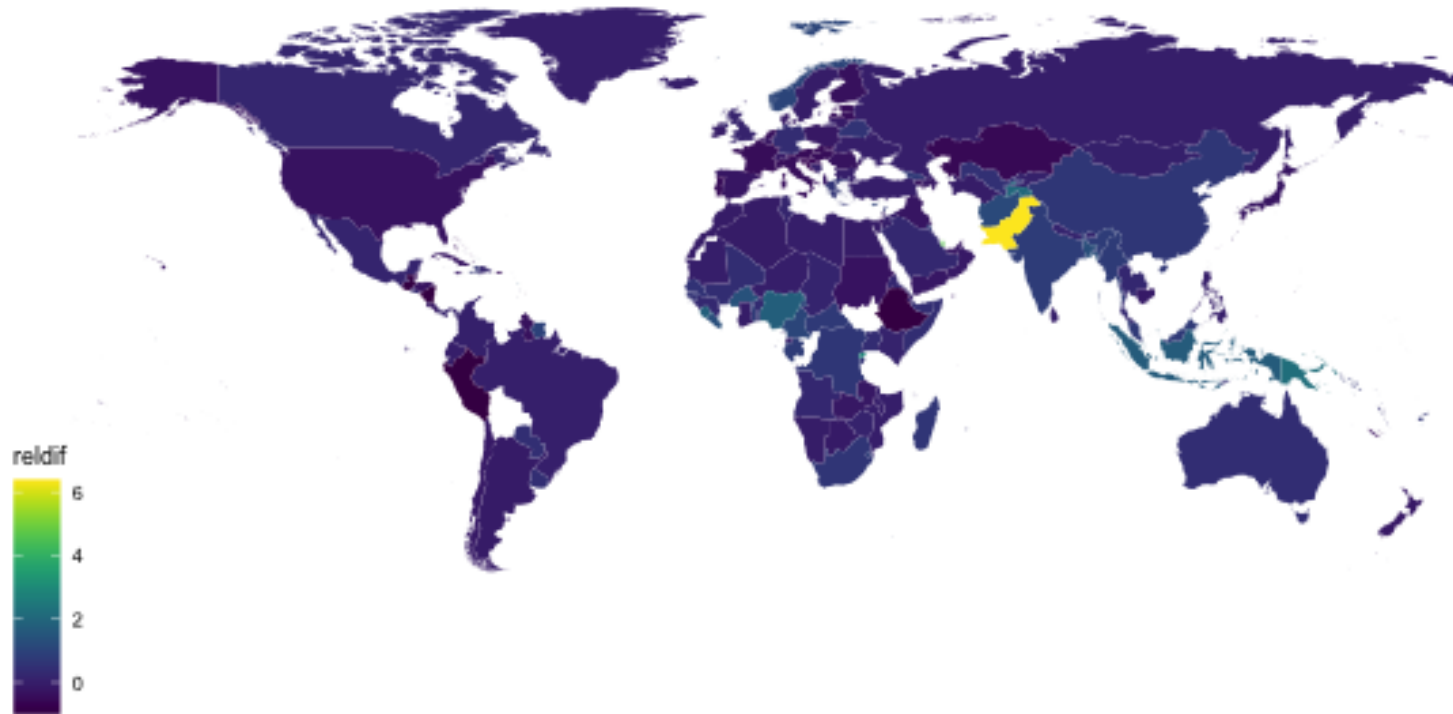
# Example: Sequential colour with default palette

```
ggplot(tb_map) + geom_polygon(aes(x = long, y = lat, group = group, fill = reldif))  
  theme_map()
```



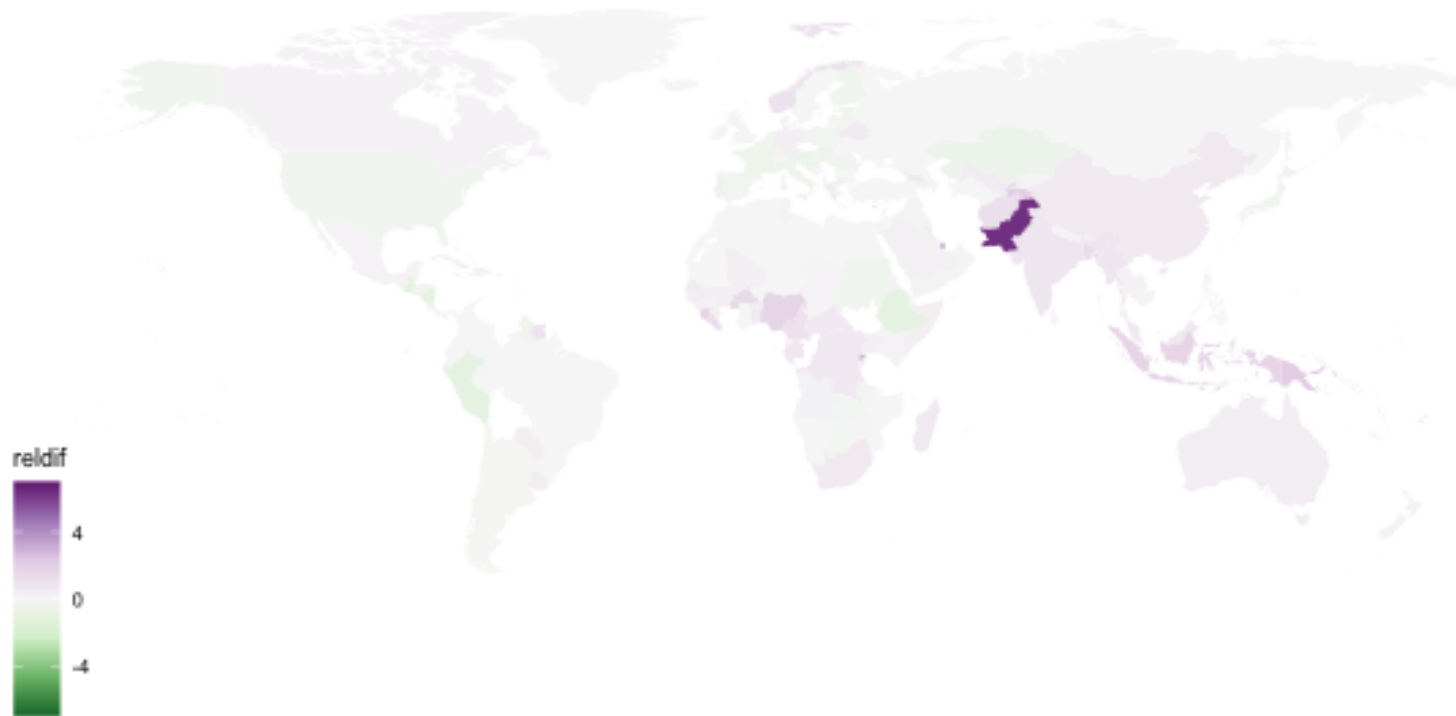
# Example: (improved) sequential colour with default palette

```
library(viridis)
ggplot(tb_map) +
  geom_polygon(aes(x = long, y = lat, group = group, fill = reldif)) +
  theme_map() + scale_fill_viridis(na.value = "white")
```



# Example: Diverging colour with better palette

```
ggplot(tb_map) +  
  geom_polygon(aes(x = long, y = lat, group = group, fill = reldif)) +  
  theme_map() +  
  scale_fill_distiller(palette = "PRGn", na.value = "white", limits = c(-7, 7))
```



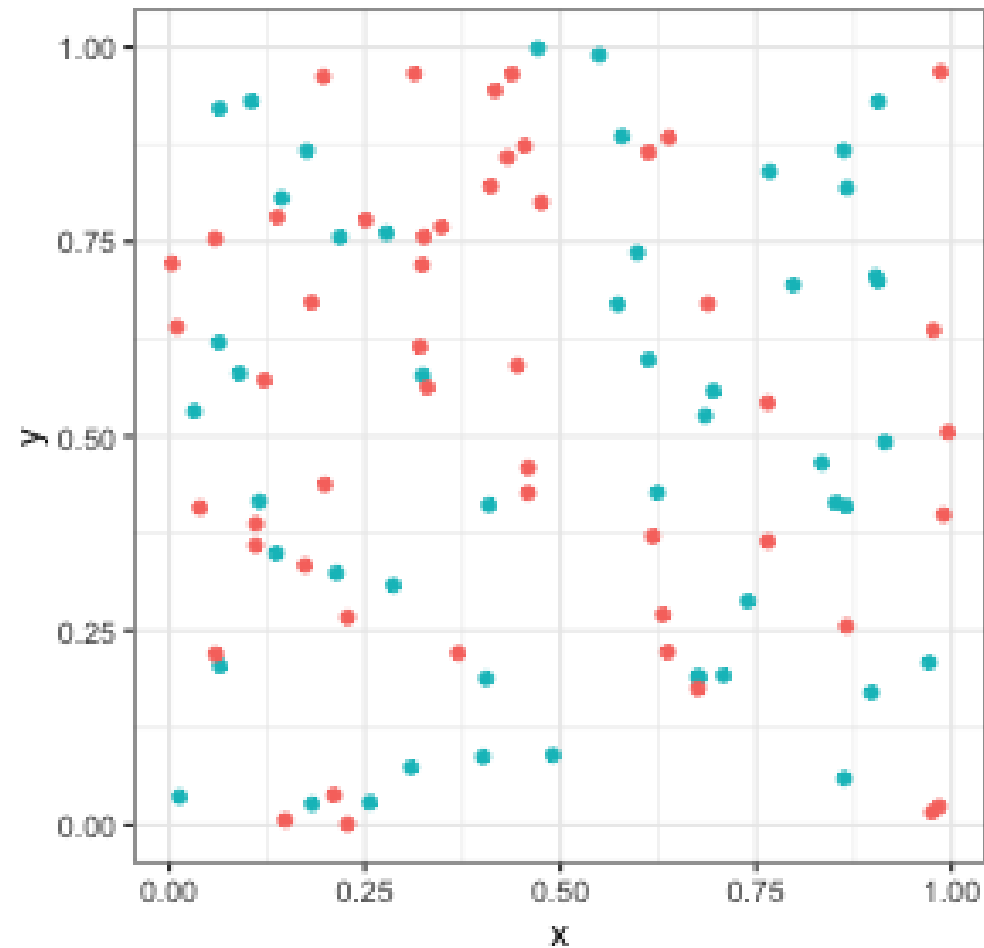
# Summary on colour palettes

- Different ways to map colour to values:
  - Qualitative: categorical variables
  - Sequential: low to high numeric values
  - Diverging: negative to positive values

# Colour blindness

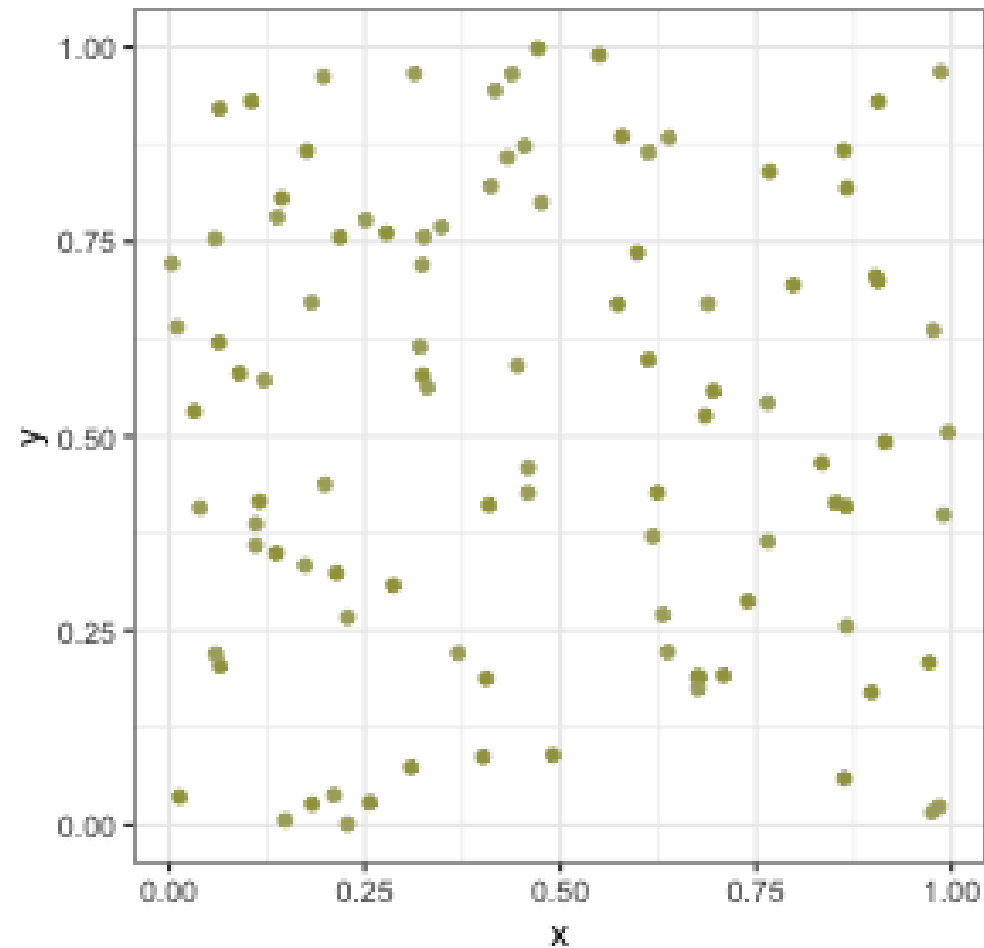
- About 8% of men (about 1 in 12), and 0.5% women (about 1 in 200) population have difficulty distinguishing between red and green.
- Several colour blind tested palettes: RColorbrewer has an associated web site [colorbrewer.org](http://colorbrewer.org) where the palettes are labelled. See also `viridis`, and `scico`.

# Plot of two coloured points: Normal Mode



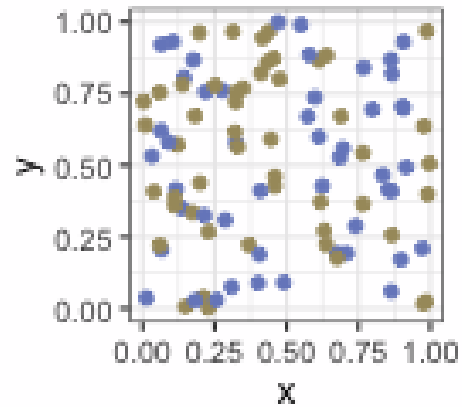


# Plot of two coloured points: dicromat mode

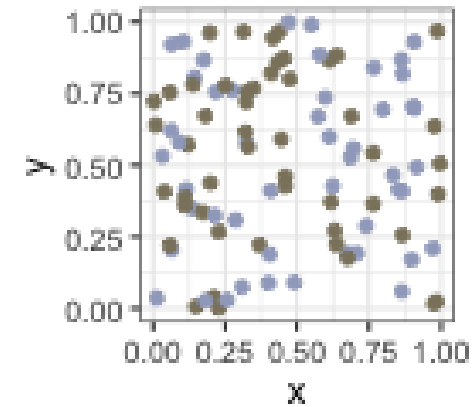


# Showing all types of colourblindness

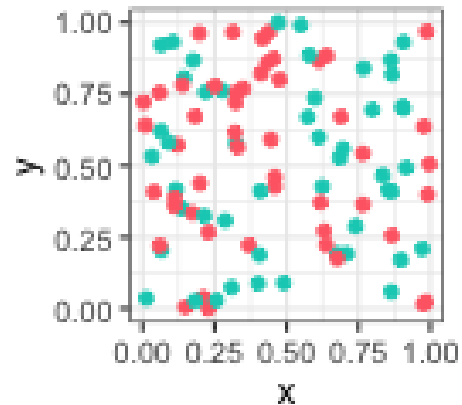
**Deutanomaly**



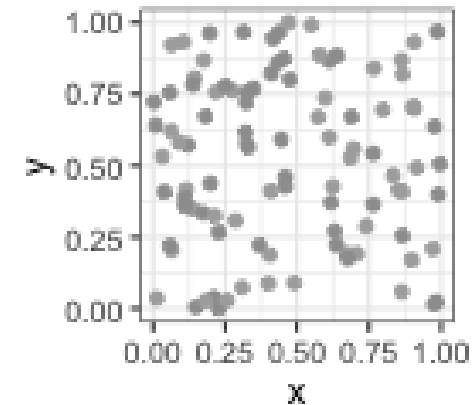
**Protanomaly**



**Tritanomaly**

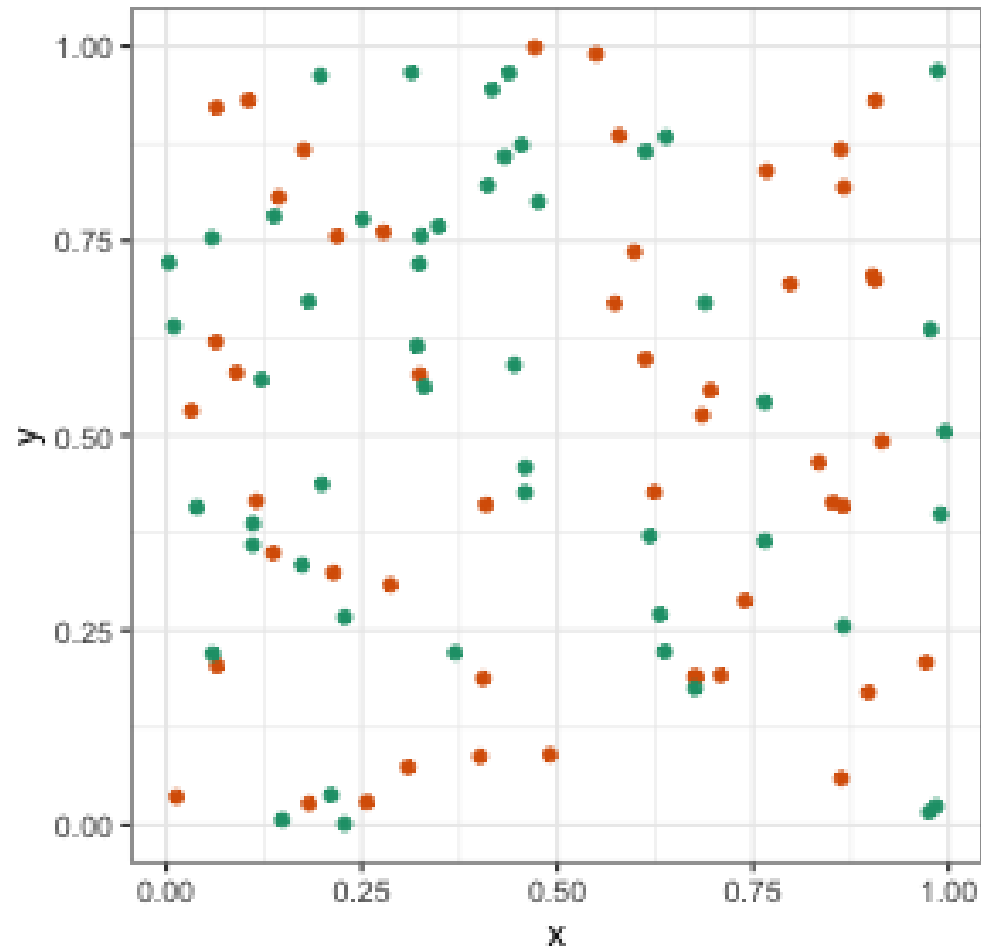


**Desaturated**



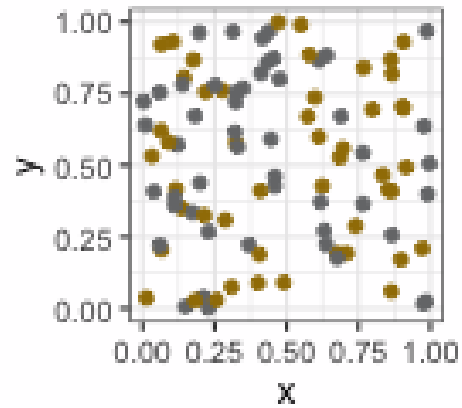
# Impact of colourblind-safe palette

```
p2 <- p + scale_colour_brewer(palette = "Dark2")  
p2
```

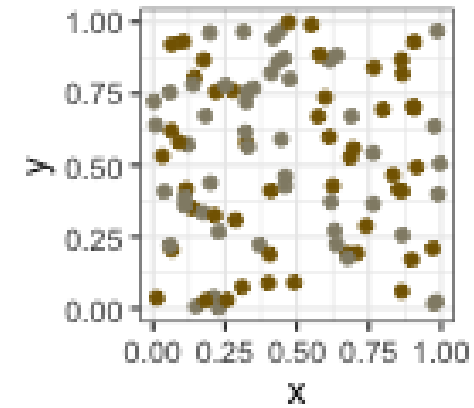


# Impact of colourblind-safe palette

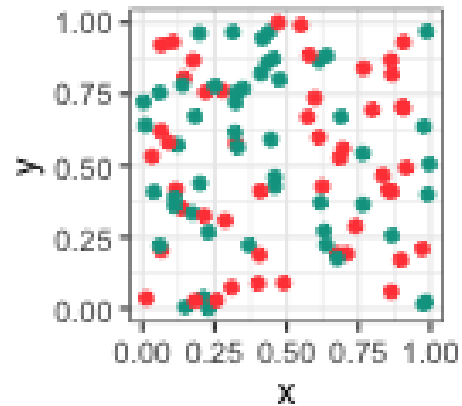
Deutanomaly



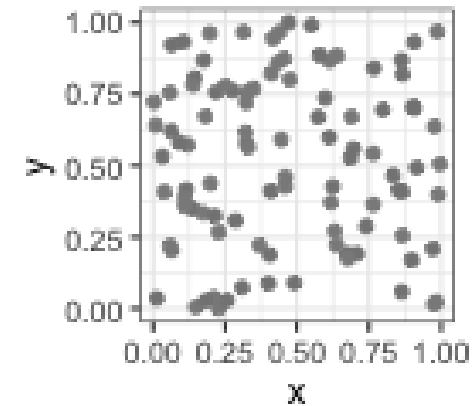
Protanomaly



Tritanomaly

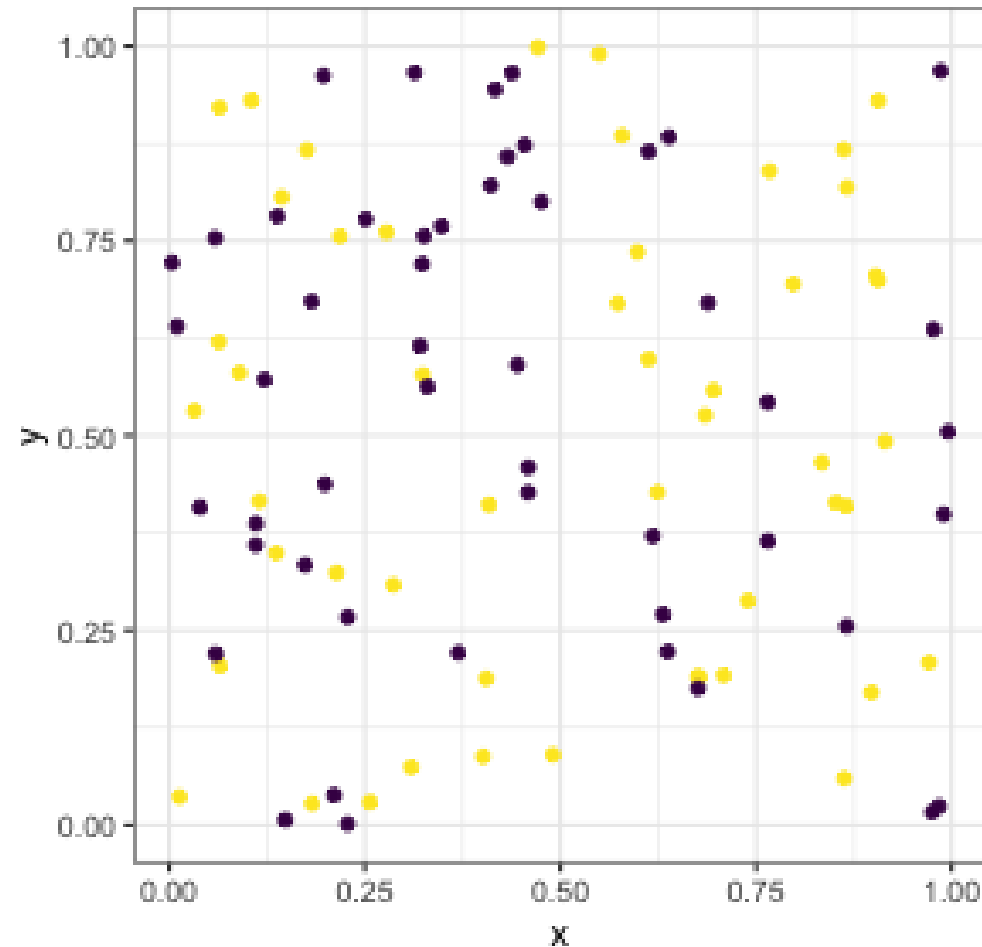


Desaturated



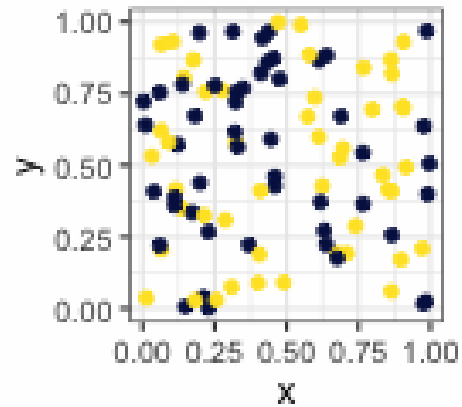
# Impact of colourblind-safe palette

```
p3 <- p + scale_colour_viridis_d()  
p3
```

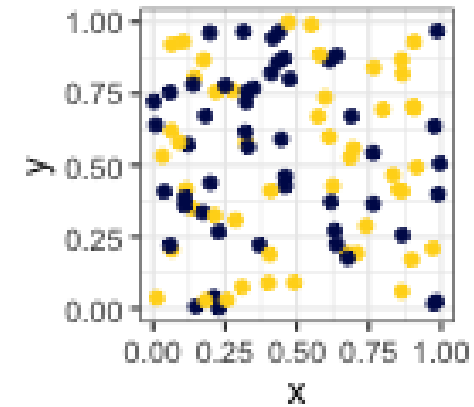


# Impact of colourblind-safe palette

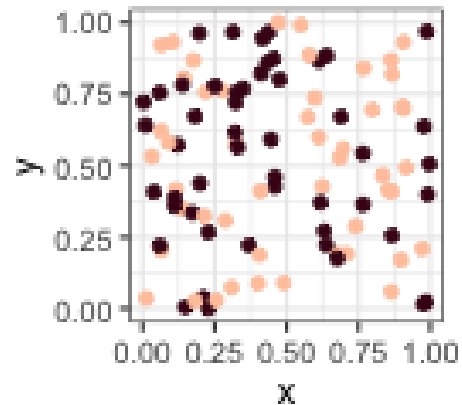
Deutanomaly



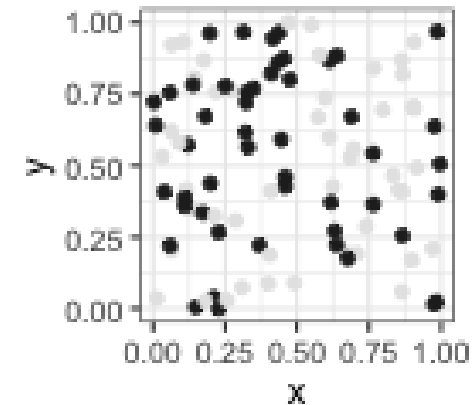
Protanomaly



Tritanomaly



Desaturated

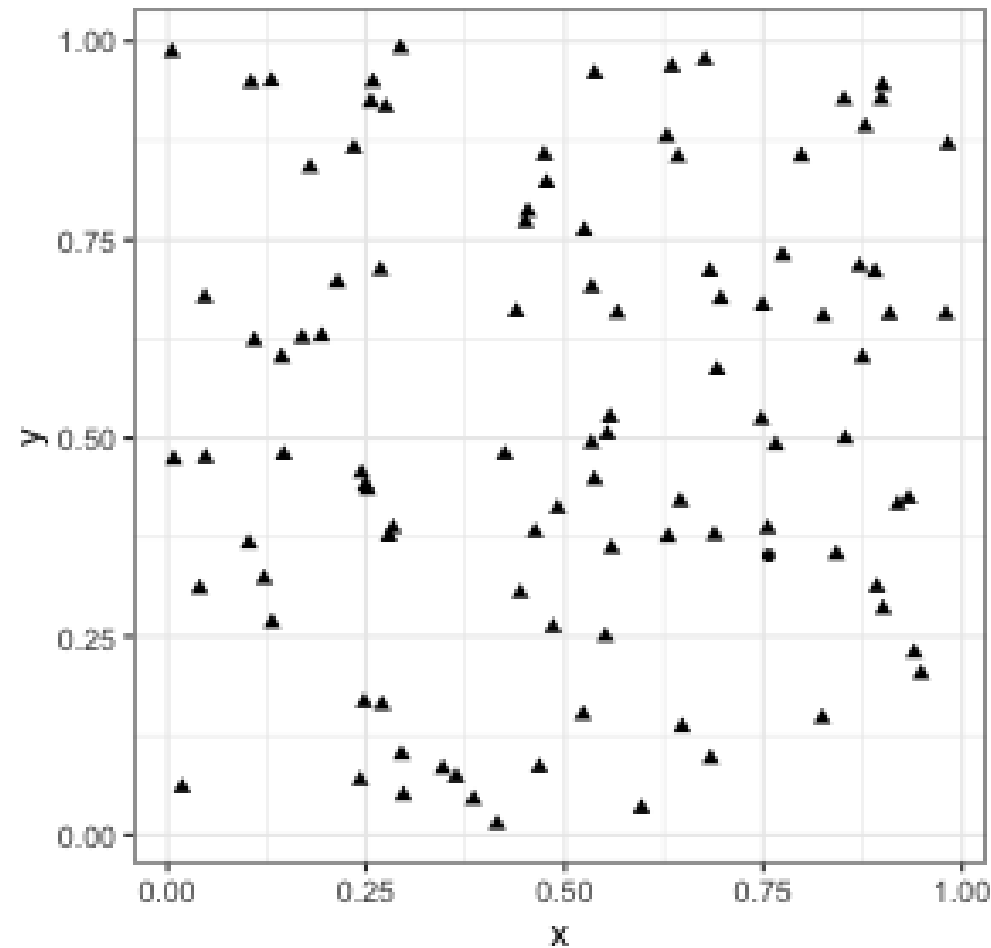




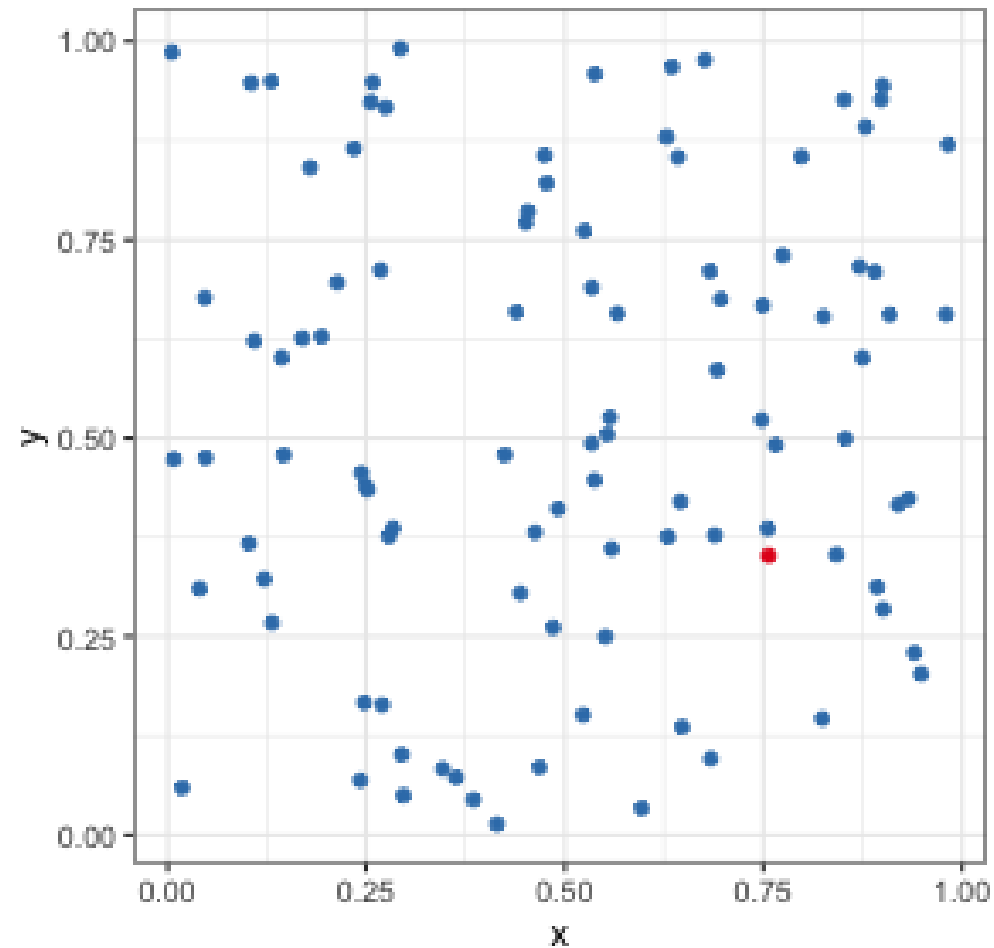
# Summary colour blindness

- Apply colourblind-friendly colourscales
  - + `scale_colour_viridis()`
  - + `scale_colour_brewer(palette = "Dark2")`
  - scico R package

# Pre-attentiveness: Find the odd one out?



# Pre-attentiveness: Find the odd one out?



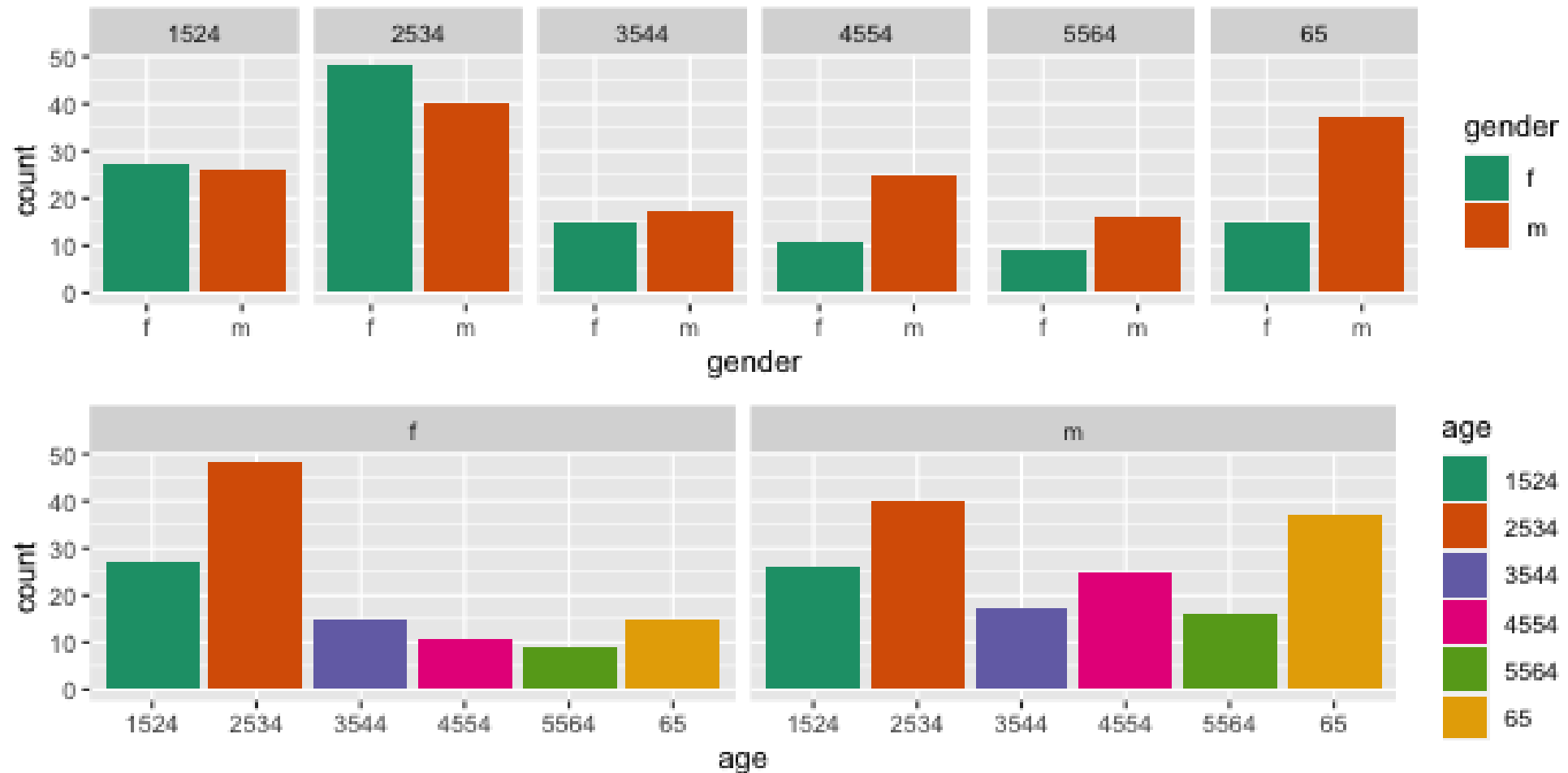
# Using proximity in your plots

Basic rule: place the groups that you want to compare close to each other

# Which plot answers which question?

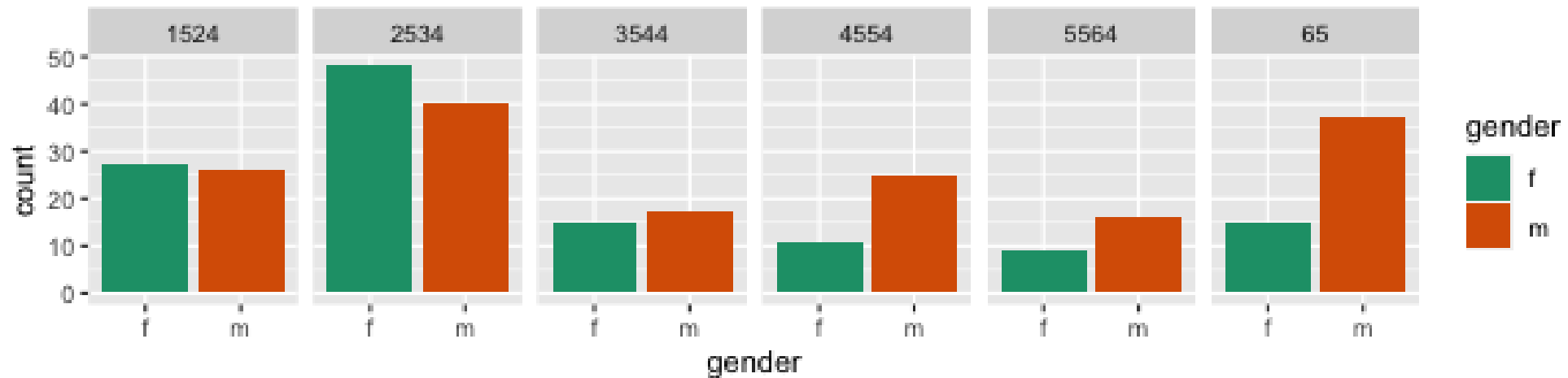
- "Is the incidence similar for males and females in 2012 across age groups?"
- "Is the incidence similar for age groups in 2012, across gender?"

# incidence similar for: (M and F) or (age, across gender) ?"



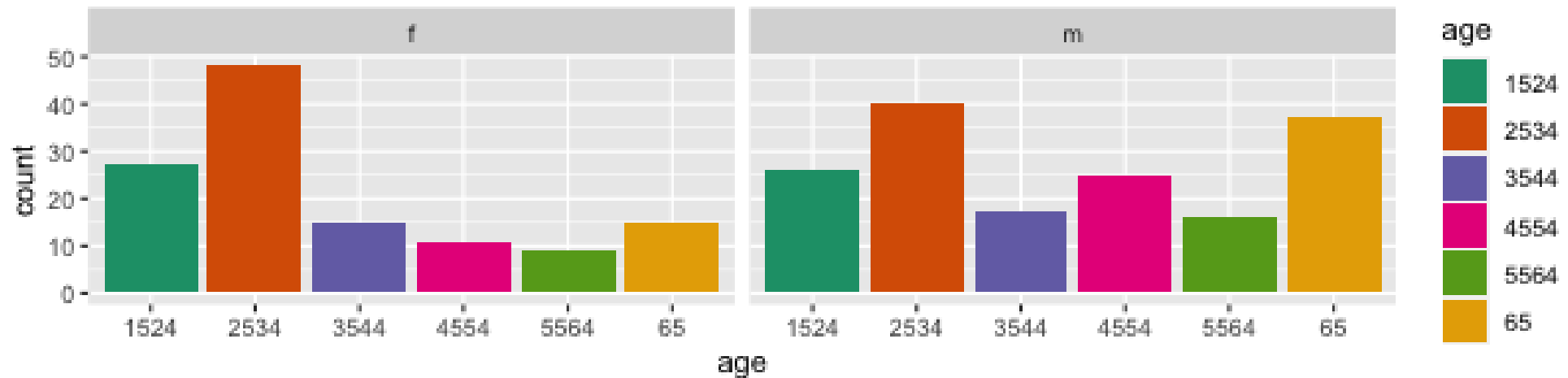


# "Incidence similar for M & F in 2012 across age?"



- Males & females next to each other: relative heights of bars is seen quickly.
- Auestion answer: "No, the numbers were similar in youth, but males are more affected with increasing age."

# "Incidence similar for age in 2012, across gender?"



- Puts the focus on age groups
- Answer to the question: "No, among females, the incidence is higher at early ages. For males, the incidence is much more uniform across age groups."

# Proximity wrap up

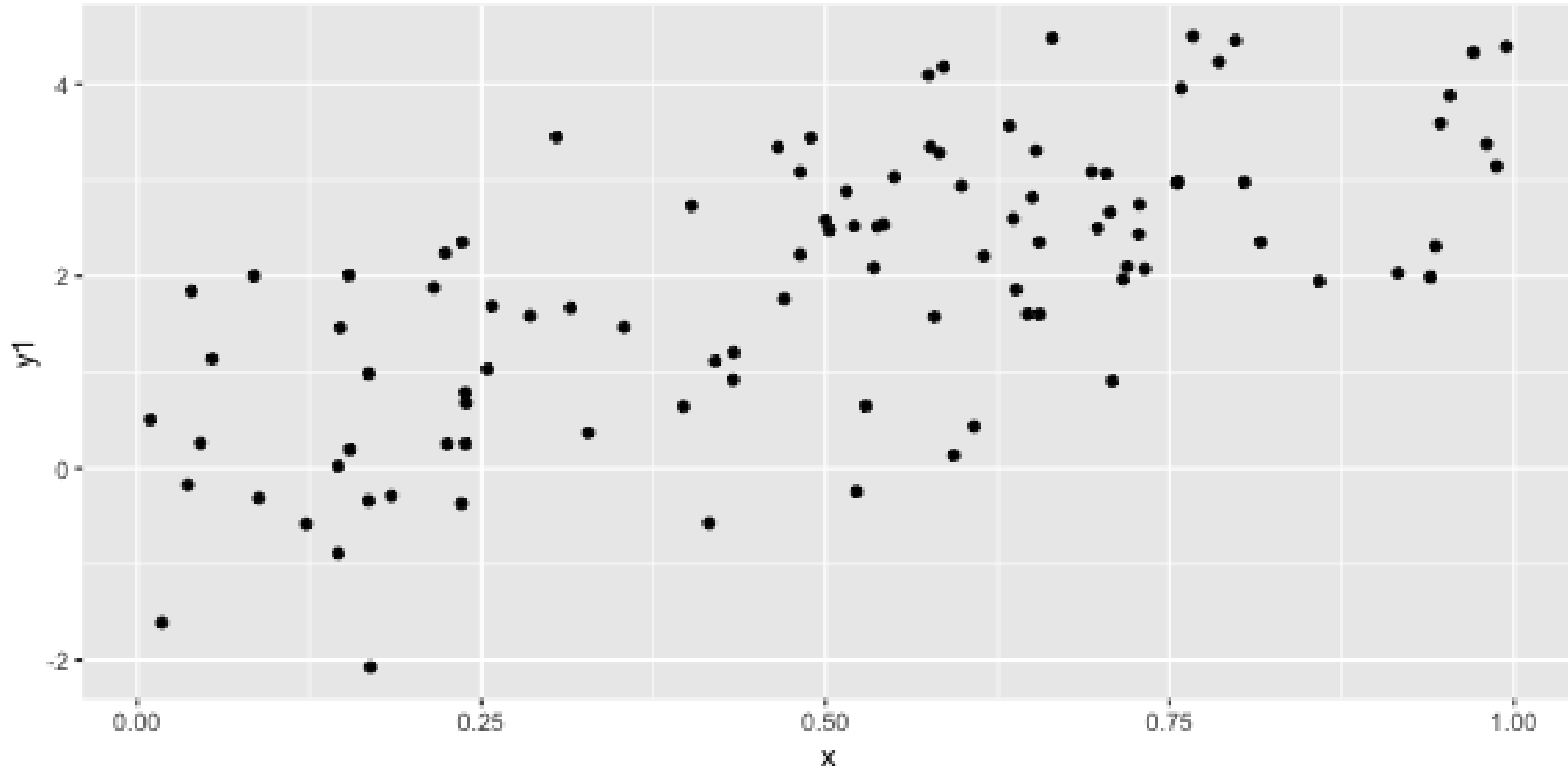
- Facetting of plots, and proximity are related to change blindness, an area of study in cognitive psychology.
- There are a series of fabulous videos illustrating the effects of making a visual break, on how the mind processes it by Daniel Simons lab.
- Here's one example:  
[The door study](#).

# Layering

- *Statistical summaries*: It is common to layer plots, particularly by adding statistical summaries, like a model fit, or means and standard deviations. The purpose is to show the **trend** in relation to the **variation**.
- *Maps*: Commonly maps provide the framework for data collected spatially. One layer for the map, and another for the data.

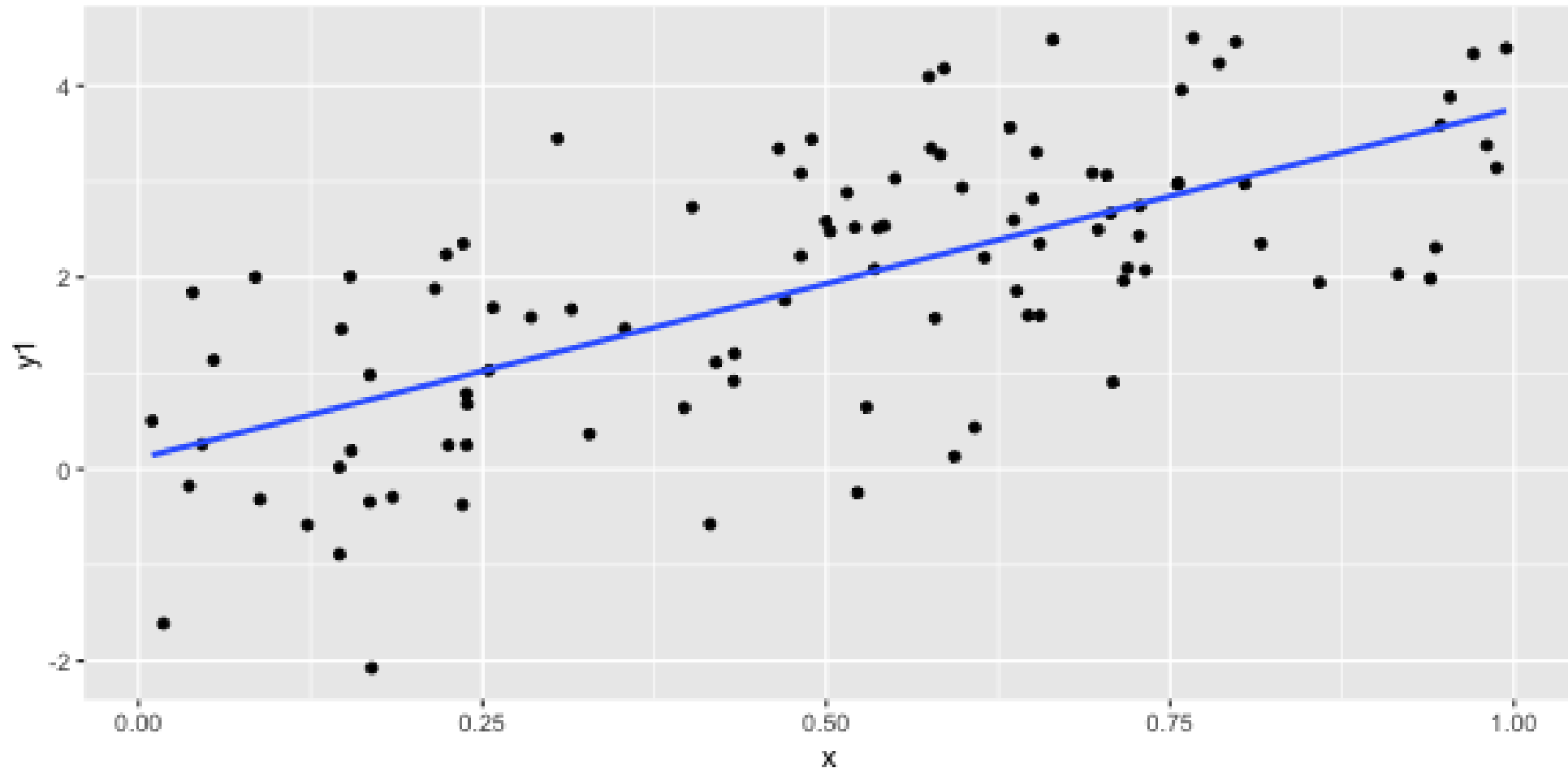
# geom\_point()

```
ggplot(df, aes(x = x, y = y1)) + geom_point()
```



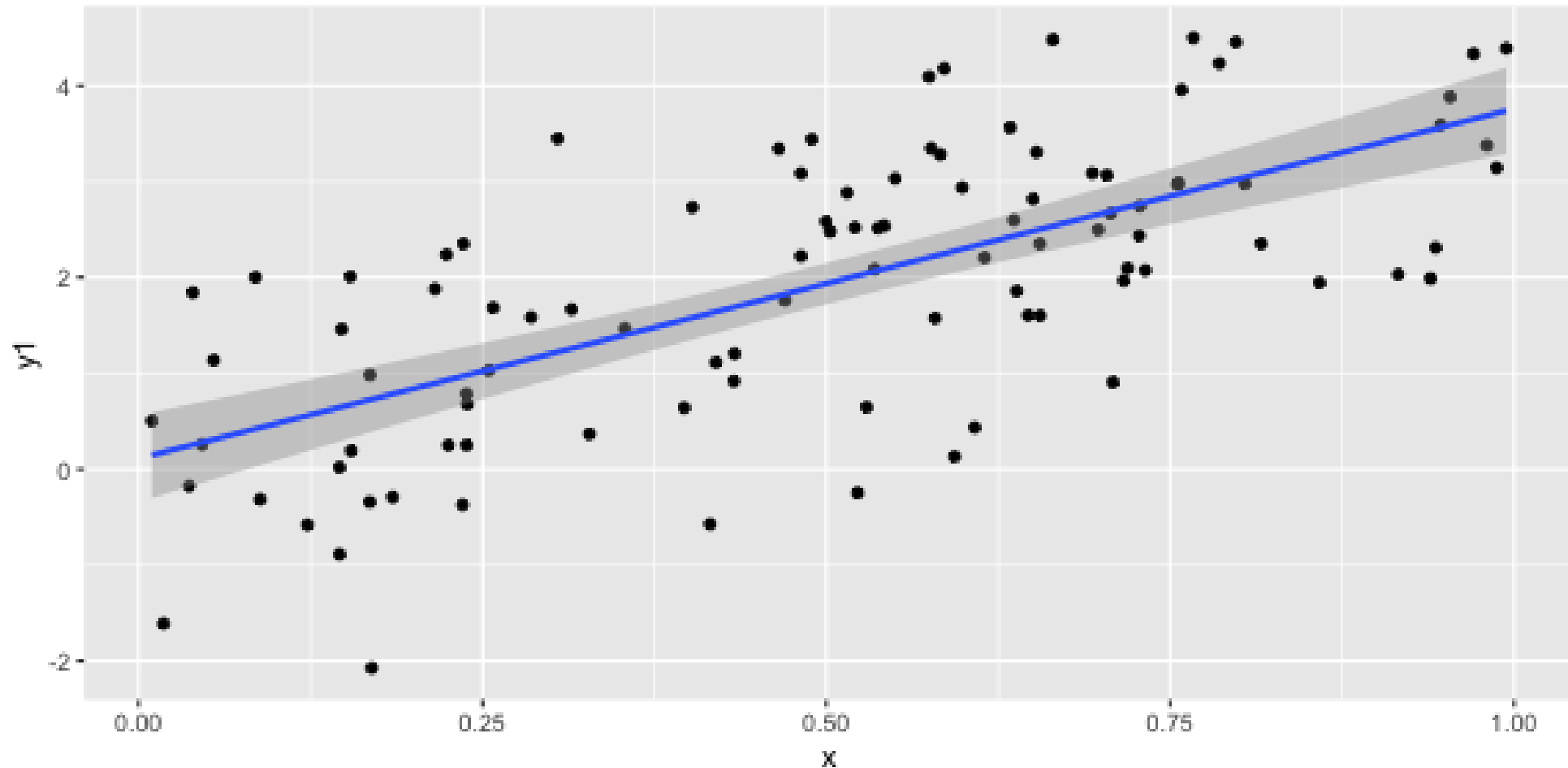
# geom\_smooth(method = "lm", se = FALSE)

```
ggplot(df, aes(x = x, y = y1)) + geom_point() +  
  geom_smooth(method = "lm", se = FALSE)
```



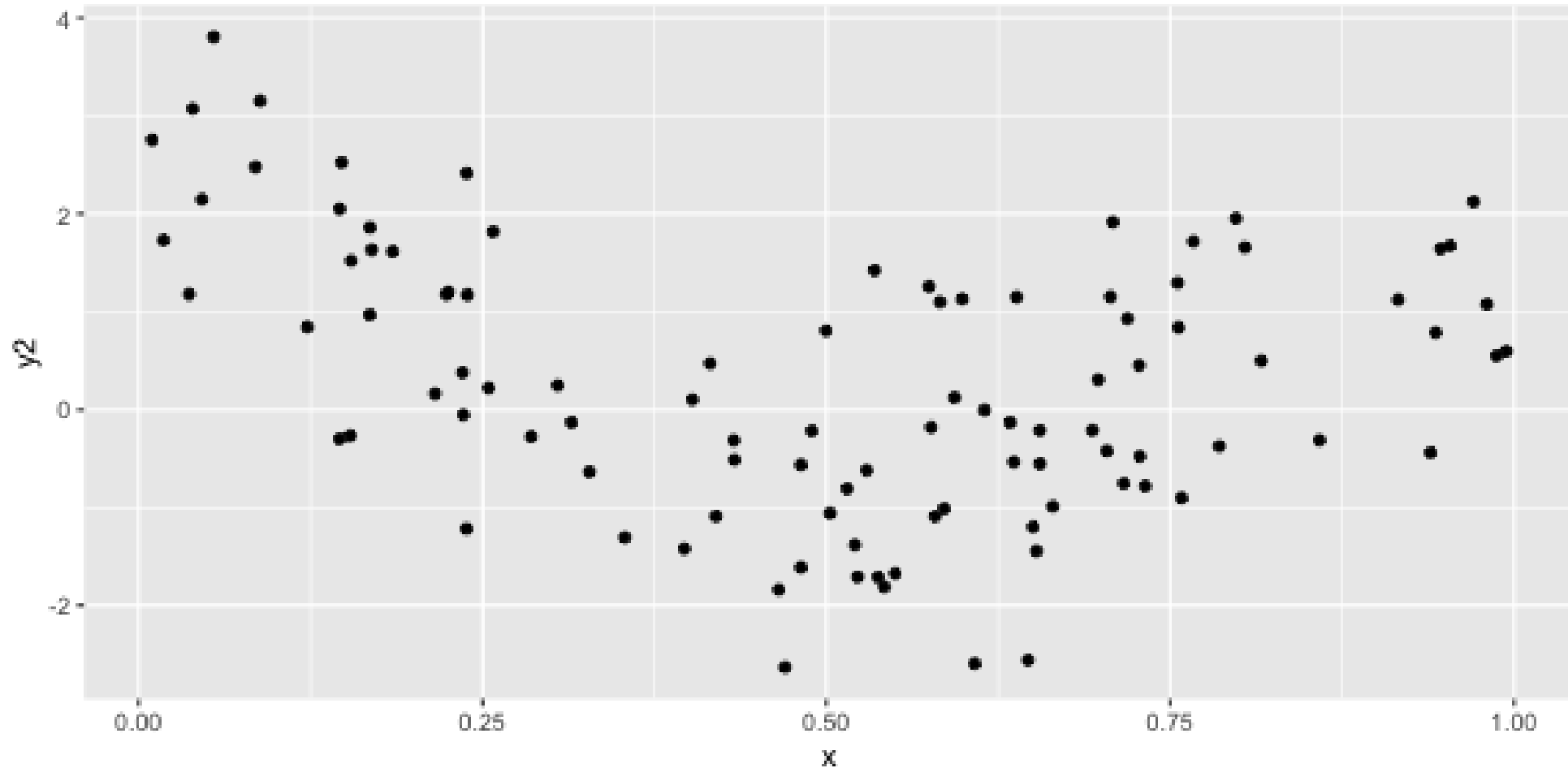
# geom\_smooth(method = "lm")

```
ggplot(df, aes(x = x, y = y1)) + geom_point() +  
  geom_smooth(method = "lm")
```



# geom\_point()

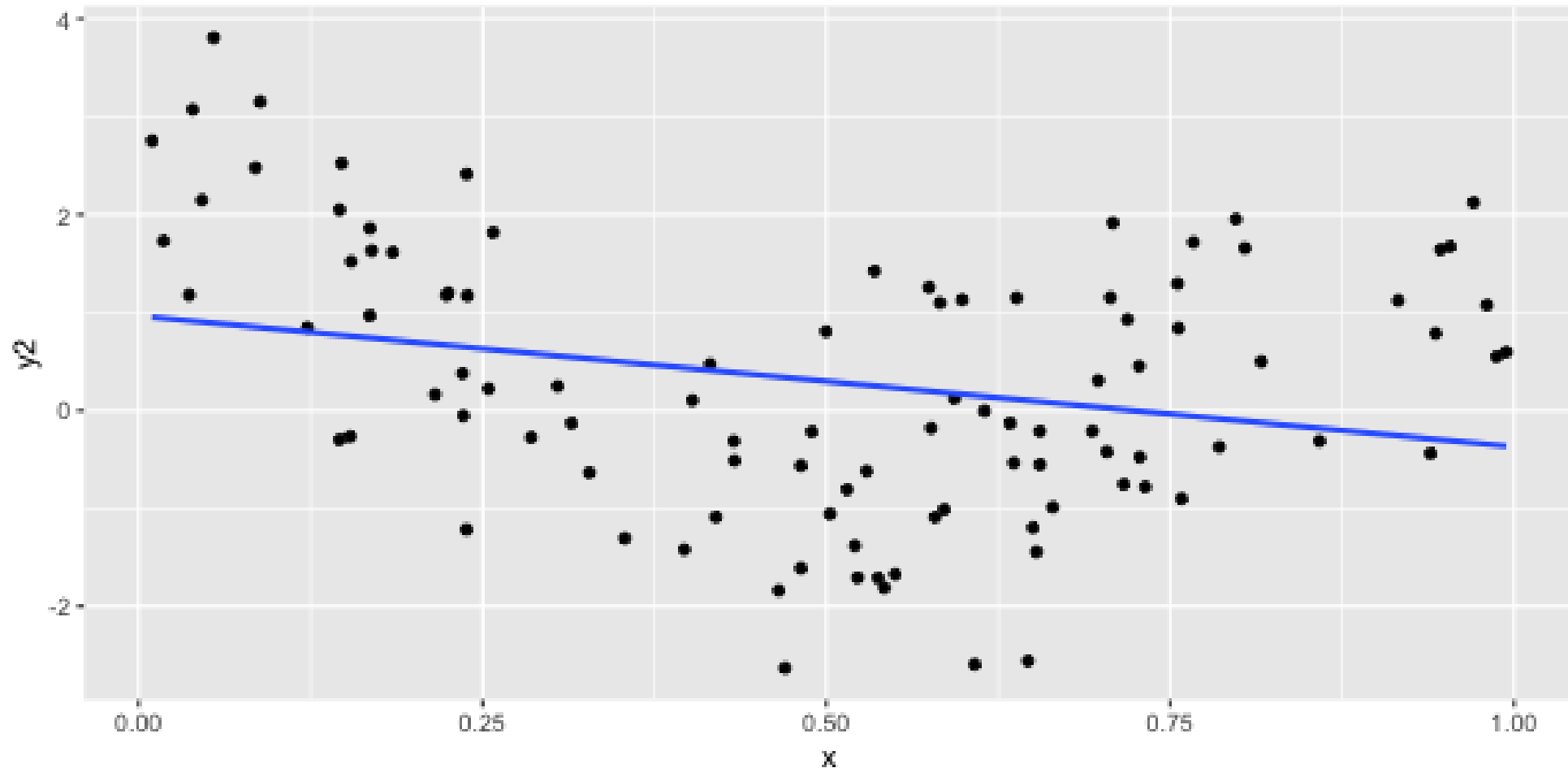
```
ggplot(df, aes(x = x, y = y2)) + geom_point()
```





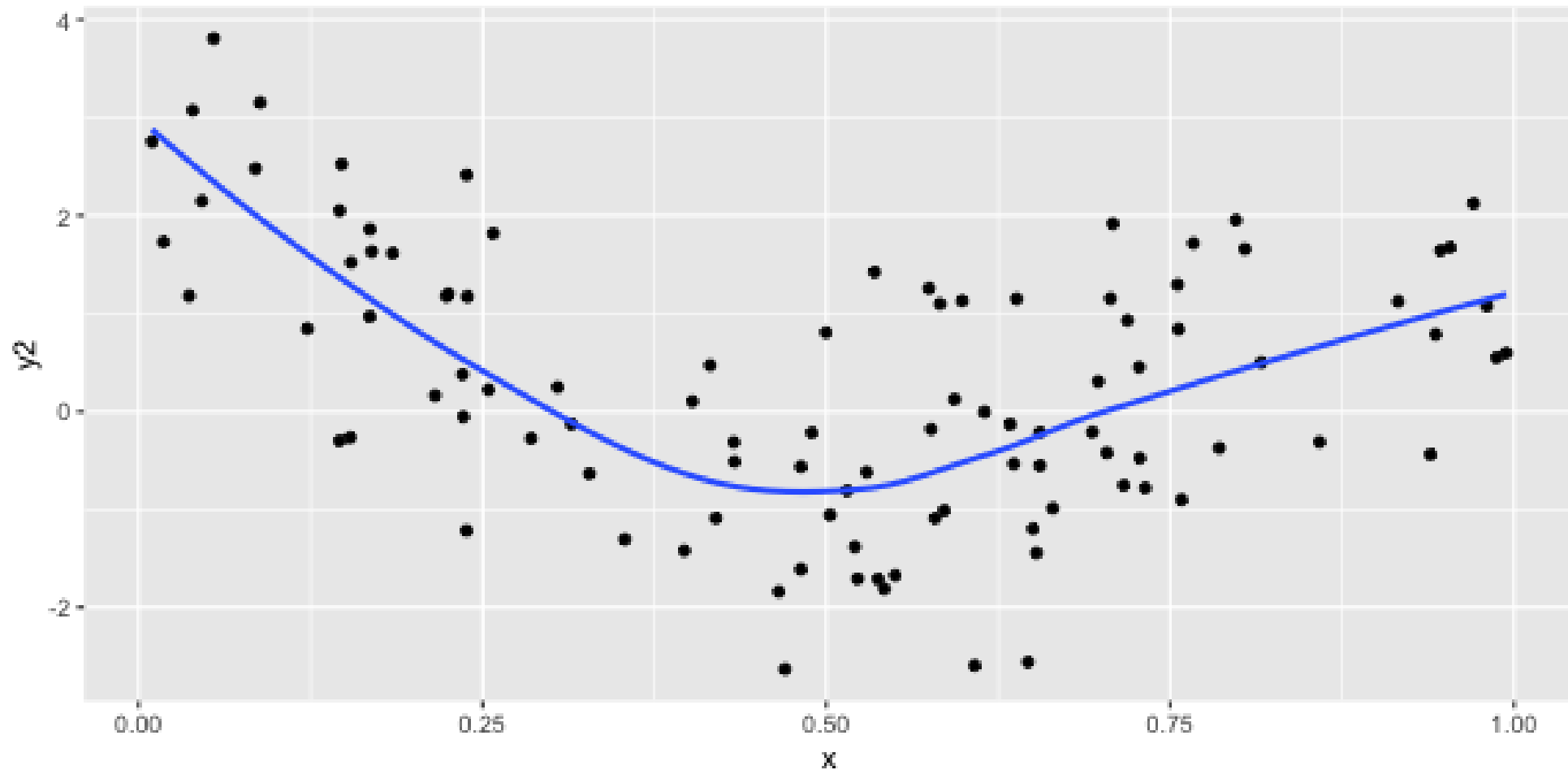
# geom\_smooth(method = "lm", se = FALSE)

```
ggplot(df, aes(x = x, y = y2)) + geom_point() +  
  geom_smooth(method = "lm", se = FALSE)
```



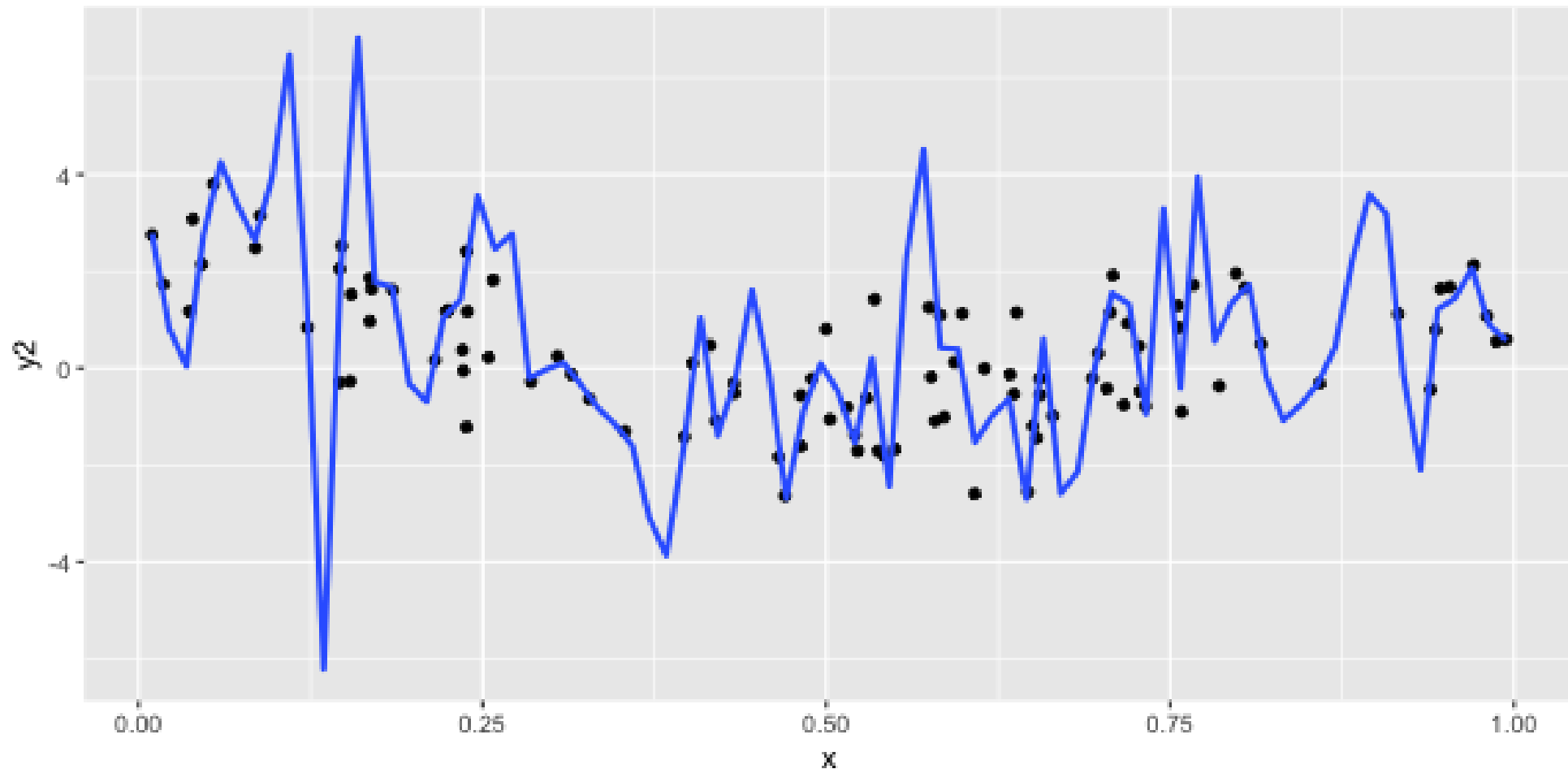
`geom_smooth(se = FALSE)`

```
ggplot(df, aes(x = x, y = y2)) + geom_point() +  
  geom_smooth(se = FALSE)
```



```
geom_smooth(se = FALSE, span = 0.05)
```

```
ggplot(df, aes(x = x, y = y2)) + geom_point() +  
  geom_smooth(se = FALSE, span = 0.05)
```



# geom\_smooth(se = FALSE, span = 0.2)

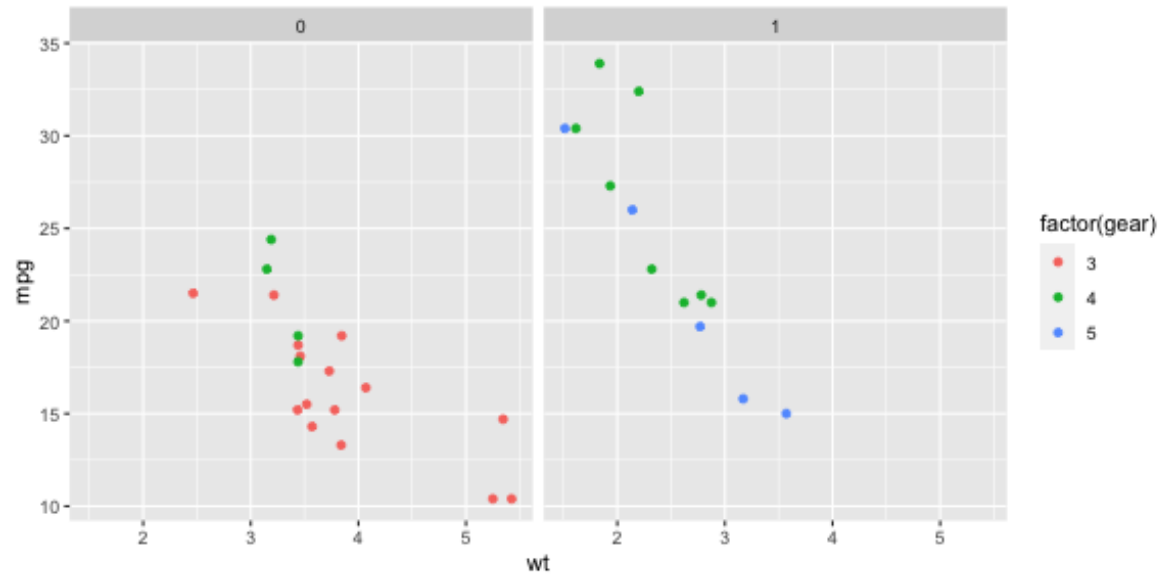
```
p1 <- ggplot(df, aes(x = x, y = y2)) + geom_point() +  
  geom_smooth(se = FALSE, span = 0.2)  
p1
```

# Interactivity with magic plotly

```
library(plotly)  
ggplotly(p1)
```

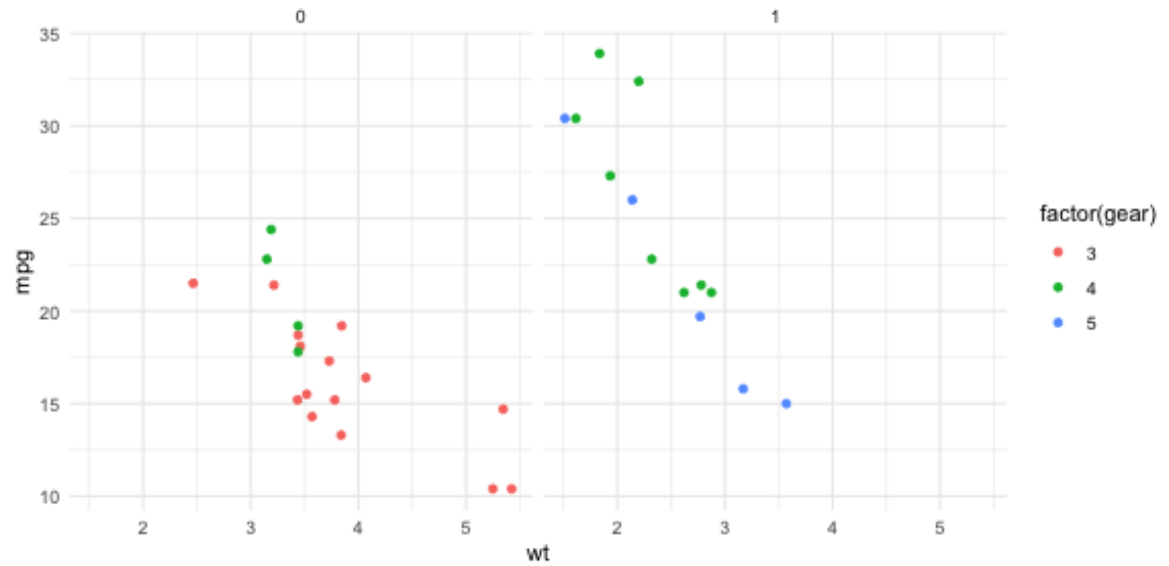
# Themes: Add some style to your plot

```
p <- ggplot(mtcars) +  
  geom_point(aes(x = wt,  
                 y = mpg,  
                 colour = factor(gear)))  
p  
p <- p + facet_wrap(~am)
```



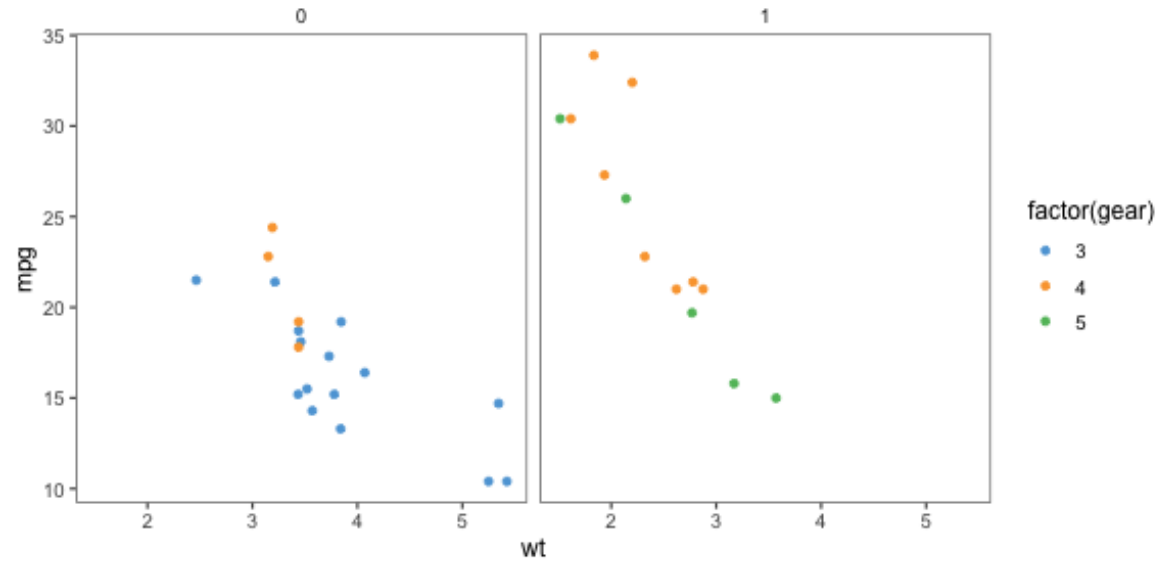
# Theme: theme\_minimal()

```
p +  
  theme_minimal()
```



# Theme: ggthemes theme\_few()

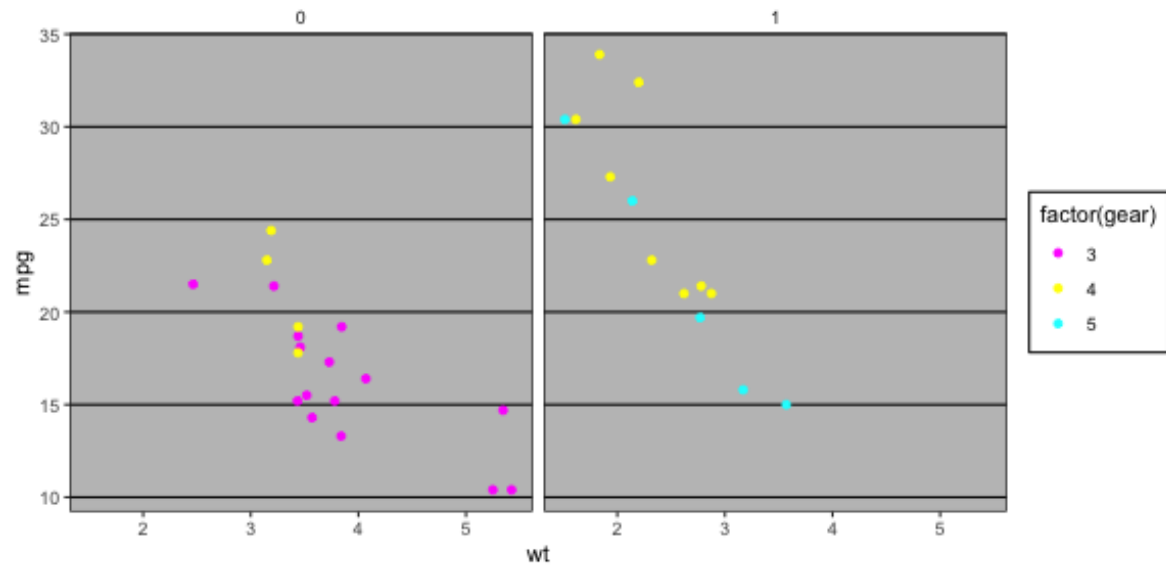
```
p +  
  theme_few() +  
  scale_colour_few()
```





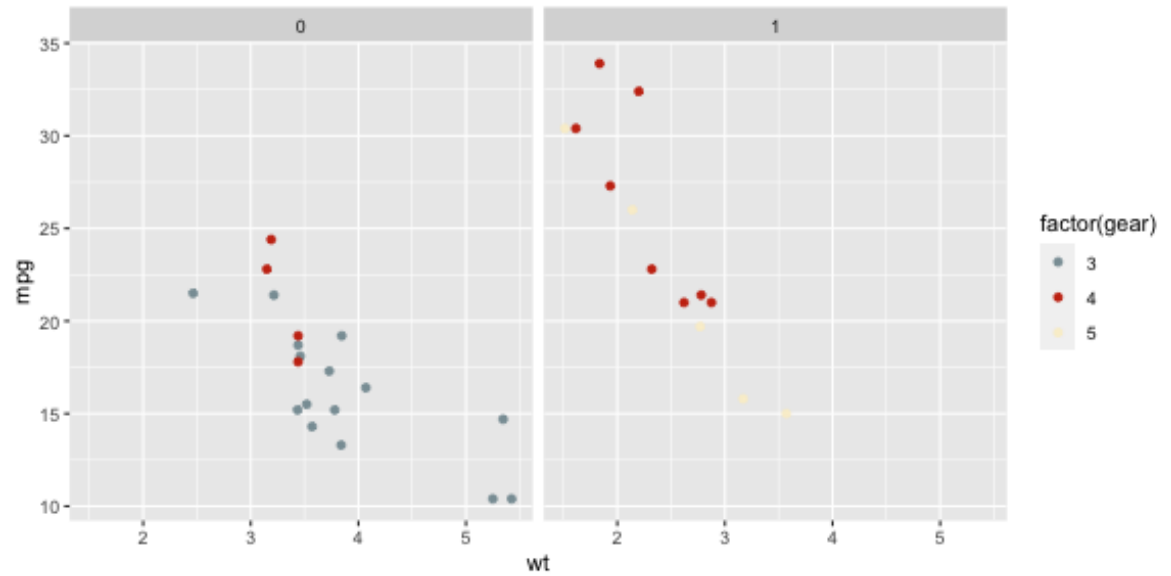
# Theme: ggthemes theme\_excel() 🤔

```
p +  
  theme_excel() +  
  scale_colour_excel()
```



# Theme: for fun

```
library(wesanderson)
p +
  scale_colour_manual(
    values = wes_palette("Royal")
  )
```



# Summary: themes

- The ggthemes package has many different styles for the plots.
- Other packages such as xkcd, skittles, wesanderson, beyonce, ochre, ....

# Hierarchy of mappings

1. Position - common scale (BEST): axis system
2. Position - nonaligned scale: boxes in a side-by-side boxplot
3. Length, direction, angle: pie charts, regression lines, wind maps
4. Area: bubble charts
5. Volume, curvature: 3D plots
6. Shading, color (WORST): maps, points coloured by numeric variable
  - [Di's crowd-sourcing expt](#)
  - Nice explanation by [Peter Aldous](#)
  - [General plotting advice and a book from Naomi Robbins](#)

# Your Turn:

- lab quiz open (requires answering questions from Lab exercise)
- go to rstudio and check out exercise 4-B
- If you want to use R / Rstudio on your laptop:
  - Install R + Rstudio (see )
  - open R
  - type the following:

```
# install.packages("usethis")  
library(usethis)  
use_course("mida.numbat.space/exercises/4b/mida-exercise-4b.zip")
```

# Resources

- Kieran Healy [Data Visualization](#)
- Winston Chang (2012) [Cookbook for R](#)
- Antony Unwin (2014) [Graphical Data Analysis](#)
- Naomi Robbins (2013) [Creating More Effective Charts](#)