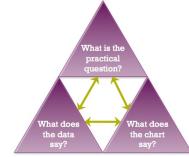


DATA

VISUAL

Week 1

	<p>Data Type:</p> <ul style="list-style-type: none">• Ratio• Interval• Nominal• Ordinal <p>DV Check-list:</p> <ul style="list-style-type: none">• Text• Arrangement• Color• Lines• Overall <p>Check-list Tool: https://datavizchecklist.stephanieevergreen.com/</p>	Junk Charts Trifecta Checkup 
<ul style="list-style-type: none">• Trifecta Check up• DV check list	<p>DV Ethical Principles:</p> <ul style="list-style-type: none">• Beneficence• Transparency• Accuracy• Objectivity• Respect• Accountability <p>Data Integrity</p> <ul style="list-style-type: none">• Permission• Security• Consent• Privacy and sensitive information• Data quality• Citation	

Basically introduction of what is DV, and how to evaluate quality of a DV. Principles of Ethical and Data Integrity is provided.

Week 2 - Story Telling

DV Story Telling Elements

Genre
Structure
Highlighting
Transition
Ordering
Interactivity

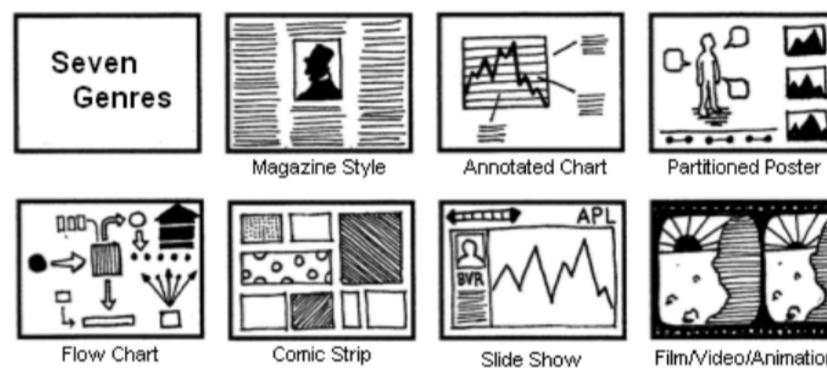
• Structure:



Intro → Supporting Facts → Main Insight → Conclusion

• Genre:

7 genres:



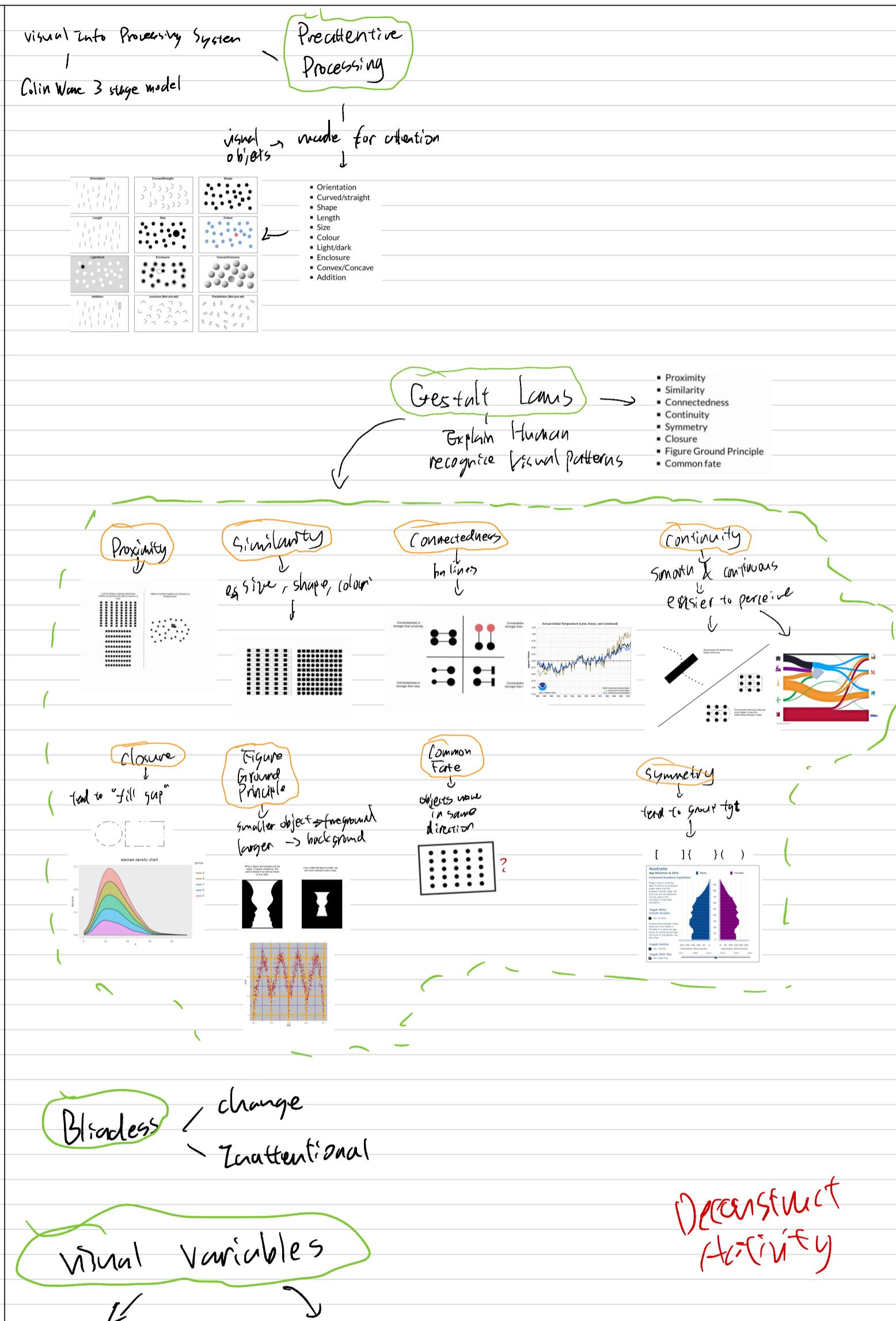
• Approach:

- | Author-driven | Mixed/hybrid | User-driven |
|--|--|--|
| <ul style="list-style-type: none">Linear order of ideasExtensive messagingNo interactivity | <ul style="list-style-type: none">Linear order to orientateMessaging to get audience startedInteractivity to dig deeper or personalise story | <ul style="list-style-type: none">Non-linear orderMinimal messagingExtensive interactivity |

• Strategies:

- | | |
|--|---|
| <ul style="list-style-type: none">1 Annotation2 Visual highlighting3 Matching content4 Progress Bars5 Consistent Visual Platform6 Multi-messaging | <ul style="list-style-type: none">7 Details on demand8 Timeline slider9 Tacit tutorial10 Semantically consistent11 Markers of interactivity12 Animated transitions |
|--|---|

WEEK 3 - Visual Perception and Color

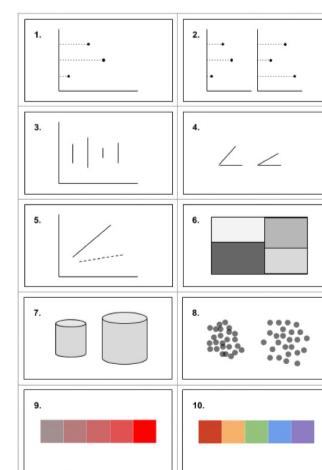


Feature/Object	Example	Notes	Nominal	Ordinal	Quantitative
Position	• • •	Powerful and accurate			✓
Length	— — —	Powerful and accurate			✓
Size	○ ○ ○	Less accurate, but still useful			✓
Angle/Slope	/ / /	Less accurate, but still useful			✓
Shape	○ □ △ ★	Can only use a few different shapes.	✓		
Colour - Hue	■ ■ ■ ■ ■	Can only use a few different hues.	✓		
Colour - Sequential	■ ■ ■ ■ ■	Can only use a few different ordered levels.		✓	
Colour - Diverging	■ ■ ■ ■ ■	Can only use a few different ordered levels.		✓	
Colour - Continuous	■ ■ ■ ■ ■	Less accurate than other quantitative encodings			✓
Texture	■ ■ ■ ■ ■	Bygone feature, but may still have applications.	✓		
Line - Type	— — — — —	Can only use a very limited number of types.	✓		
Line - Weight/Boldness	— — — — —	Not very accurate, but still useful.		✓	

Visual Comparison Accuracy

Visual Comparison Accuracy (Cleveland and McGill, 1985)

Rank (Most - Least)	Feature
1	Position - common scale
2	Position - same scale, but unaligned
3	Length
4	Angle
5	Slope
6	Area
7	Volume
8	Density
9	Colour saturation
10	Colour hue



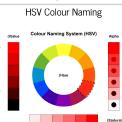
Week 3 - continued

• color:

visible light:

400nm - 700nm

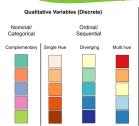
RGB



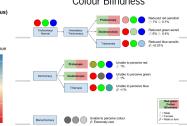
Color Brewer

- Select effective colour palettes based on:
 - Sequential
 - Diverging
 - Qualitative data
 - Single or multi-hue
- Colour-blind and print-safe palettes
- Web tool : <http://colorbrewer2.org/>

color scale



color Mind



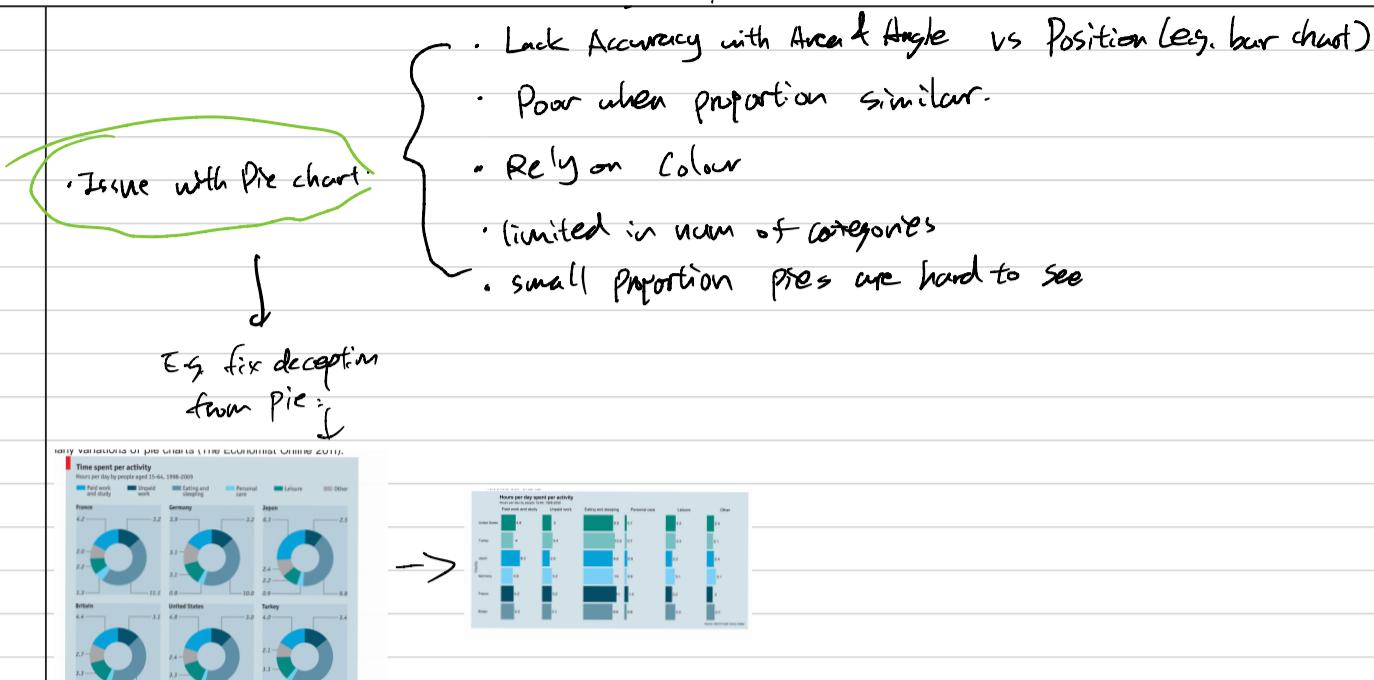
color Associations

Colour Associations

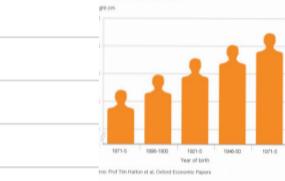
Color	Positive	Negative
Red	Romantic, strength, energy, heat, passion.	Blood, violence, dragon, anger, fire.
Green	Nature, spring, fertility, safety, growth, health.	Impression, decay, envy.
Yellow	Sun, summer, gold, forests, trees, energy, optimism.	Caveats, increase hazard, illness.
Blue	Sky, sea, stability, peace, calm, trust, reliability.	Depression, melancholy, sadness.
White	Snow, purity, peace, cleanliness, innocence, purity, light.	Cold, clinical, summer, death.
Grey	Intelligence, dignity, restraint, elegance, maturity.	Shadow, monotony, dramatic.
Black	Cold, power, humidity, depth, softness, style.	Fear, void, night, mystery, evil, uncertainty.

Week 4 - avoiding deception

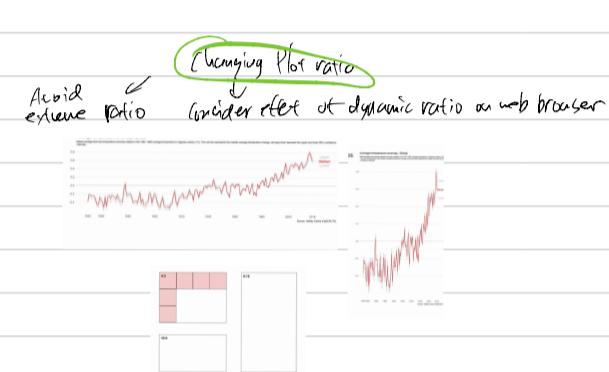
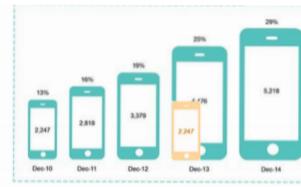
Issues



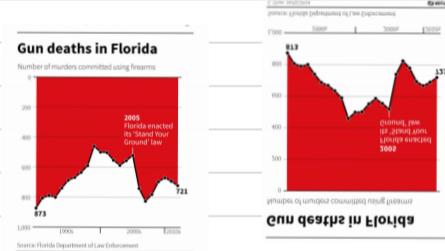
Trivariate Axes



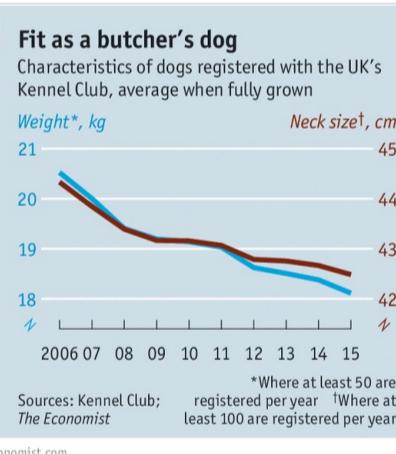
Use Area/size to depict quantity



Ignore convention



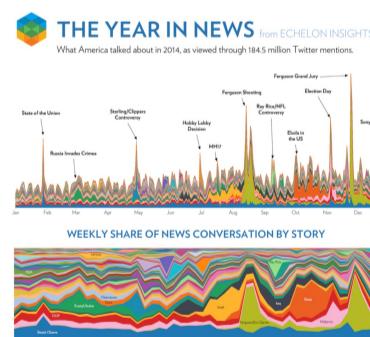
Dual Axes?



Solutions

Google Slides

Visual bombardment



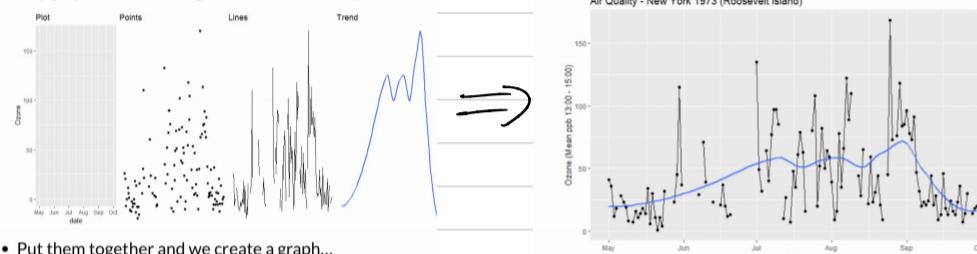
Week 5 - Grammar & vocabulary

• Layered Grammar of Graphics

Wickham proposed: a graphic is a series of layers consists of:

- a default dataset and set of mappings from variable: aesthetics,
- one or more layers, with each layer having
 - one geometric object
 - one statistical transformation,
 - one position adjustment,
 - and optionally, one dataset and set of aesthetic mappings
 - one scale for each aesthetic mapping used
- a coordinate system
- an optional facet specification.

- Any graphic can be thought of as a series of layers...



Represent Data / Statistical Transformation

- Boxes used in boxplots
- Bins used in histograms
- Bars used in barchart
- Points in a scatter plot
- Lines in a line chart



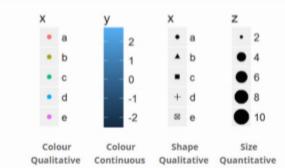
Mapping Aesthetics

- Properties of geometric objects, such as points, lines, colours, and shapes, are referred to as **aesthetics**
- The process of assigning variables from a dataset to aesthetics is known as **mapping**.
- For example in the air quality example $x = \text{Date}$ and $y = \text{Ozone}$.
- Points are geometric objects and the position they are drawn on the plot is determined by the mappings.

Scale

control mapping between a variable & aesthetic

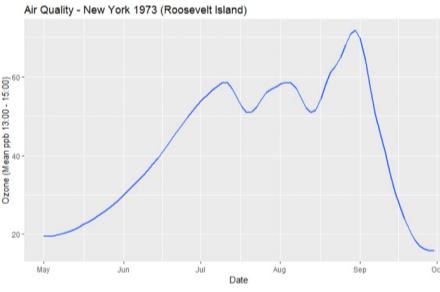
- For example, changing the range of the x axis, or changing the colour of a colour scale.



Statistical Transformation

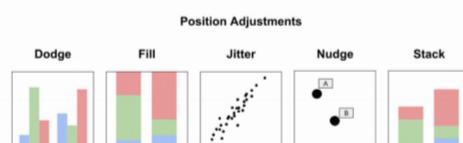
- Examples of stats transformations include the following:
 - quartiles of box plots
 - means
 - error bar/confidence intervals
 - binning in histograms and dot plots
 - tallies, counts, proportions, percentages in bar charts
 - lines of best fit for linear regression.
- Statistical transformations are the reason why statistics is so important for data visualisation.

The smoothed trend line in the Air Quality visualisation was estimated using a non-parametric, locally weighted regression model



Position Adjustments

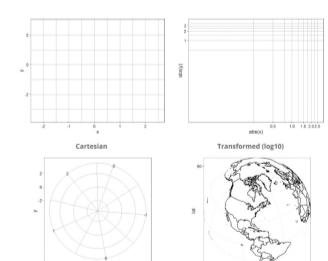
- Position adjustments aim to avoid overlapping elements by either dodging, filling, jittering, nudging or stacking. Layers can incorporate multiple position adjustments.



Coordinate System

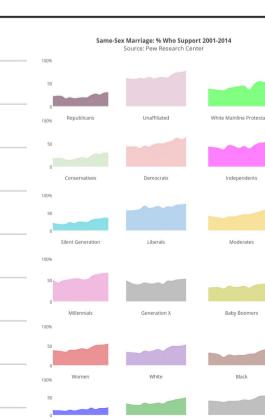
Determine placement of geometric objects within a plot

- Cartesian
- Transformed
- Polar
- Map



Faceting (latticing / trellising)

Break DV into small multiples



Week 5 - Grammar & vocabulary

- continue

R-package based on
Layered Grammar of Graphics

ggplot2.

Powerful grammar
for building DV.

Build DV layer
by layer

1. Install & Import:

```
install.packages("ggplot2")
library(ggplot2)
```

=>

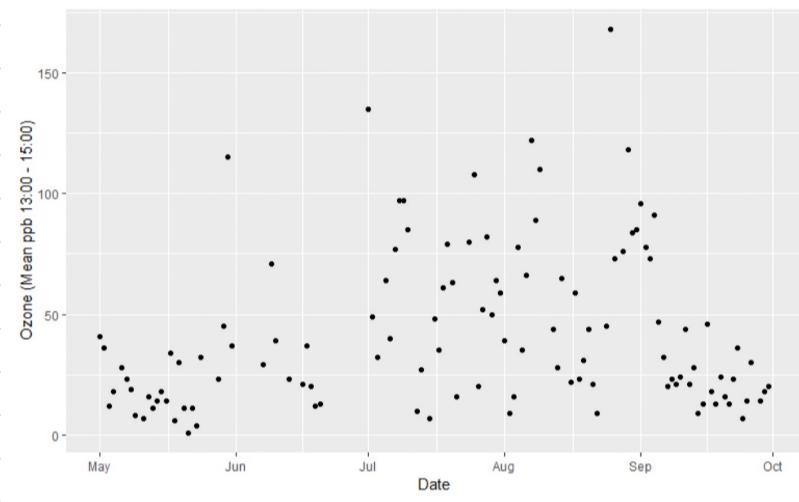
Nothing because we have not defined any layers.



2. Create ggplot2 object: (define coordinate system & select scales)

```
ggplot() +
  coord_cartesian() +
  scale_x_date(name = "Date") +
  scale_y_continuous(name = "Ozone (Mean ppb 13:00 - 15:00)")
```

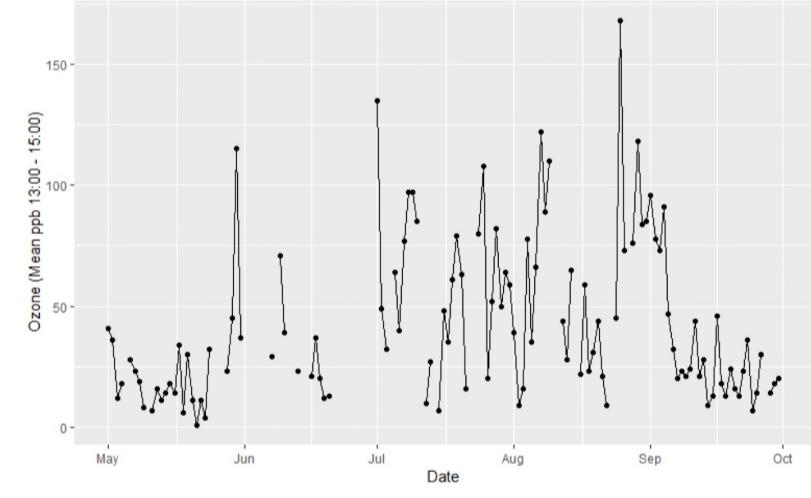
=>



3. Add layer for points

```
ggplot() +
  coord_cartesian() +
  scale_x_date(name = "Date") +
  scale_y_continuous(name = "Ozone (Mean ppb 13:00 - 15:00)") +
  layer(
    data = airquality,
    mapping = aes(x = date, y = Ozone),
    stat = "identity",
    geom = "point",
    position = position_identity()
  )
```

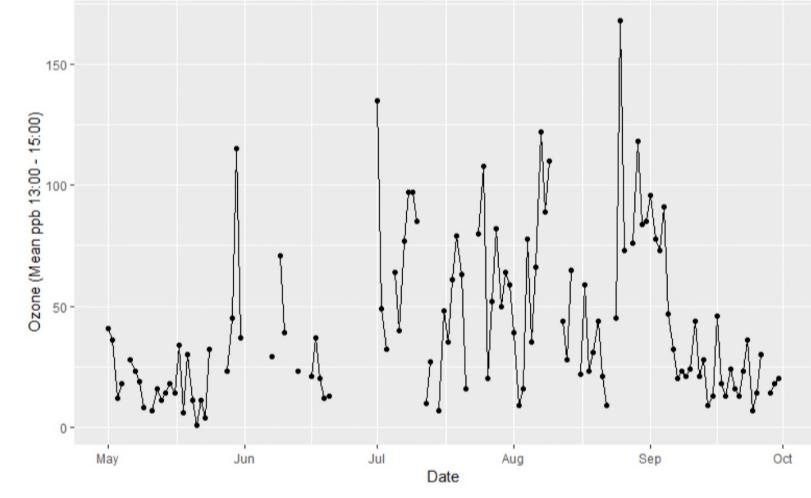
=>



4. Add lines

```
ggplot() +
  coord_cartesian() +
  scale_x_date(name = "Date") +
  scale_y_continuous(name = "Ozone (Mean ppb 13:00 - 15:00)") +
  layer(
    data = airquality,
    mapping = aes(x = date, y = Ozone),
    stat = "identity",
    geom = "point",
    position = position_identity()
  ) +
  layer(
    data = airquality,
    mapping = aes(x = date, y = Ozone),
    stat = "identity",
    geom = "line",
    position = position_identity()
  )
```

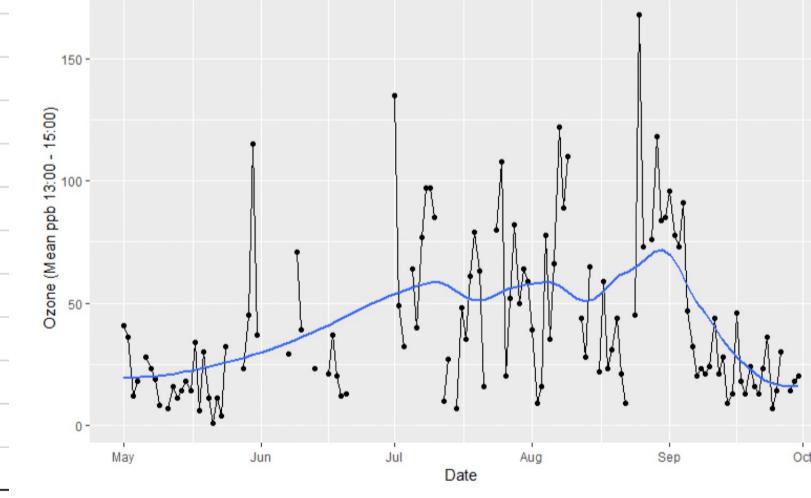
=>



5. Trend lines

```
ggplot() +
  coord_cartesian() +
  scale_x_date(name = "Date") +
  scale_y_continuous(name = "Ozone (Mean ppb 13:00 - 15:00)") +
  layer(
    data = airquality,
    mapping = aes(x = date, y = Ozone),
    stat = "identity",
    geom = "point",
    position = position_identity()
  ) +
  layer(
    data = airquality,
    mapping = aes(x = date, y = Ozone),
    stat = "identity", → Trend ?
    geom = "line",
    position = position_identity()
  )
```

=>

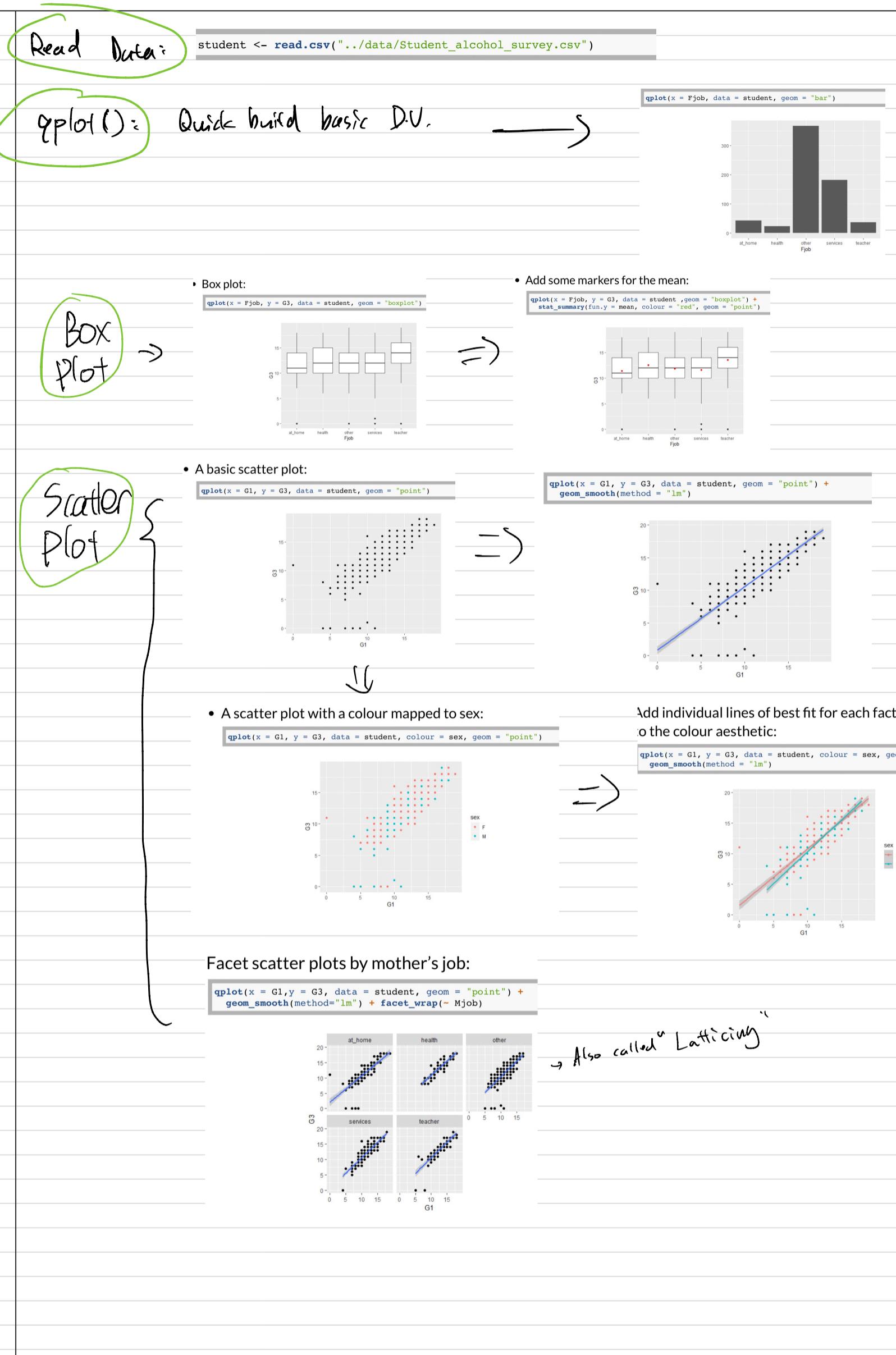


- The following code will reproduce the same visualisation

```
p <- ggplot(data = airquality, aes(x = date, y = Ozone))
p + geom_point() +
  geom_line(aes(group = 1)) +
  geom_smooth(se = FALSE, span = 0.4) +
  labs(
    title = "Air Quality - New York 1973 (Roosevelt Island)",
    x = "Date",
    y = "Ozone (Mean ppb 13:00 - 15:00"
  )
```

Week 5 - Grammar & vocabulary

- continue,



Week 5 - Grammar & vocabulary

- continue

ggplot()

: more powerful tool building DV.

1. Define ggplot object

```
p1 <- ggplot(data = student, mapping = aes(x = as.factor(Dalc), y = G3))
```

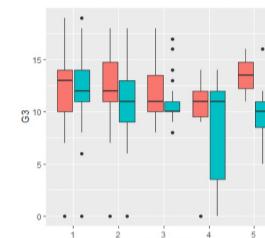
- aes are the aesthetic mappings. Any layers added will map the data to the corresponding aesthetics in the geom. For example, a box plot.

2. Map variables : (geom)

- We can map additional variables to other scales...

• Note how we mapped a fill aesthetic to sex.

```
p2 <- ggplot(student, aes(x = as.factor(Dalc), y = G3, fill = sex))  
p2 + geom_boxplot()
```

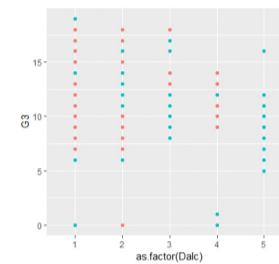


- Box plots hide sample size. Use a different geom that conveys sample size...

```
p3 <- ggplot(student, aes(x = as.factor(Dalc), y = G3, colour = sex))  
p3 + geom_point()
```



- Note colour refers to the outline or solid colour of an object. fill refers to the inside colour of a large object, e.g. bar or box.

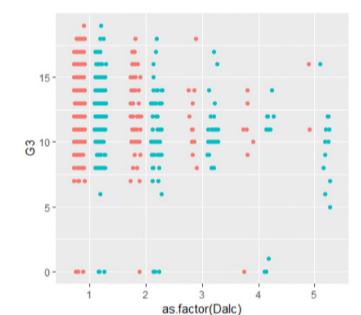


3. Points are overlapping

Solve: overlapping =>

Use position adjustments to avoid over-plotting

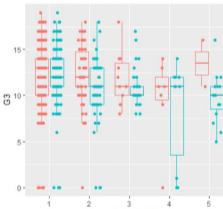
```
p3 + geom_jitter(position = position_jitterdodge())
```



4. Additional geoms

Note how we use outlier.shape = NA to suppress the outliers plotted for the boxplot

```
p3 + geom_jitter(position = position_jitterdodge()) +  
  geom_boxplot(fill = NA, outlier.shape = NA)
```

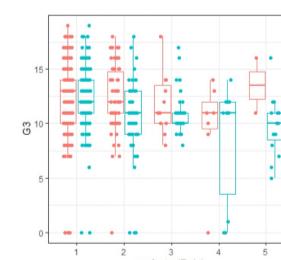


5. Themes



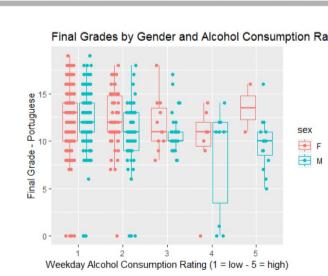
You can change default themes:

```
p3 + geom_jitter(position = position_jitterdodge()) +  
  geom_boxplot(fill = NA, outlier.shape = NA) + theme_bw()
```



6. Adding Titles & Labels

```
p3 + geom_jitter(position = position_jitterdodge()) +  
  geom_boxplot(fill = NA, outlier.shape = NA) +  
  labs(title = "Final Grades by Gender and Alcohol Consumption Ratings",  
       x = "Weekday Alcohol Consumption Rating (1 = low - 5 = high)",  
       y = "Final Grade - Portuguese")
```



7. Saving

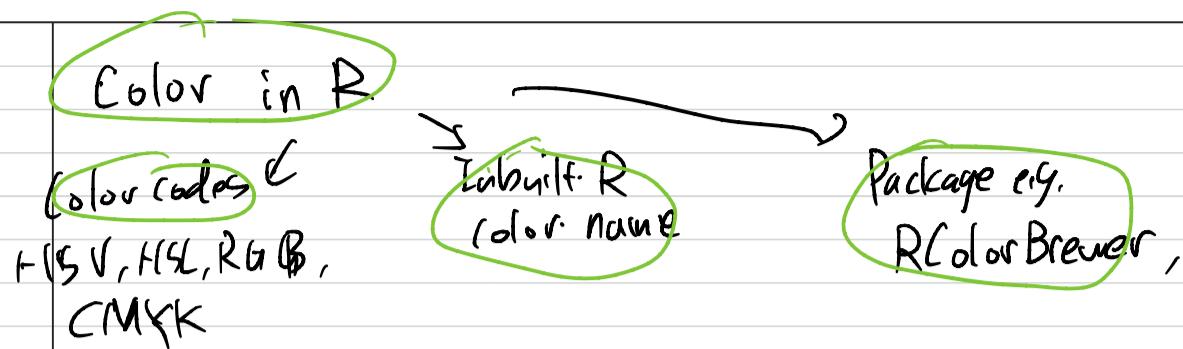
as image, PDF, or Clipboard

ggsave() → save as image with sensible defaults
to your working directory.

```
ggsave("Box_plot_Grades_Gender_Consumption_01.png",  
       width = 18, height = 12, units = "cm")
```

Week 5 - Grammar & vocabulary

- continue



colourpicker package:

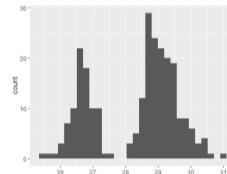
```
install.packages("colourpicker")
library(colourpicker)
```

Access Colour Picker
in Addins menu in RStudio

Basic Colour Assignment:

- Produce a histogram showing the distribution of pizza diameter.

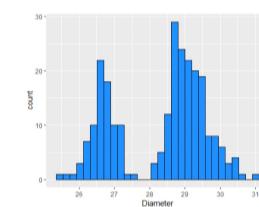
```
Pizza <- read.csv("../data/Pizza.csv")
p1 <- ggplot(data = Pizza, aes(x = Diameter))
p1 + geom_histogram()
```



- Change colour using R colour names

```
p1 + geom_histogram(fill = "dodgerblue", colour = "black")
```

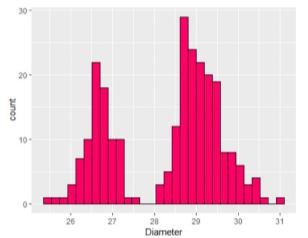
=>



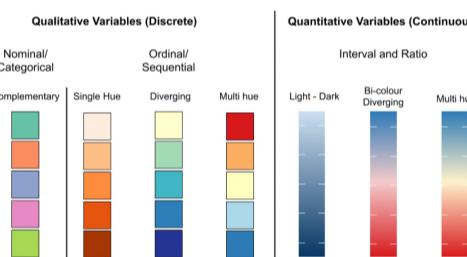
- Change colour using hex codes

```
p1 + geom_histogram(fill = "#ff0066", colour = "#000000")
```

=>



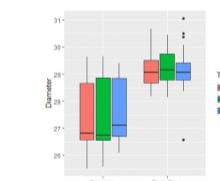
Color Scales



Nominal Color Scale E.g.

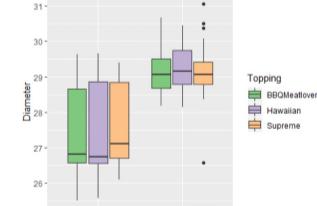
- Side-by-side box plot comparing pizza diameter by Store and Topping

```
p2 <- ggplot(data = Pizza, aes(x = Store, y = Diameter, fill = Topping))
p2 + geom_boxplot()
```



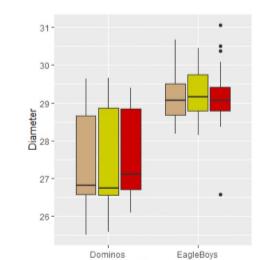
- Change ColourBrewer palette...

```
p2 + geom_boxplot() + scale_fill_brewer(palette = "Accent")
```



- Set manual colour scale...

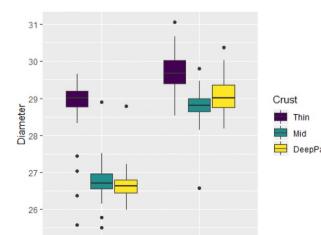
```
p2 + geom_boxplot() + scale_fill_manual(
  values = c("burlywood3", "yellow3", "red3"))
```



Ordinal Color Scale E.g.

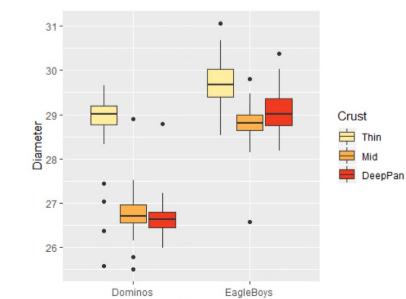
- Boxplot comparing pizza diameter by Store and Crust.

```
Pizza$Crust<-factor(Pizza$Crust, levels = c("Thin", "Mid", "DeepPan"),
ordered = T)
p3 <- ggplot(data = Pizza, aes(x = Store, y = Diameter, fill = Crust))
p3 + geom_boxplot()
```



- Change palette(...)

```
p3 + geom_boxplot() + scale_fill_brewer(palette = "YlOrRd")
```



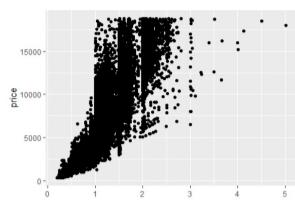
Week 5 - Grammar & vocabulary

- continue

Continuous Color Scale

- Explore the relationship between diamond carat and price

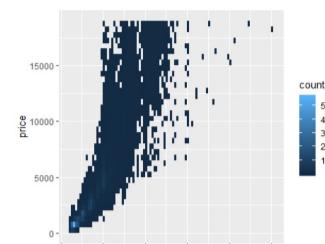
```
Diamonds <- read.csv("../data/Diamonds.csv")
p4 <- ggplot(data = Diamonds, aes(x = carat, y = price))
p4 + geom_point()
```



- Hard to see data density as $n = 53940$.

- Use a continuous colour scale to represent data density in a 2d histogram

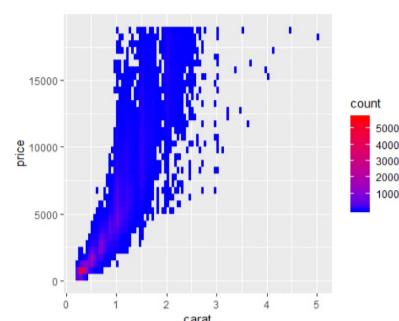
```
p4 + geom_bin2d(binwidth = c(0.05, 500))
```



⇒

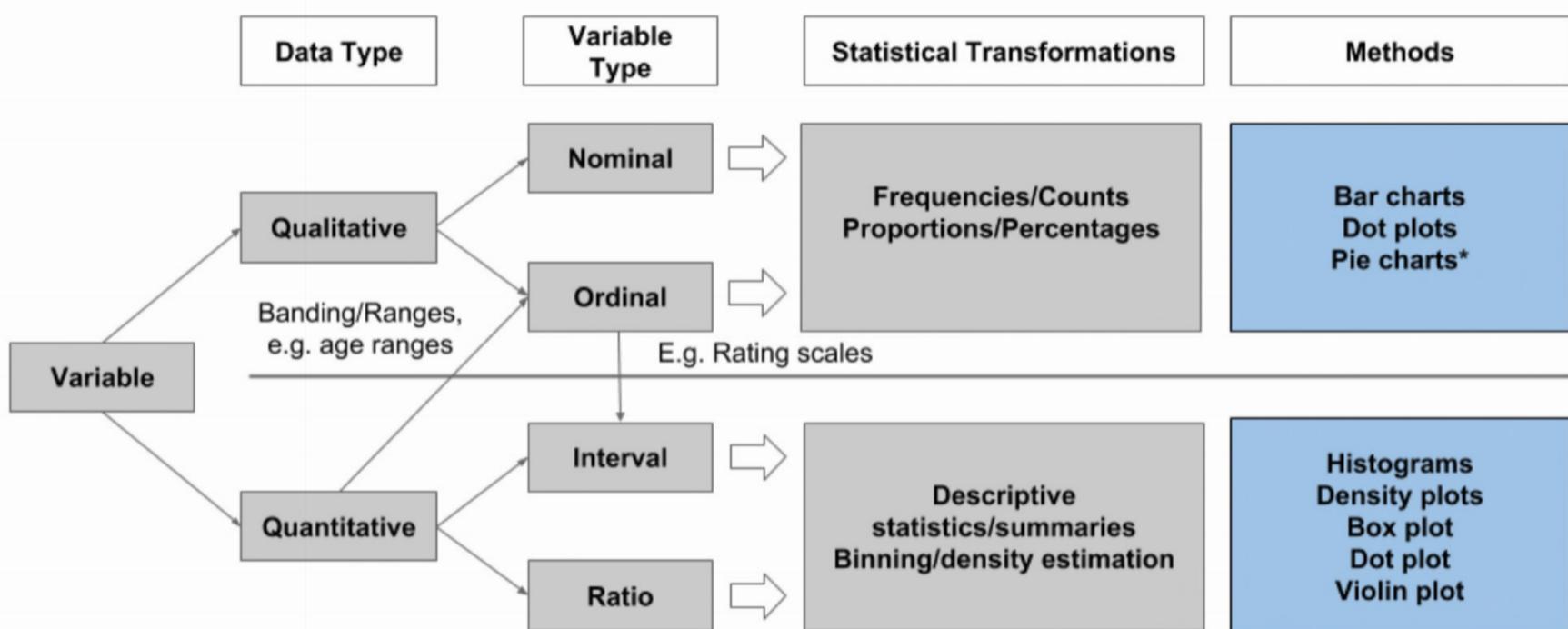
- Change colour gradient...

```
p4 + geom_bin2d(binwidth = c(0.05, 500)) +
  scale_fill_gradient(low="blue", high="red")
```

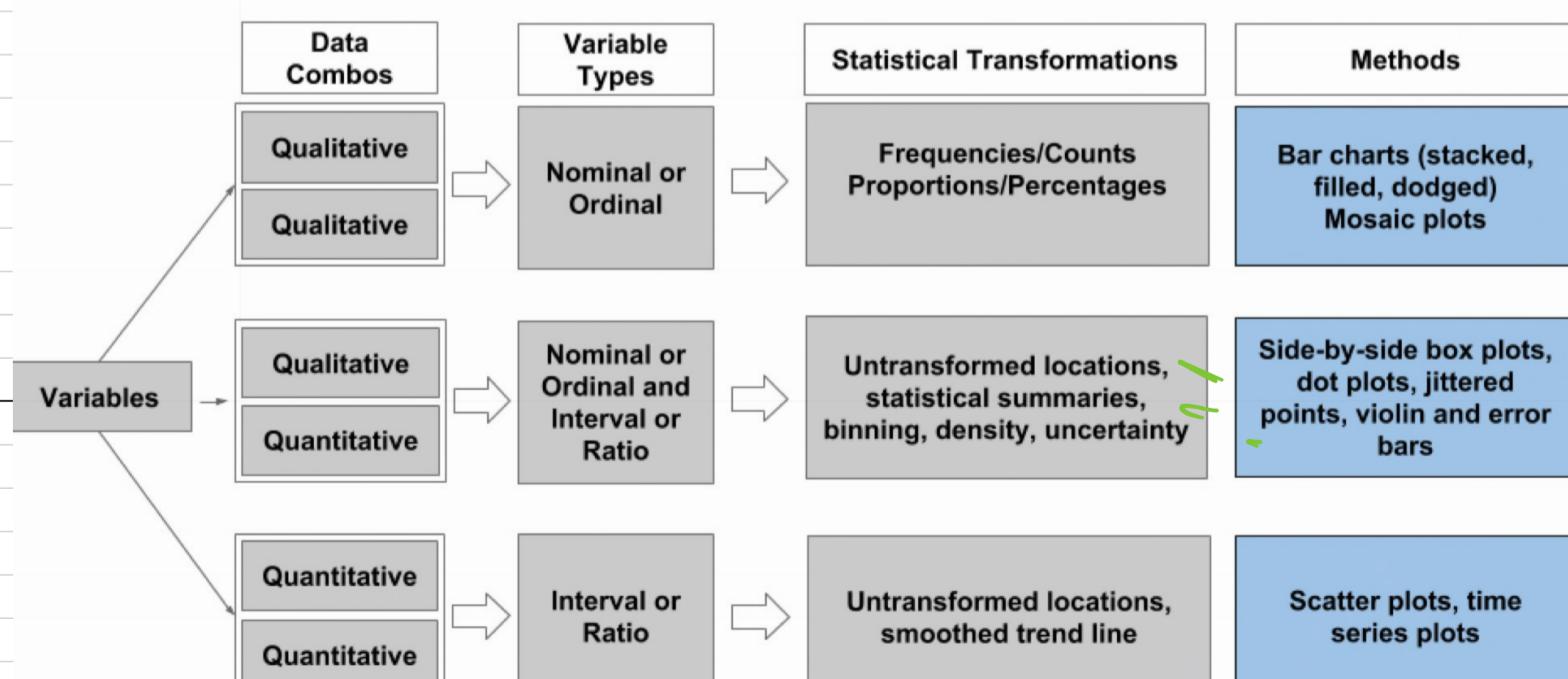


⇒

Common Univariate Methods



Common Bivariate Methods



Week 6 - Multivariate Strategies

- Mapping Additional Aesthetics
- Faceting
- 3D Scatter

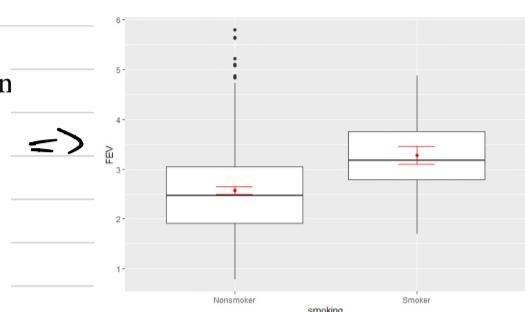
- Force Expiratory Volume Data:

- Do children or adolescent smokers tend to have lower FEV readings than non-smokers?

Bivariate Thinking:

- Let's start with a side-by-side box plot with mean and error bars (95% CI)

```
FEV = read.csv("../data/FEV.csv")
p1 <- ggplot(data = FEV, aes(x = smoking, y = FEV))
p1 + geom_boxplot() + stat_summary(fun.y = "mean", geom = "point",
                                    colour = "red") +
  stat_summary(fun.data = "mean_cl_boot", colour = "red",
              geom = "errorbar", width = .2)
```

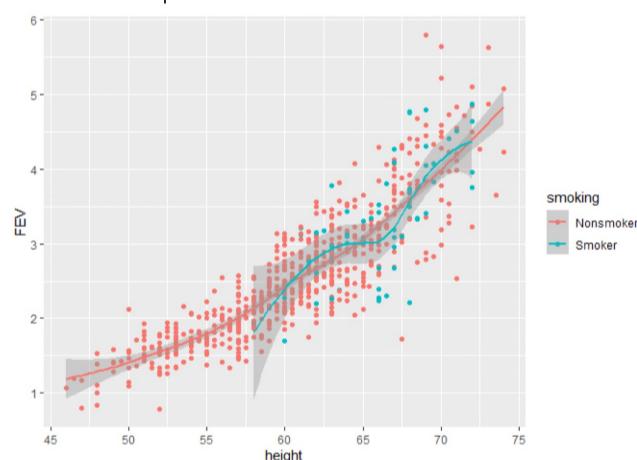


Explore Covariation:

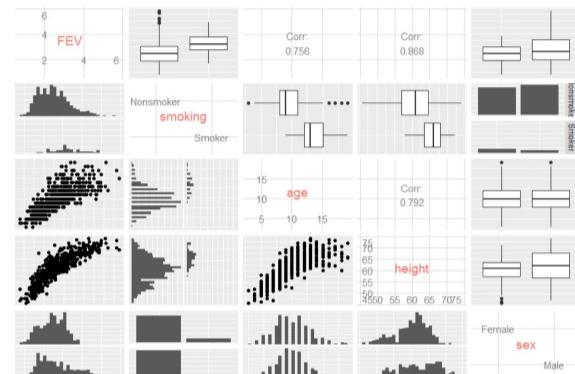
3 variables

Does height explain difference in FEV for smokers & non-smoker?

```
p2 <- ggplot(data = FEV, aes(x = height, y = FEV, colour = s
p2 + geom_point() + geom_smooth()
```

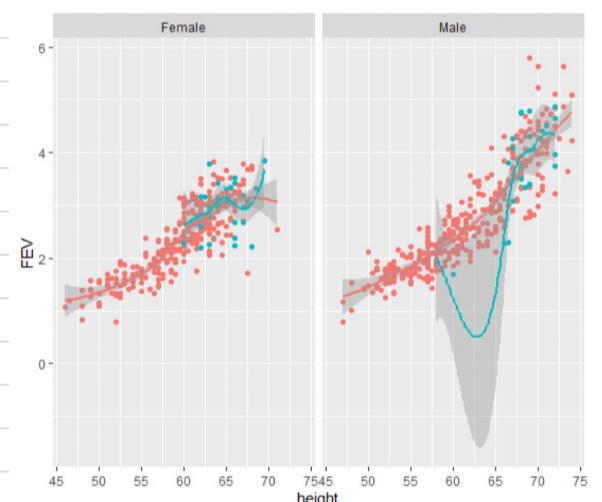


```
library(GGally)
ggpairs(FEV, columns = 1:5, axisLabels = "internal")
```



Does it also depend on gender?

```
p3 <- ggplot(data = FEV, aes(x = height, y = FEV, colour = s
p3 + geom_point() + geom_smooth() + facet_grid(. ~ sex)
```



Visualize all in a plot:

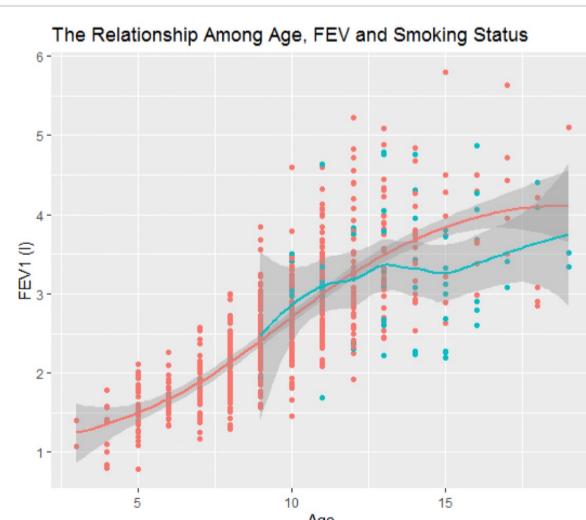
```
library(plotly)
plot_ly(data = FEV, x = ~height, y = ~FEV, z = ~age, color =
width = 800*1.5, height = 600*1.5) %>%
add_markers()
```

No pic in slides offline.

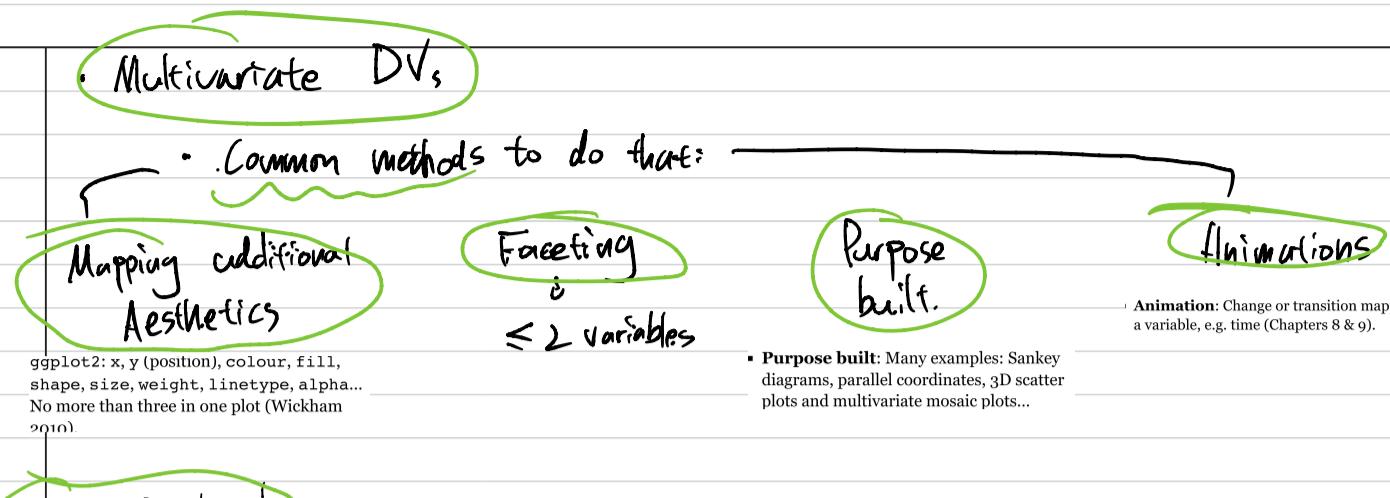
- The world is not bivariate.
- Our world is far more interesting and complex.
- The aim of multivariate data visualisation is to see through this complexity

```
p4 <- ggplot(data = FEV, aes(x = age, y = FEV, colour = smok
p4 + geom_point() + geom_smooth() +
  labs(x = "Age", y = "FEV1 (l)",
       title = "The Relationship Among Age, FEV and Smoking Status
       scale_colour_discrete(name = "")
```

=>



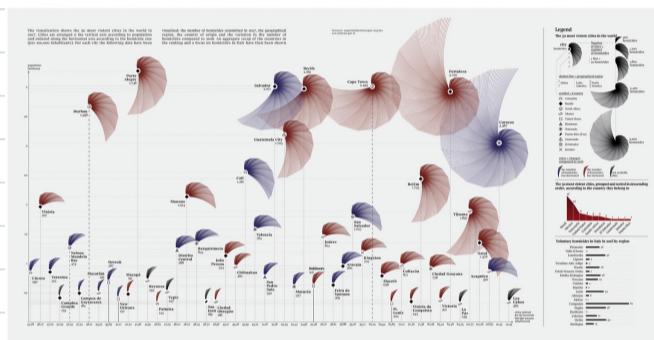
Week 6 - Multivariate Strategies - continue



• Animation: Change or transition mapped to a variable, e.g. time (Chapters 8 & 9).

Multivariate Dilemma

Violent Cities (Fragapane 2018)



1. Clarify Objective

- To educate the audience about the cities with the highest rates of homicide.
 - Rank cities by homicide rate
 - Show population size
 - Show total homicides
 - Show country/region
 - Change (Increased or decreased from 2016) - no data*

2. Data

- Data was sourced from Citizen Council for Public Security and Criminal Justice AC (2018).
- Download a clean copy [here](#).

```
library(readxl)
top50 <- read_excel("../data/top_50_cities_by_homicide_2018")
```

- Original ranks were incorrect

```
top50$Rank = rank(top50$Rate)
```

- Reorder ranks into ascending order

```
library(forcats)
top50$City <- fct_reorder(top50$City, top50$Rank, desc = FALSE)
```

3. Data - collapse countries

- Convert population unit to millions - easier to plot

```
top50$Population <- top50$Population/1000000
```

⇒

- Collapse low tally countries into "other" category

```
table(top50$Country)
```

⇒

```
##          Brazil    Colombia   El Salvador    Guatemala
##          14             2            1             1
##          Jamaica        Mexico       South Africa United States
##          1              15            3              5
```

```
top50$Country_2 <- fct_collapse(top50$Country,
Brazil = "Brazil",
Mexico = "Mexico",
United_States = "United States",
Venezuela = "Venezuela",
South_Africa = "South Africa",
Other = c("Colombia", "El_Salvador", "Guatemala", "Honduras"))
```

```
table(top50$Country_2)
```

```
##          Brazil    Other    Mexico    South_Africa
##          14           7          15            3
##          Venezuela
##          6
```

4. ordered labels levels

```
top50$Country_2 <- factor(top50$Country_2,
levels = c("Mexico",
"Brazil",
"Venezuela",
"Other",
"United_States",
"South_Africa"),
labels = c("Mexico",
"Brazil",
"Venezuela",
"Other*",
"United States",
"South Africa"))
```

```
table(top50$Country_2)
```

```
##          Mexico    Brazil    Venezuela    Other*
##          15           14            6            7
##          South Africa
##          3
```