# Data Wrangling and Data Analysis
# **Data Visualization**

**Daniel L. Oberski & Erik-Jan van Kesteren**

Department of Methodology & Statistics

Utrecht University

Utrecht University

*VERSION: 2023-09-25*

# Coordinators

Erik-Jan

Ayoub
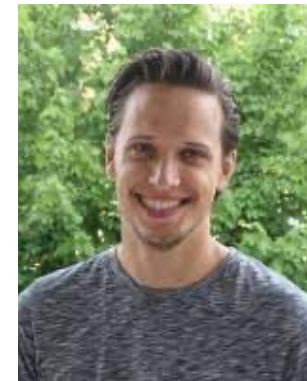
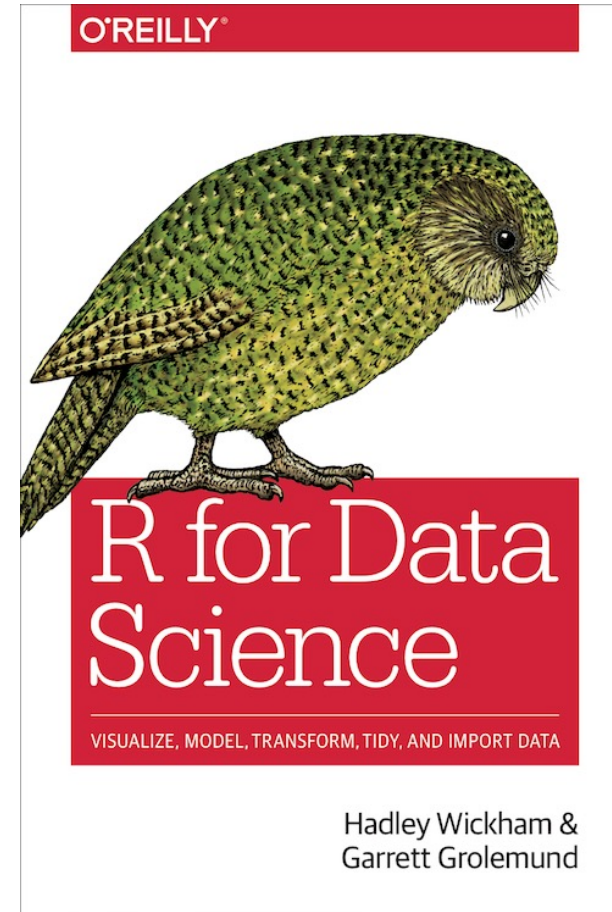# Lecturer

Daniel

# Tutors

Laura

Javier

Mahdi

Anastasia

Jelle

Thom

# This week

1. (Data preparation 2/2)
2. (Cloud computing guest lecture)
3. **Data visualization principles & Grammar of graphics**

Utrecht University

# Reading materials for this week

- Chapters from **R for Data Science (R4DS), open access book at:**

- https://r4ds.had.co.nz

- Today: ch 3 visualization

- Tomorrow: chh 3, 5, 7

# Carte Figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812–1813.

Dressée par M. Minard, Inspecteur Général des Ponts et Chaussées en retraite.

Paris, le 20 Novembre 1869.

Les nombres d'hommes présents sont représentés par les largeurs des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en travers des zones. Le rouge désigne les hommes qui entrent en Russie, le noir ceux qui en sortent. _____ Les renseignements qui ont servi à dresser la carte ont été puisés dans les ouvrages de M.M. Thiers, de Ségur, de Fezensac, de Chambray et le journal inédit de Jacob, pharmacien de l'Armée depuis le 28 Octobre.

Pour mieux faire juger à l'œil la diminution de l'armée, j'ai supposé que les corps du Prince Jérôme et du Maréchal Davoust qui avaient été détachés sur Minsk et Mohilow et ont rejoint vers Orscha et Witebsk, avaient toujours marché avec l'armée.
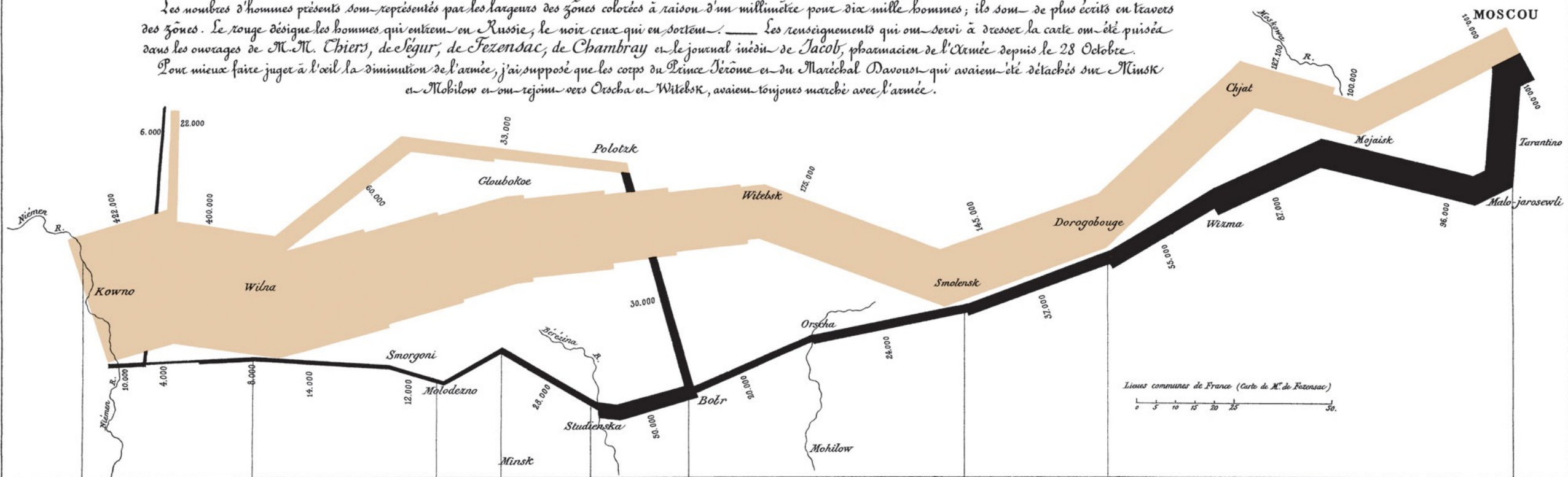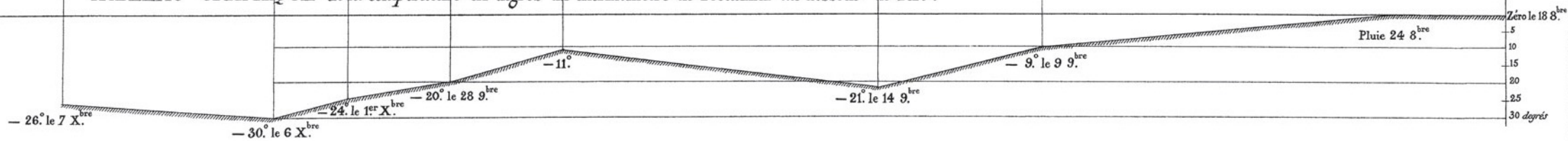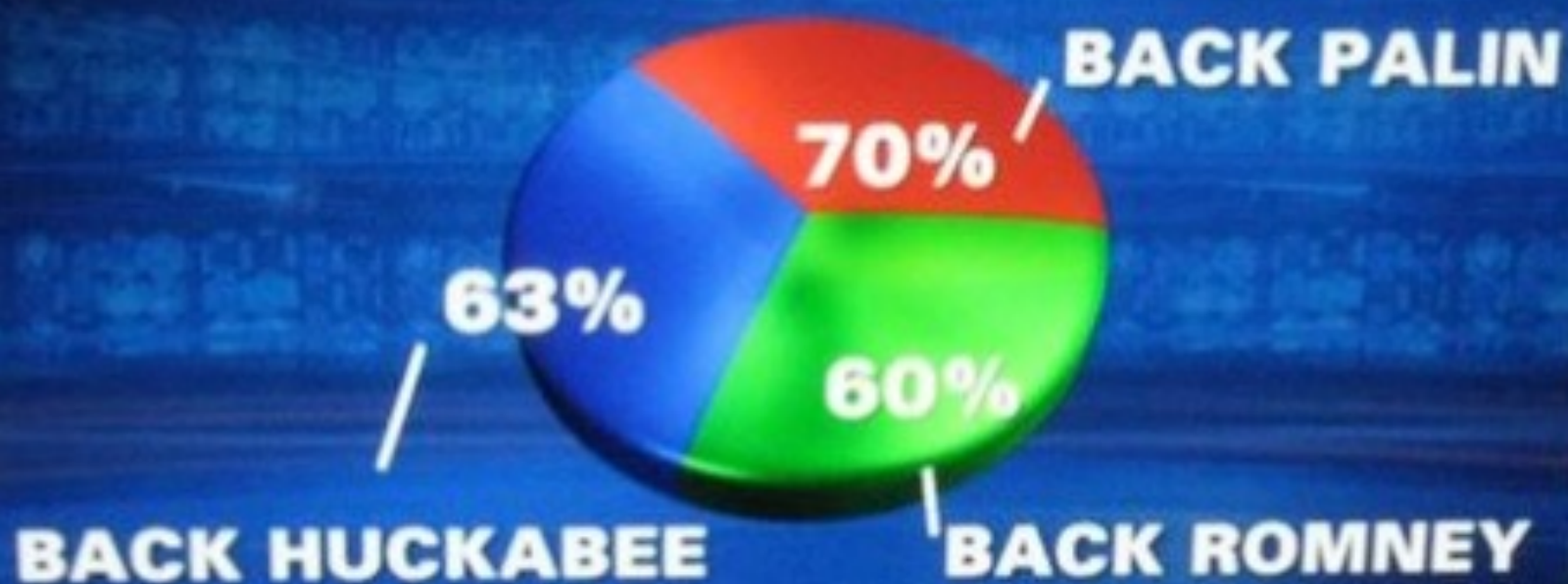
MOSCOU

Moskowa R.

100.000

100.000

Chjat  127.100

100.000

Tarantino

6.000    22.000

Polotzk

Gloubokoe    53.000

Mojaisk

Malo-jarosewli

Niemen R.    +22.000    400.000    60.000    Malo-jarosewli

Witebsk    175.000

87.000    96.000

Wizma    Dorogobouge    145.000    55.000    37.000    24.000

Kowno    Wilna    Smolensk

Orscha

Smorgoni    Bérézina R.

10.000    4.000    8.000    14.000    12.000    Molodezno    28.000    30.000    50.000    Botr    20.000    Mohilow

Niemen R.    Studienska

Minsk

Lieues communes de France (Carte de M. de Fezensac)

0  5  10  15  20  25                    50.

## TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

Les Cosaques passent au galop le Niémen gelé.

Zéro le 18 8.bre

Pluie 24 8.bre

— 11°.

— 9°. le 9 9.bre

— 20°. le 28 9.bre

— 24°. le 1.er X.bre

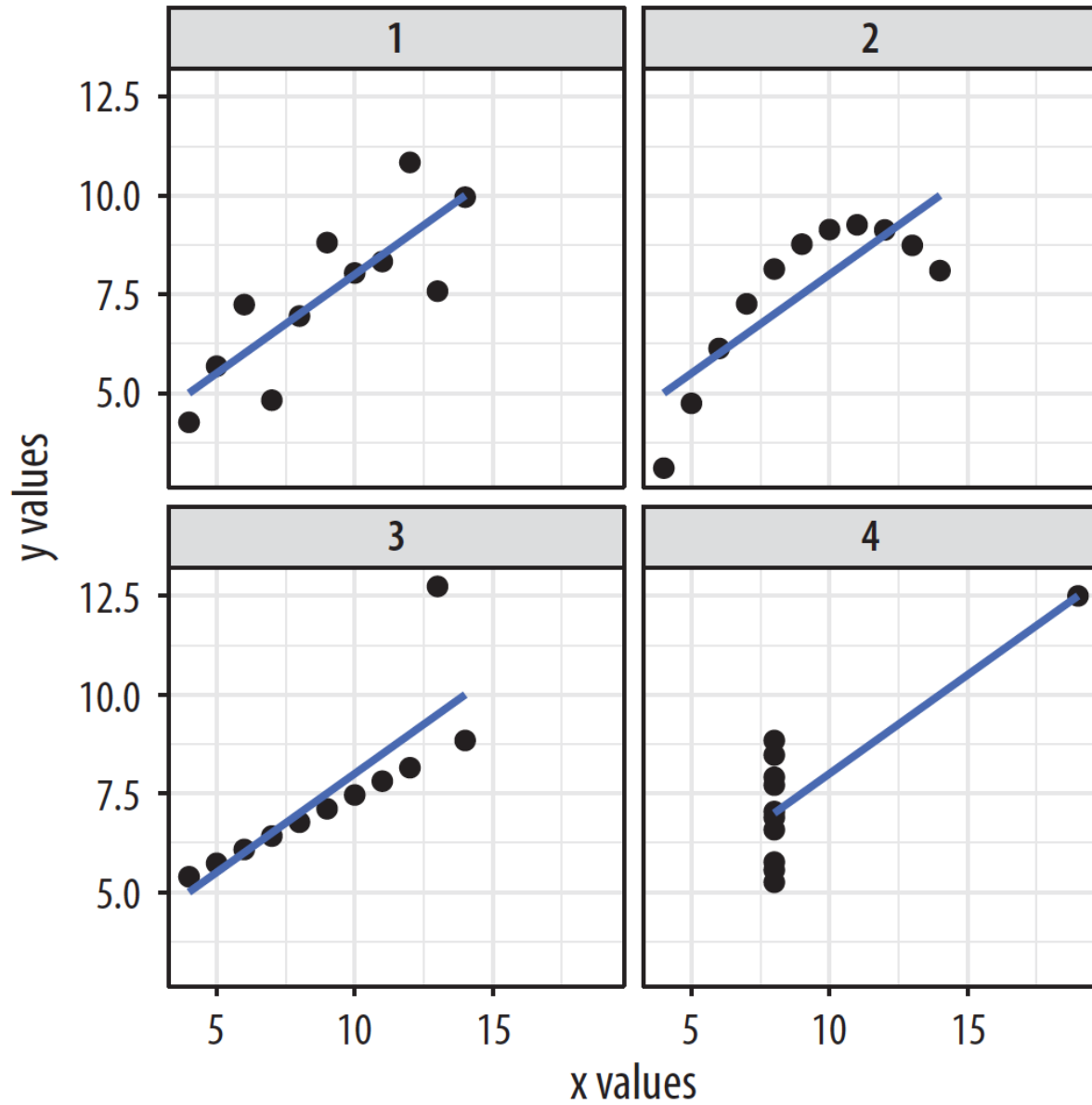— 21°. le 14 9.bre

— 26°. le 7 X.bre

— 30°. le 6 X.bre

30 degrés

# Today: visualization principles

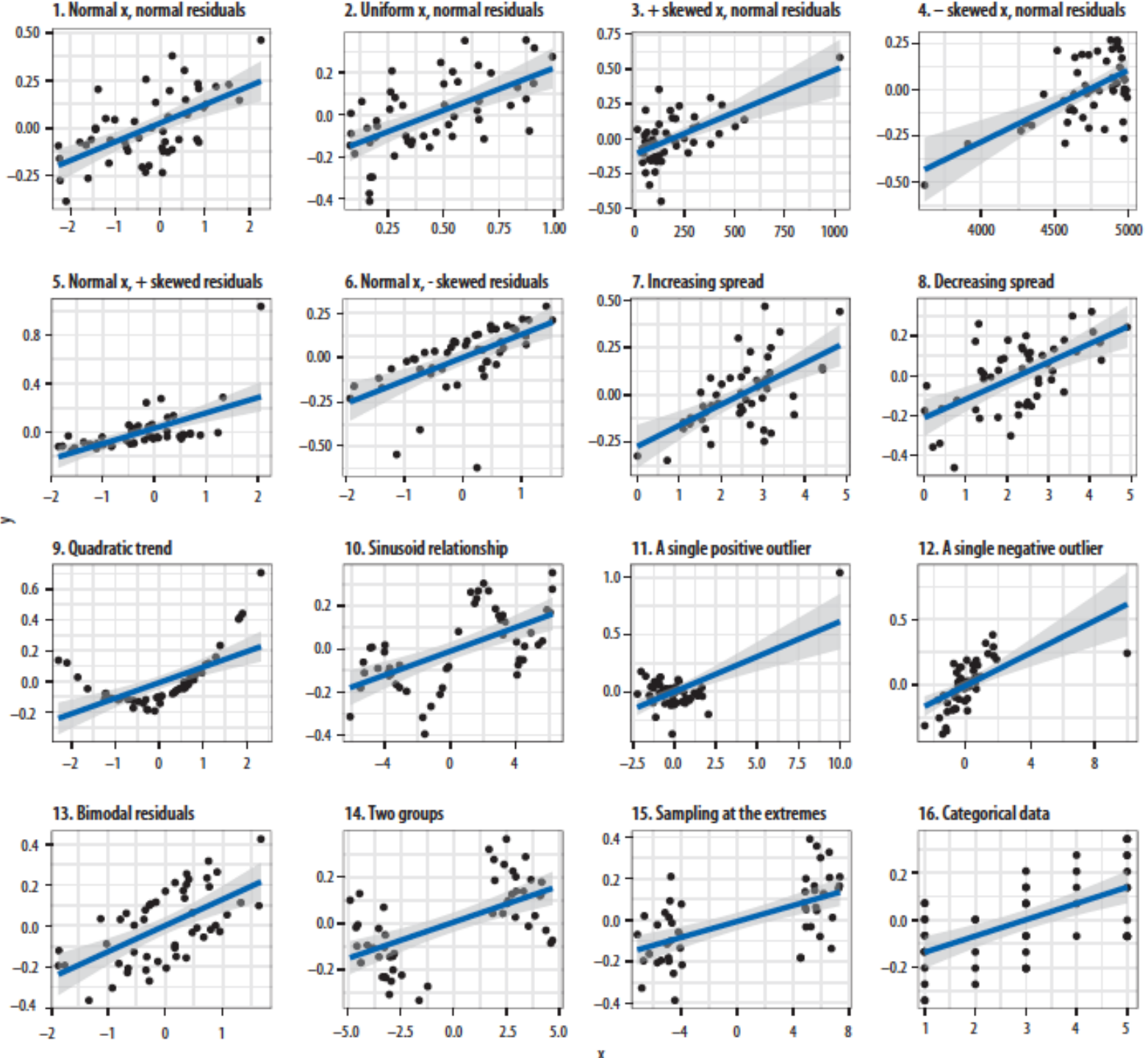# Data visualization

- **For exploration, data analysis ←**
- For communication
- For entertainment

Anscombe's quartet

Source: Healy (2019)

Utrecht University

1. Normal x, normal residuals   2. Uniform x, normal residuals   3. + skewed x, normal residuals   4. – skewed x, normal residuals

5. Normal x, + skewed residuals   6. Normal x, - skewed residuals   7. Increasing spread   8. Decreasing spread

9. Quadratic trend   10. Sinusoid relationship   11. A single positive outlier   12. A single negative outlier

13. Bimodal residuals   14. Two groups   15. Sampling at the extremes   16. Categorical data

Source: Healy (2019)

# Graphics for data analysis

- The **human retina** can transfer around $10^6$ or $10^7$ bits per second to the brain;

- **Reading** transfers about 3 words, so $\sim 10^2$ or $10^3$ bits/s;

- Potentially (!) visualization is about 4 orders of magnitude more powerful.

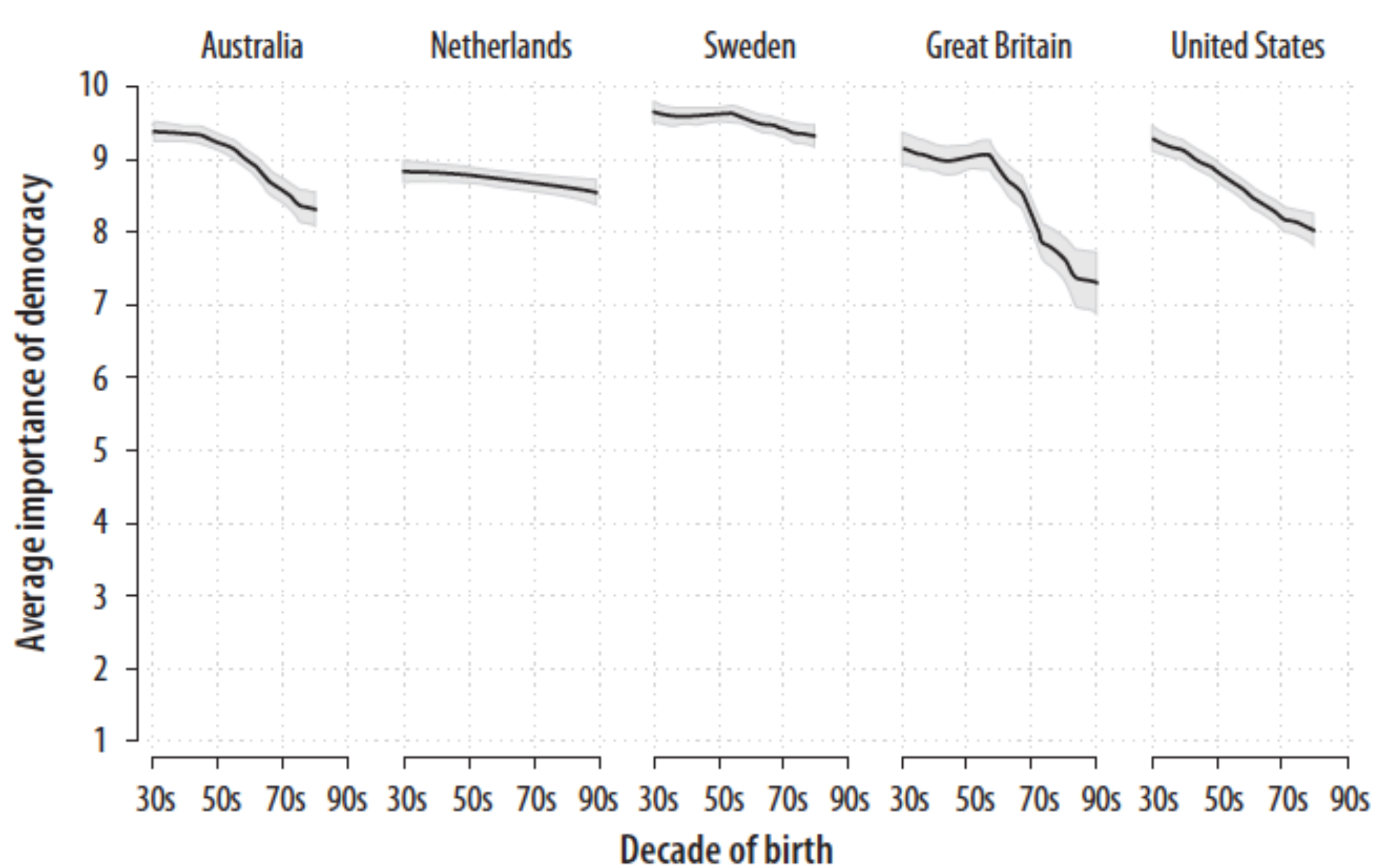**How can we leverage the human visual system to analyze data?**

Utrecht University

# Plotting the right thing

Most common problems:

- (Accidentally) misrepresenting what is being plotted
- Omitting baselines

# Percentage of people who say it is "essential" to live in a democracy



Source: Healy (2019)

Graph by Erik Voeten, based on WVS 5

Source: Healy (2019)

# Case fatality rate of the ongoing COVID-19 pandemic

The Case Fatality Rate (CFR) is the ratio between confirmed deaths and confirmed cases. During an outbreak of a pandemic the CFR is a poor measure of the mortality risk of the disease. We explain this in detail at OurWorldInData.org/Coronavirus



Italy
Jun 25: methodology change

Mexico

Netherlands
Indonesia
Germany
Brazil
United States
Jun 26: probable/earlier deaths added & Jul 1: probable/earlier deaths added

South Africa
Norway
Congo
South Korea
New Zealand
India
Jun 17: earlier deaths added

14%
12%
10%
8%
6%
4%
2%
0%

Feb 7, 2020    Mar 11    Apr 30    Jun 19    Aug 8    Sep 27, 2020

# Death rates have climbed far above historical averages in many countries that have faced Covid-19 outbreaks

Number of deaths per week from all causes, 2020 / vs recent years:

Shading indicates total excess deaths during outbreak

**Good example (FT)**

### UK
67,500 excess deaths (+37%)
25,000
12,500
0
Jan — Sep 11 — Dec
LATEST DATA

### Austria
2,100 (+7%)
2,000
Historical average 1,000
0
Jan — Sep 13 — Dec

### Belgium
4,500
10,700 (+31%)
2,250
0
Jan — Aug 30 — Dec

### Chile
4,500
12,700 (+30%)
2,250
0
Jan — Sep 16 — Dec

### Denmark
500 (+5%)
1,500
750
0
Jan — Aug 12 — Dec

### Ecuador
7,500
35,900 (+104%)
3,750
0
Jan — Sep 16 — Dec

### France
26,600 (+15%)
9,500
9,750
0
Jan — Sep 6 — Dec

### Germany
17,800 (+6%)
21,000
10,500
0
Jan — Aug 23 — Dec

### Iceland
100
No excess deaths
50
0

### Israel
1,500
1,700 (+17%)
750
0

### Italy
50,500 (+38%)
23,500
11,750
0

### Netherlands
5,500
11,300 (+19%)
2,750
0

https://www.ft.com/content/a29
8-5eb7-4633-b89c-cbdf5b38693

OUR SITE'S USERS

SUBSCRIBERS TO *MARTHA STEWART LIVING*

CONSUMERS OF FURRY PORNOGRAPHY

THE BUSINESS IMPLICATIONS ARE CLEAR.

PET PEEVE #208:
GEOGRAPHIC PROFILE MAPS WHICH ARE
BASICALLY JUST POPULATION MAPS

# Data



# Picture



# Brain



Utrecht University

encode → view → decode

Source: **Michael Friendly**, http://euclid.psych.yorku.ca/www/psy6135/#Graphical_Perception

# Making pictures that help analyze data

- We'd like to make, not just any kind of picture or graph, but one that transfers some part of the data to our brain

- How do we make sure that the graphs we make transfer:
  - The right part of the data, and;
  - As much of it as possible?


This is where the **"grammar of graphics"** comes in.


Goal is to **specify how data map to picture**, so the correct type and largest amount possible is transferred

Utrecht University

# Grammar of graphics (Wickham version)

- http://r4ds.had.co.nz/visualize.html
- Map raw data to following elements:
  - Aesthetics (position, shape, color, …)
  - Geometric objects (points, lines, bars, …)
  - Scales (continuous, discrete, …)
  - Facets (small multiples)
- Additionally, can apply:
  - Statistical transformation (identity, binning, median, …)
  - Coordinate system (Cartesian, polar, parallel, …

# Grammar of graphics (Wickham version)

In R, grammar of graphics is implemented in `ggplot`, a function in the `ggplot2` package.

# Example data set: cars

```
mpg
#> # A tibble: 234 × 11
#>   manufacturer model displ  year   cyl       trans   drv   cty   hwy    fl
#>          <chr> <chr> <dbl> <int> <int>       <chr> <chr> <int> <int> <chr>
#> 1         audi    a4   1.8  1999     4    auto(l5)     f    18    29     p
#> 2         audi    a4   1.8  1999     4  manual(m5)     f    21    29     p
#> 3         audi    a4   2.0  2008     4  manual(m6)     f    20    31     p
#> 4         audi    a4   2.0  2008     4    auto(av)     f    21    30     p
#> 5         audi    a4   2.8  1999     6    auto(l5)     f    16    26     p
#> 6         audi    a4   2.8  1999     6  manual(m5)     f    18    26     p
#> # ... with 228 more rows, and 1 more variables: class <chr>
```
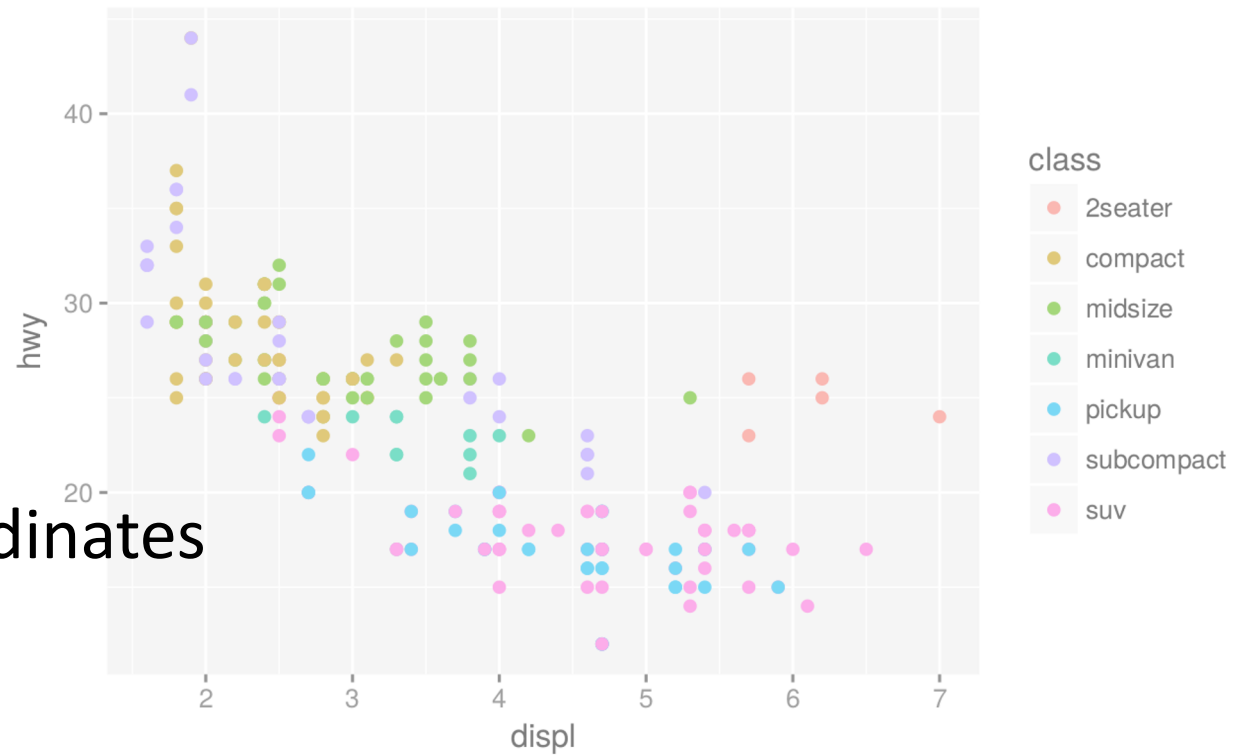
Utrecht University

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy,
                           color = class))
```

- Aesthetics:
  - x-position mapped to engine size
  - y-position mapped to fuel efficiency
  - color mapped to car type

- Geometric objects: points

- Transformation: identity

- Scales: continuous, cartesian coordinates

- No facets



Utrecht University

# Facets

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_wrap(~ class, nrow = 2)
```

# Transformation (stats)



1. **geom_bar()** begins with the **diamonds** data set

2. **geom_bar()** transforms the data with the "count" stat, which returns a data set of cut values and counts.

3. **geom_bar()** uses the transformed data to build the plot. cut is mapped to the x axis, count is mapped to the y axis.

| carat | cut | color | clarity | depth | table | price | x | y | z |
|-------|---------|-------|---------|-------|-------|-------|------|------|------|
| 0.23 | Ideal | E | SI2 | 61.5 | 55 | 326 | 3.95 | 3.98 | 2.43 |
| 0.21 | Premium | E | SI1 | 59.8 | 61 | 326 | 3.89 | 3.84 | 2.31 |
| 0.23 | Good | E | VS1 | 56.9 | 65 | 327 | 4.05 | 4.07 | 2.31 |
| 0.29 | Premium | I | VS2 | 62.4 | 58 | 334 | 4.20 | 4.23 | 2.63 |
| 0.31 | Good | J | SI2 | 63.3 | 58 | 335 | 4.34 | 4.35 | 2.75 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |

stat_count()

| cut | count | prop |
|-----------|-------|------|
| Fair | 1610 | 1 |
| Good | 4906 | 1 |
| Very Good | 12082 | 1 |
| Premium | 13791 | 1 |
| Ideal | 21551 | 1 |

Utrecht University

# What should I choose?

# LES VARIABLES DE L'IMAGE

|  | POINTS | LIGNES | ZONES |
|---|---|---|---|

**XY 2 DIMENSIONS DU PLAN**

**Z**

**TAILLE**

**VALEUR**

# LES VARIABLES DE SÉPARATION DES IMAGES

**GRAIN**

**COULEUR**

**ORIENTATION**

**FORME**

Steven's Psychophysical Power Law: $S = I^N$

Electric Shock (3.5)
Saturation (1.7)
Length (1)
Area (0.7)
Depth (0.67)
Brightness (0.5)

Perceived Sensation

Physical Intensity

*Source*: **Tamara Munzer** (2014). Visualization Analysis and Design.

# Channels: Expressiveness Types and Effectiveness Ranks

→ **Magnitude** Channels: **Ordered** Attributes

Position on common scale

Position on unaligned scale

Length (1D size)

Tilt/angle

Area (2D size)

Depth (3D position)

Color luminance

Color saturation

Curvature

Volume (3D size)

Most

Effectiveness

Least

Same

Same

→ **Identity** Channels: **Categorical** Attributes

Spatial region

Color hue

Motion

Shape

# Color: hue-saturation-brightness (HSB)



Hue Changes

Saturation Changes

Brightness Changes

Cleveland & McGill's Results

Crowdsourced Results

Source: **Tamara Munzer** (2014). Visualization Analysis and Design.

How many 5s in this display?

156132120365841307651037 4627
417312752732759273 2990709742
170370777417952793174 9270973
401974321790937094517 9279417

How many 5s in this display?

156132120365841307651037 4627
41731275273275927322990709742
170370777417952793174 9270973
401974321790937094517 9279417

Numerals differ only in shape, and are high-level symbols
You have to literally scan them all & count the 5s.
The distinction of color is immediate & pre-attentive
You only have to scan & count the 5s.

This is why color is an important visual attribute for a categorical variable in graphs

Source: **Michael Friendly**, http://euclid.psych.yorku.ca/www/psy6135/#Graphical_Perception

# Gestalt principles of relatedness

- **Proximity**: Things that are spatially near to one another seem to be related.
- **Similarity:** Things that look alike seem to be related.
- **Connection**: Things that are visually tied to one another seem to be related.
- **Continuity**: Partially hidden objects are completed into familiar shapes.
- **Closure**: Incomplete shapes are perceived as complete.
- **Figure and ground**: Visual elements are taken to be either in the foreground or in the background.
- **Common fate**: Elements sharing a direction of movement are perceived as a unit.
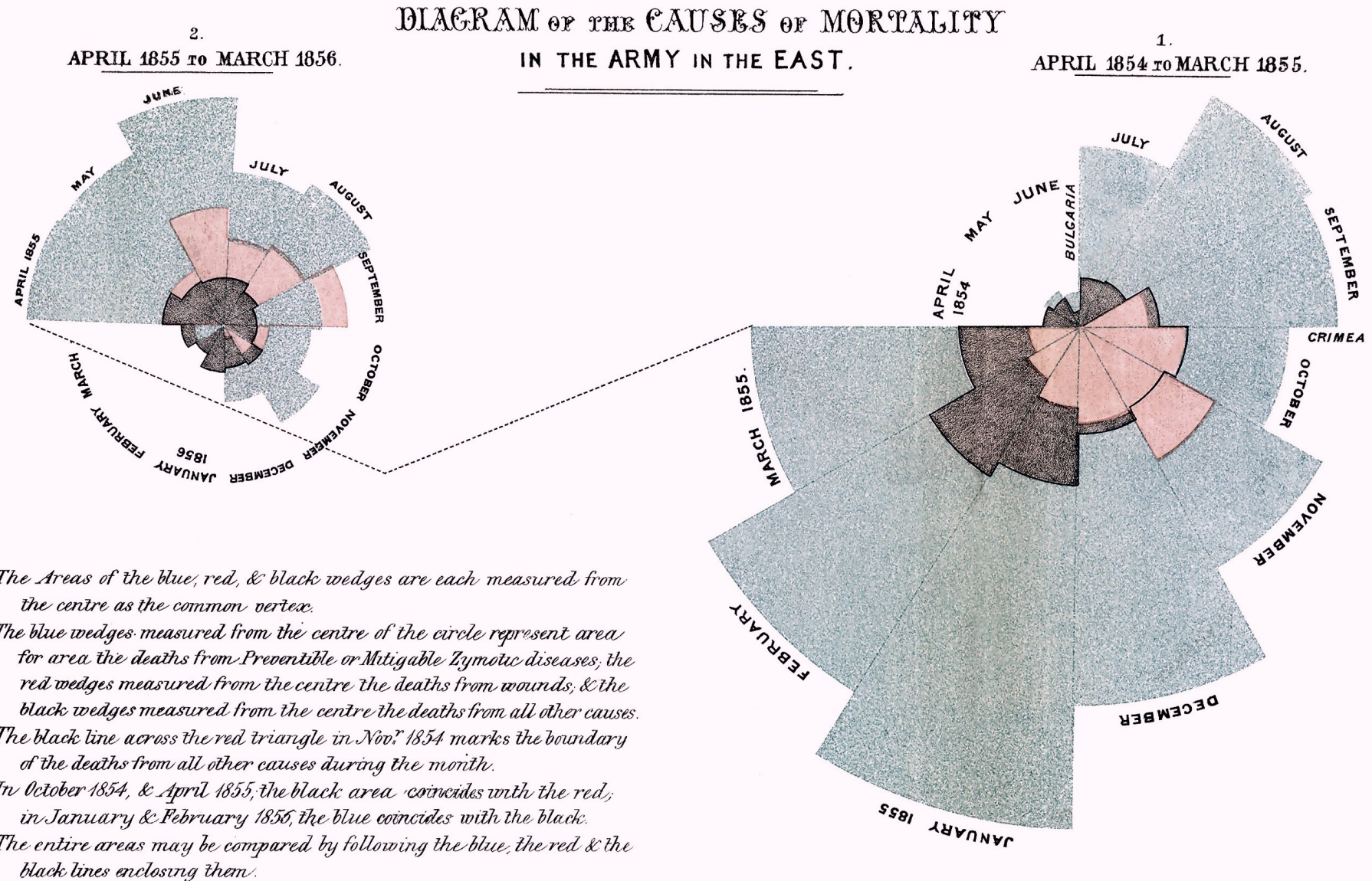
Utrecht University

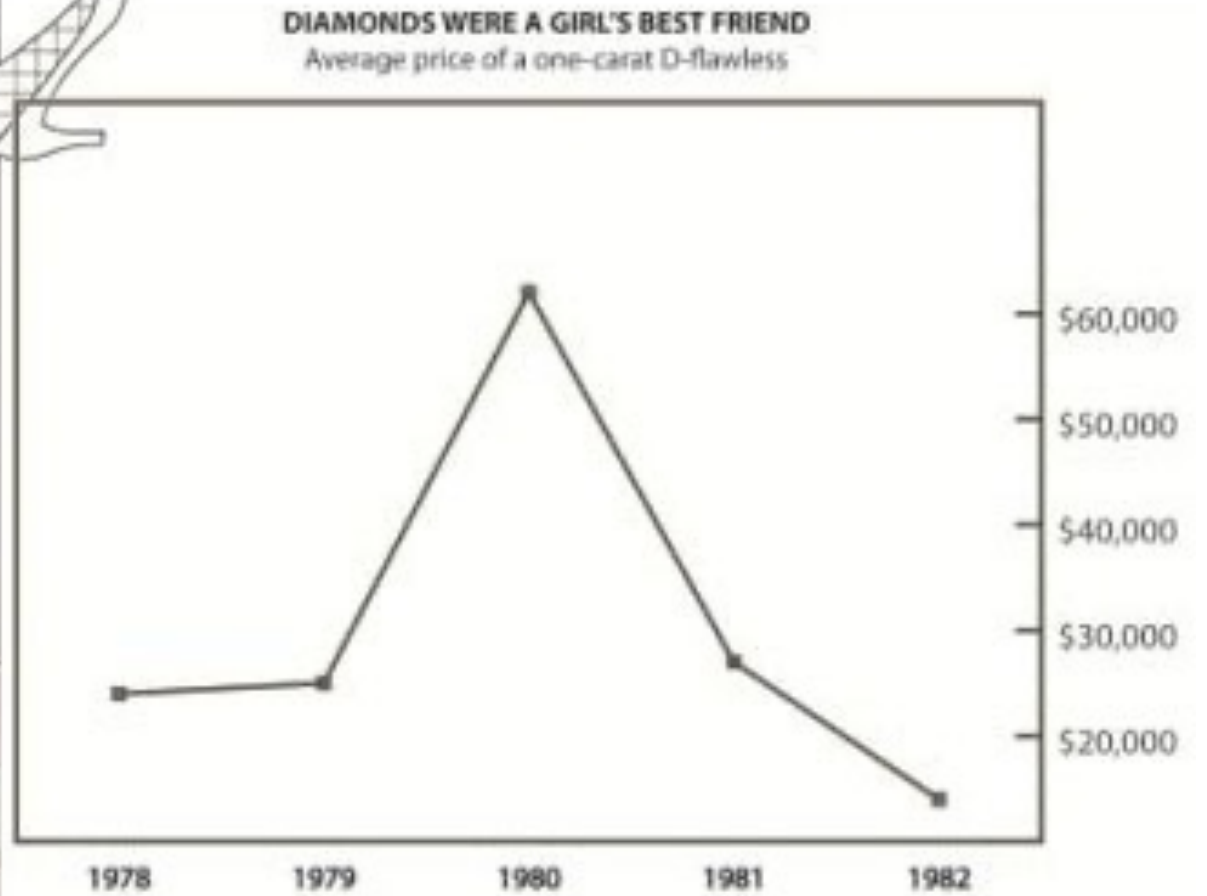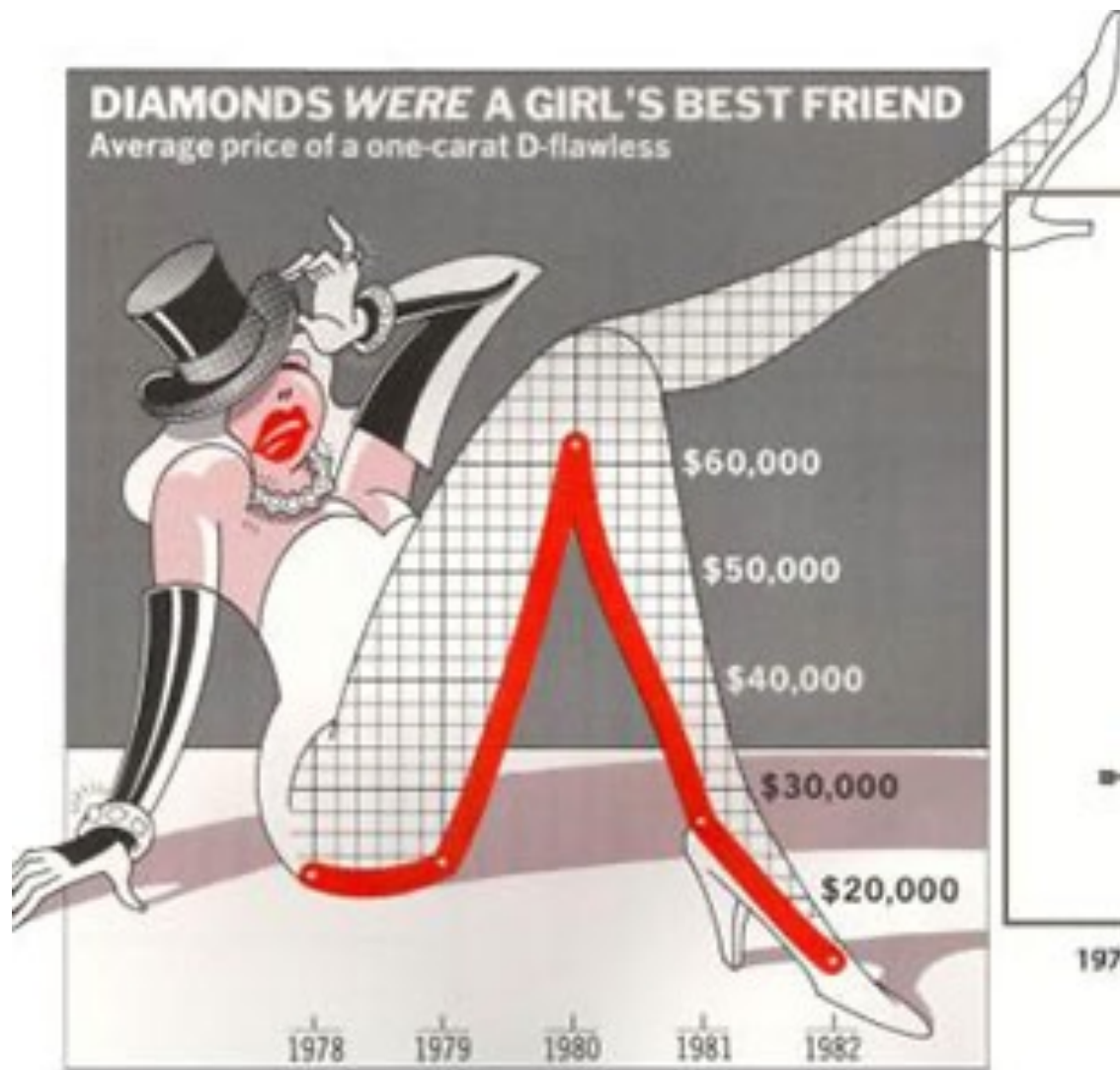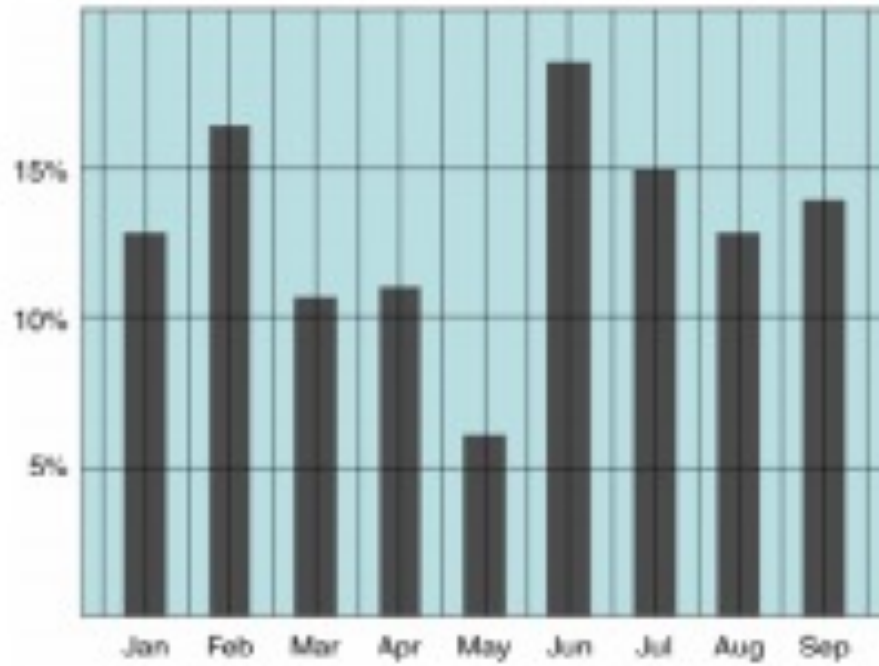Source: Healy (2019)

# Some (distilled) principles from Tufte

- Ask how data maps to perception
- Ask which comparisons you want, guide eye to those
- Maximize data-to-ink ratio
- Present more data (without losing interpretability)
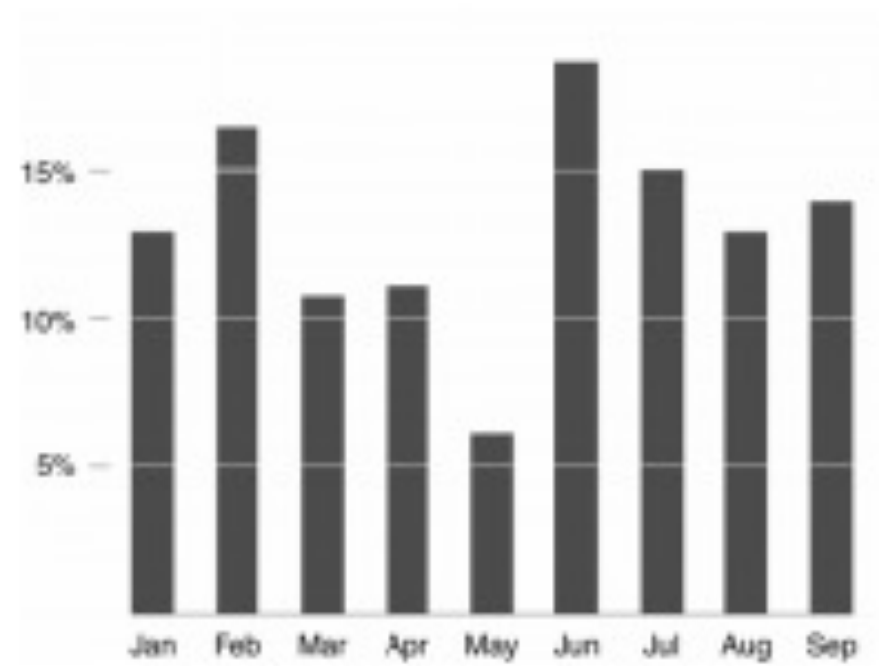- Use levels of detail
- (Remember narrative)

Utrecht University

Mastercam 13%
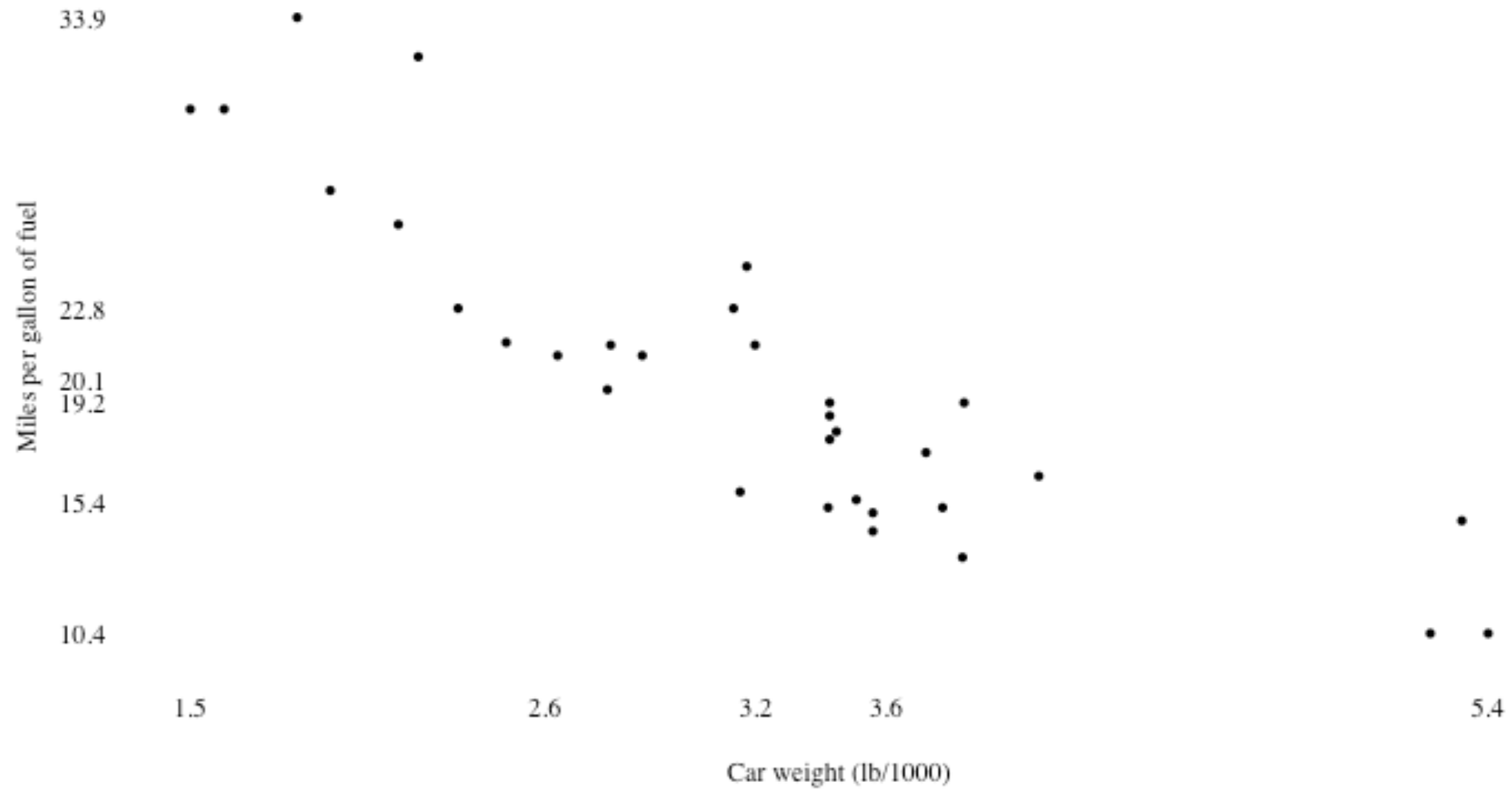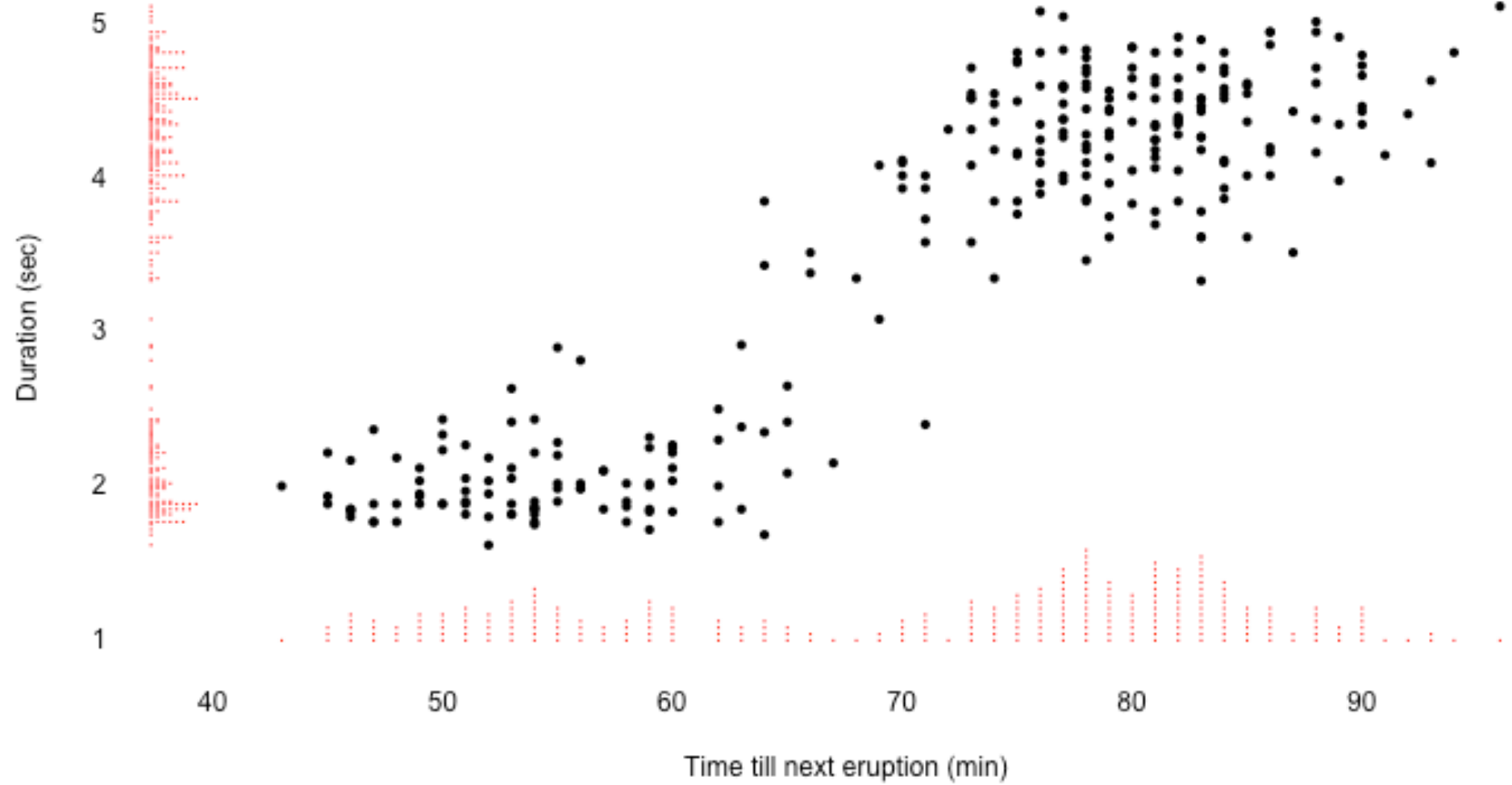
All Other

UGS PLM

PTC

Pathtrace

Delcam

IBM

Software

TekSoft

DP Tech

Planit

# Nightingale Rose / Coxcomb chart

DIAMONDS *WERE* A GIRL'S BEST FRIEND
Average price of a one-carat D-flawless

$60,000
$50,000
$40,000
$30,000
$20,000

1978    1979    1980    1981    1982

DIAMONDS WERE A GIRL'S BEST FRIEND
Average price of a one-carat D-flawless

$60,000
$50,000
$40,000
$30,000
$20,000

1978    1979    1980    1981    1982

Utrecht University
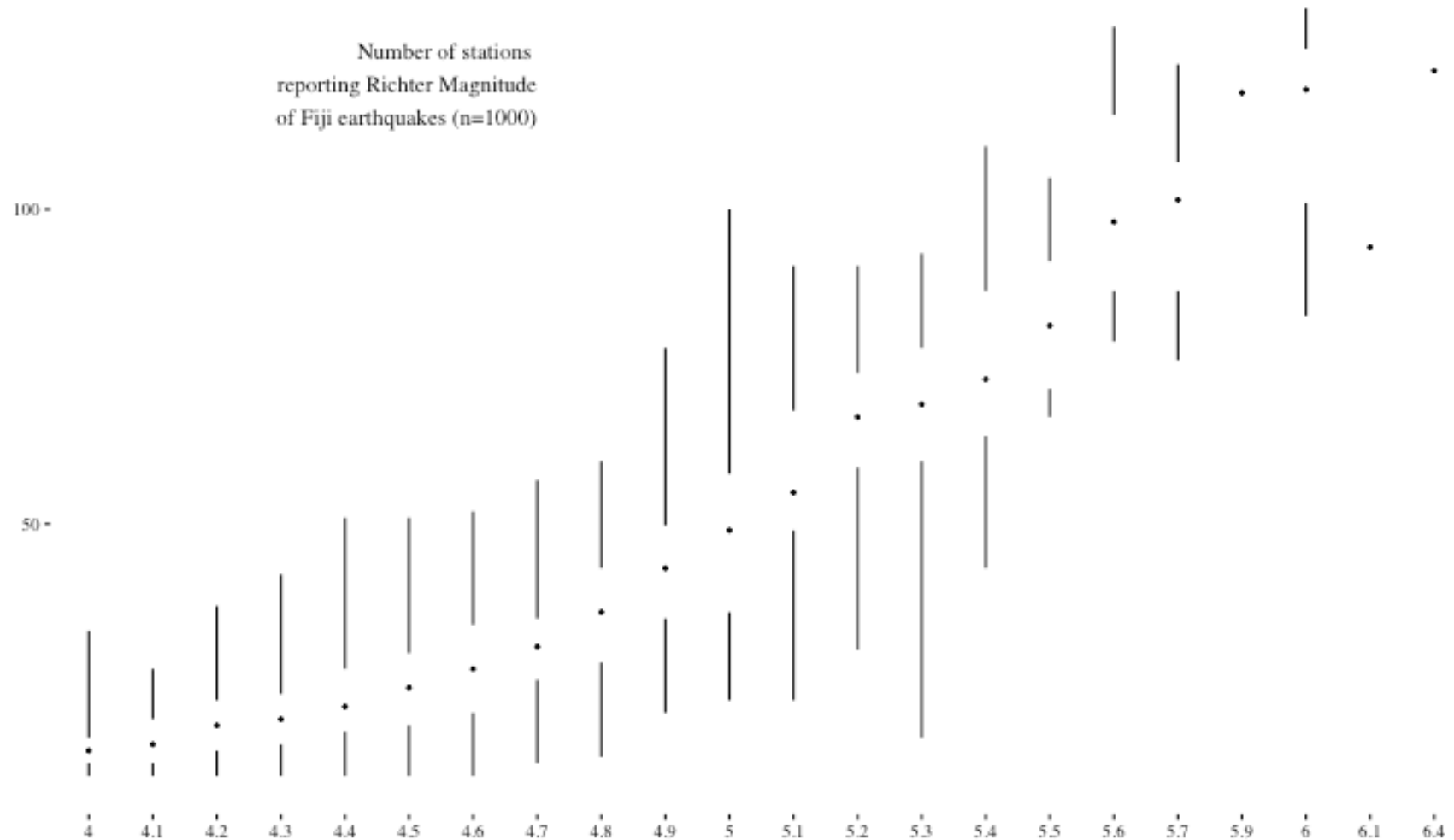
Low Data/Ink          High Data/Ink

Car weight (lb/1000)

```
ggplot(quakes, aes(factor(mag), stations)) +
  theme_tufte() +
  geom_tufteboxplot(outlier.colour = "transparent") +
  theme(axis.title = element_blank())
```



Number of stations
reporting Richter Magnitude
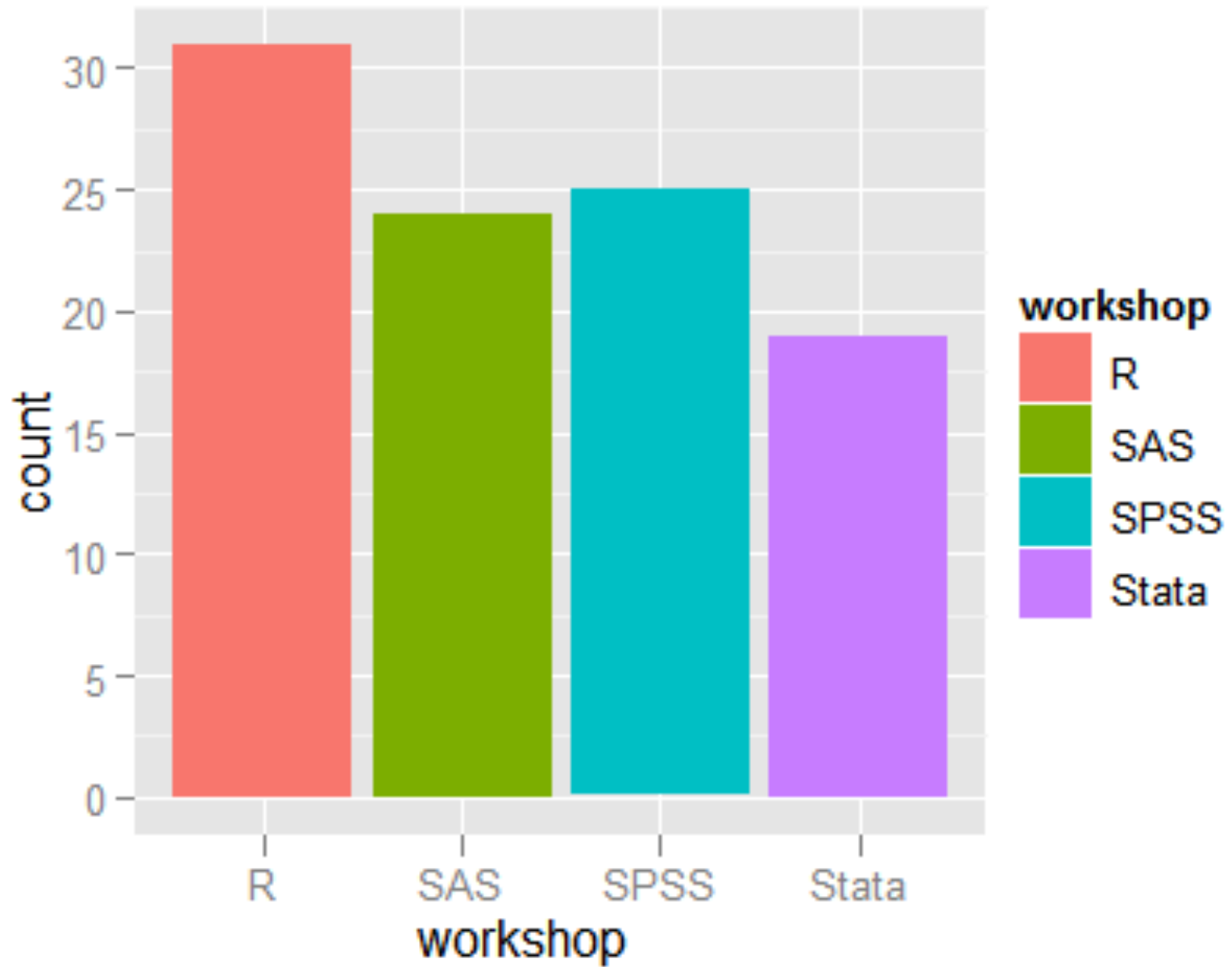of Fiji earthquakes (n=1000)
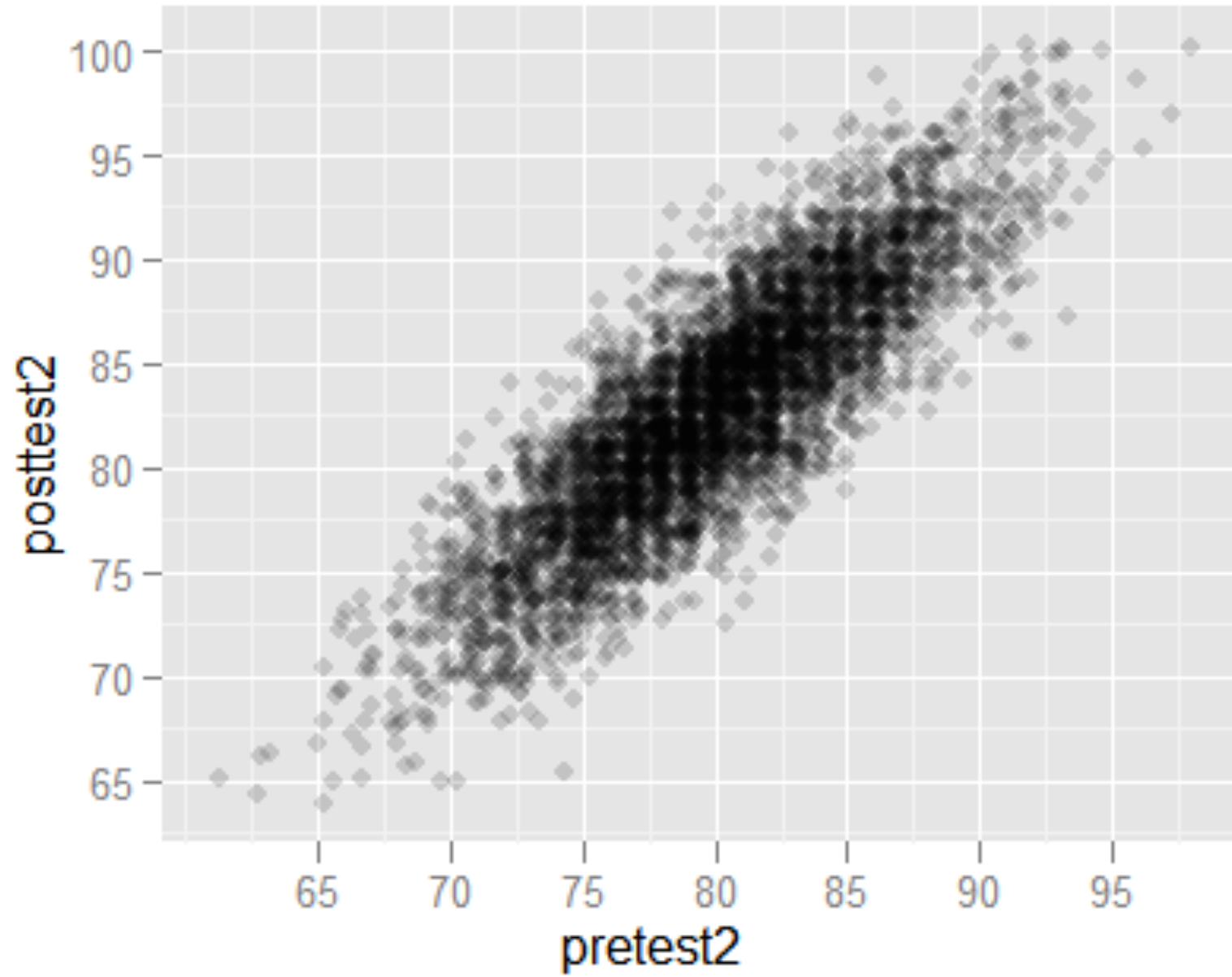
# Tufte wisdom

- Tufte's principles are more oriented to communication and can be taken too far


- Better data/ink → display more information without overload;
- Thinking about perception can help you choose better geoms, aesthetics.
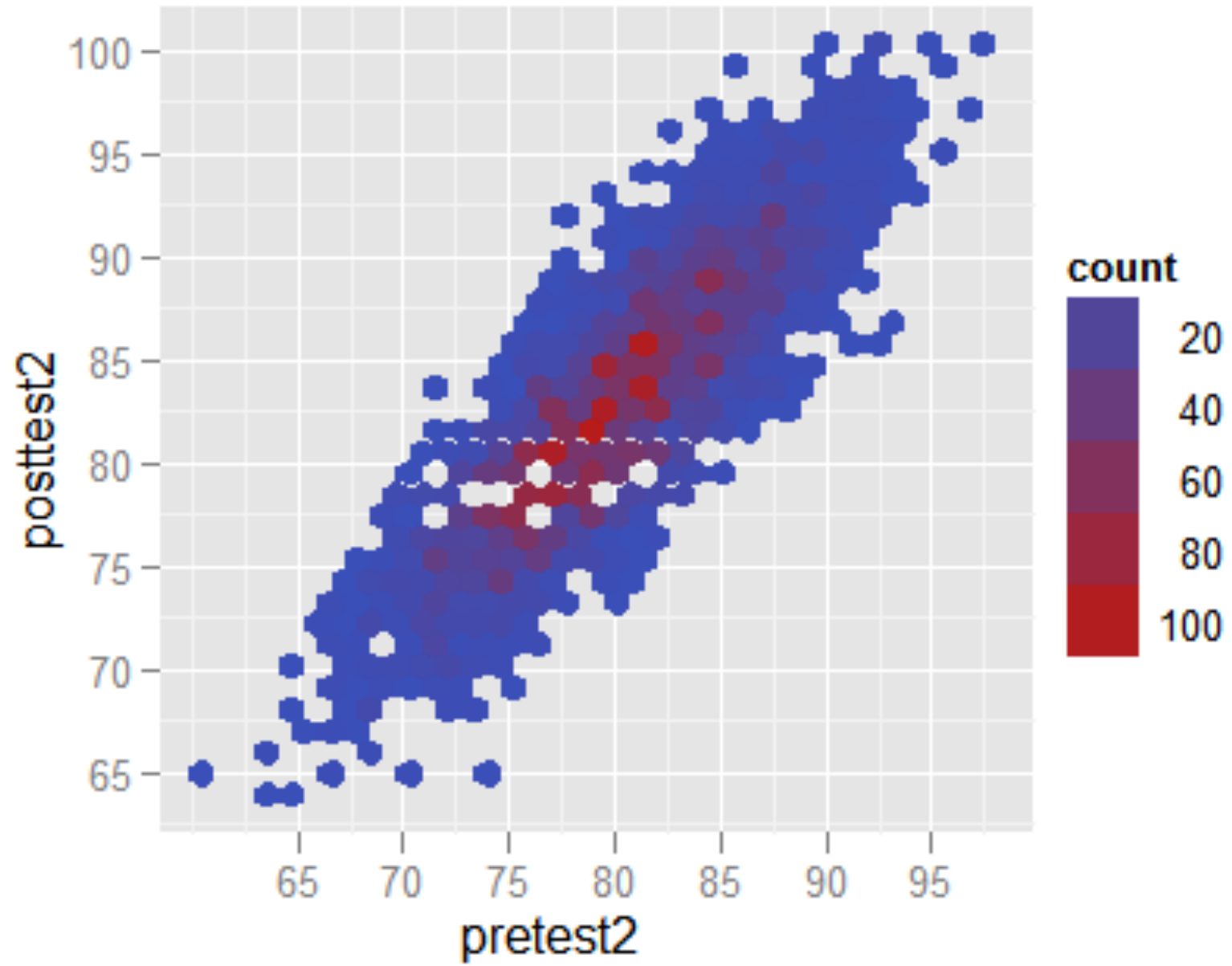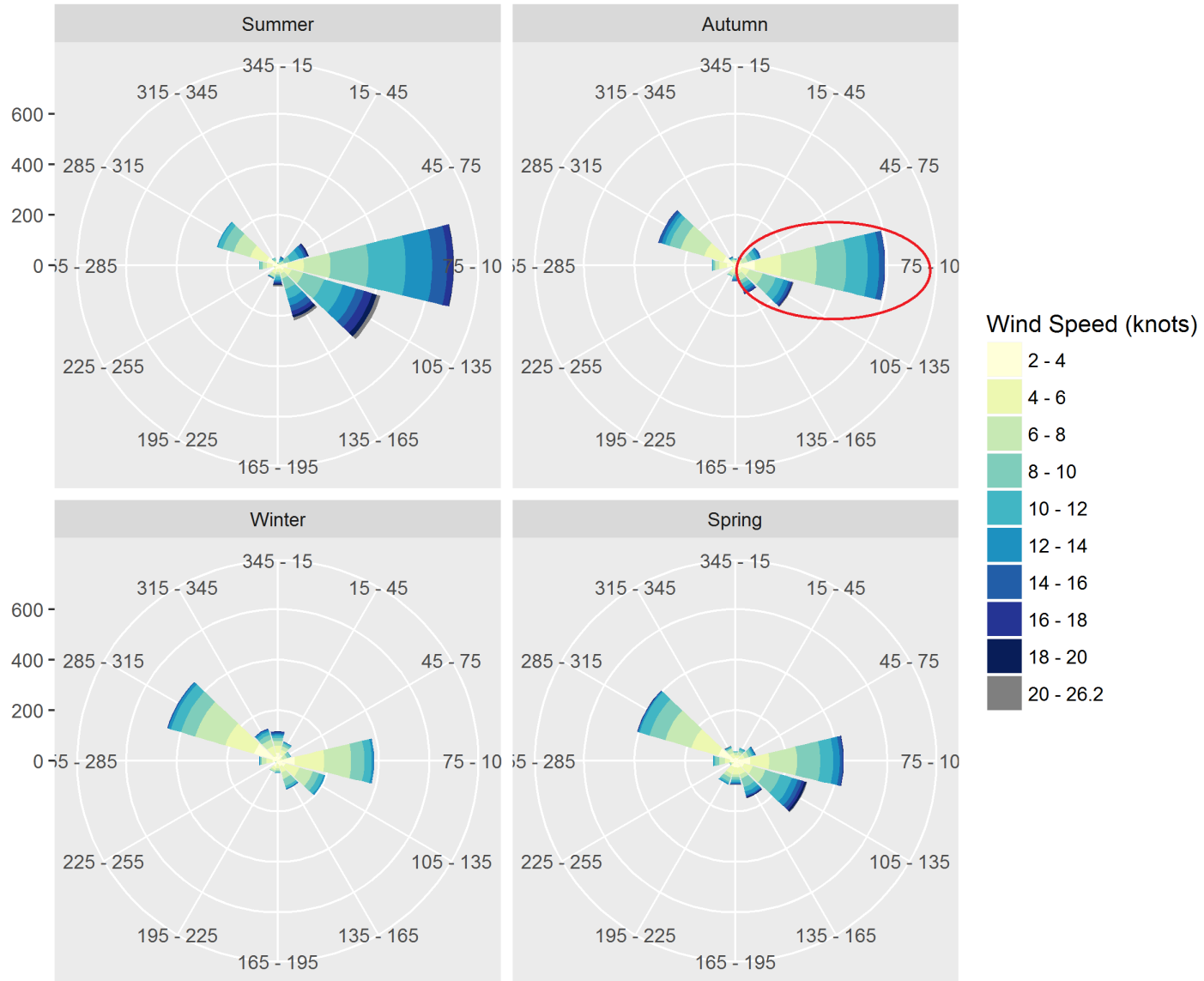
Utrecht University

# Some practice

# Answer these questions:

- Are we plotting the right thing?
- What are: aesthetics, geom, scale, facets, transformation, coordinate system
- How is data/ink?
- Is perception considered optimally?
- Can you think of questions you can't answer from this plot which are in the data?

Utrecht University

# Conclusion

# Conclusion

- Data visualization is data analysis + psychology;
- Sticking to **basic principles** helps:
  - **Map data** to aesthetics, geoms, scales, facets;
  - Perception research guides choices;
  - **Which comparisons** do I want?
  - Maximize **data-ink** (within reason).
- Some standard recipes (e.g. "barplot", "histogram", "line graph"), but pros do not need the cookbook…
- Don't believe everything you hear ("do's and don'ts")

Utrecht University