

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.

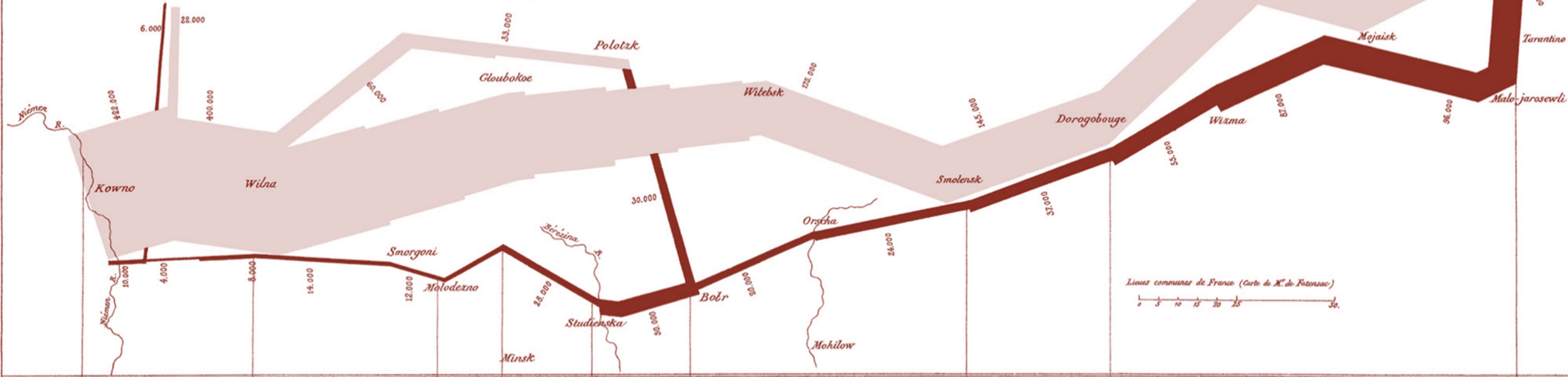


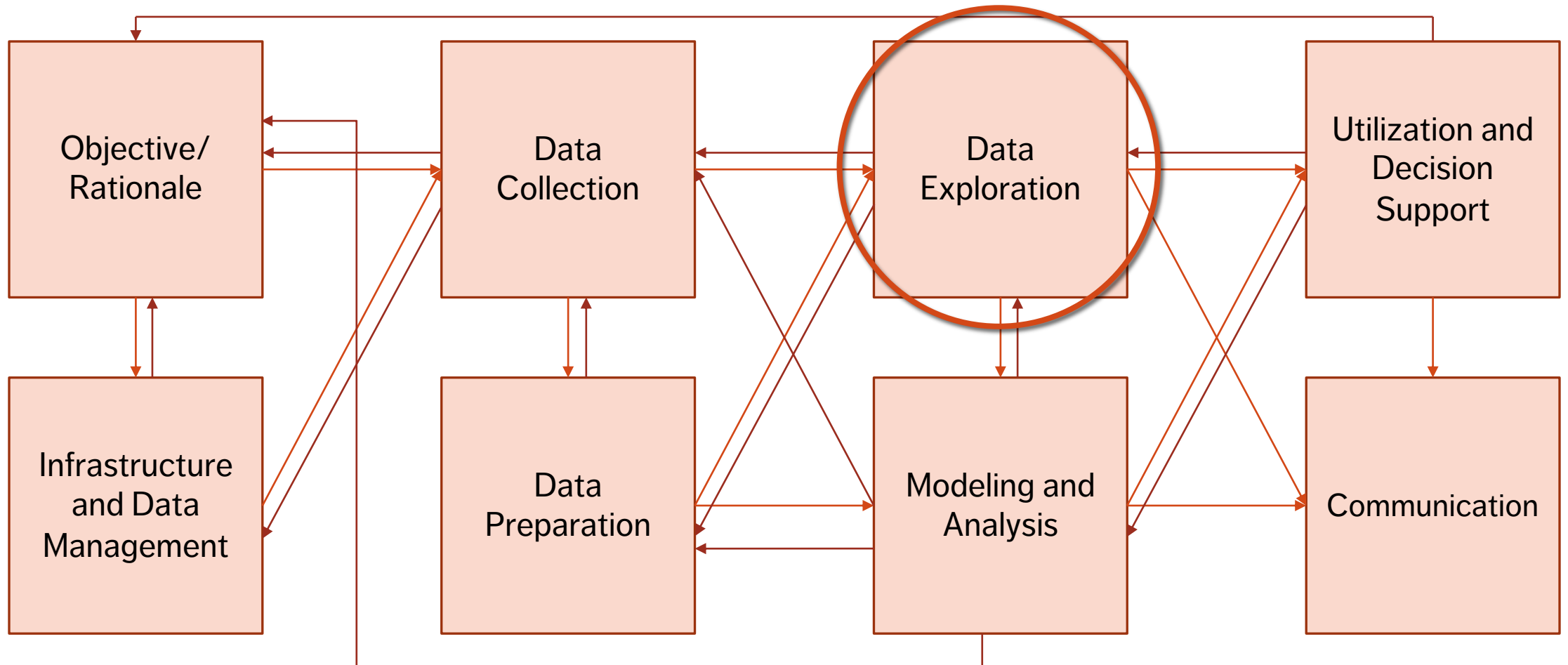
TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

Les Cosaques passent au galop le Niémen, gelé.



# 1. Exploratory Data Analysis

# The (Messy) Analysis Process



# Some Basic Questions

---

What **system** does your data represent – **objects, attributes, relationships**?

How does it represent this system – i.e., the **data model**?

Who made this dataset? When? For what purpose?

Assuming a flat file – what do the rows and columns represent?

Do you have enough information (e.g., **metadata**) to answer these questions?

Where can you find more information?

# Non-Visualization Summaries

---

	CL	N03	NH4
Min.	: 0.222	Min. : 0.000	Min. : 5.00
1st Qu.:	10.994	1st Qu.: 1.147	1st Qu.: 37.86
Median :	32.470	Median : 2.356	Median : 107.36
Mean :	42.517	Mean : 3.121	Mean : 471.73
3rd Qu.:	57.750	3rd Qu.: 4.147	3rd Qu.: 244.90
Max. :	391.500	Max. : 45.650	Max. : 24064.00
NA's :	16	NA's : 2	NA's : 2

```

season
Length:340
Class :character      autumn spring summer winter
Mode  :character      80      84      86      90

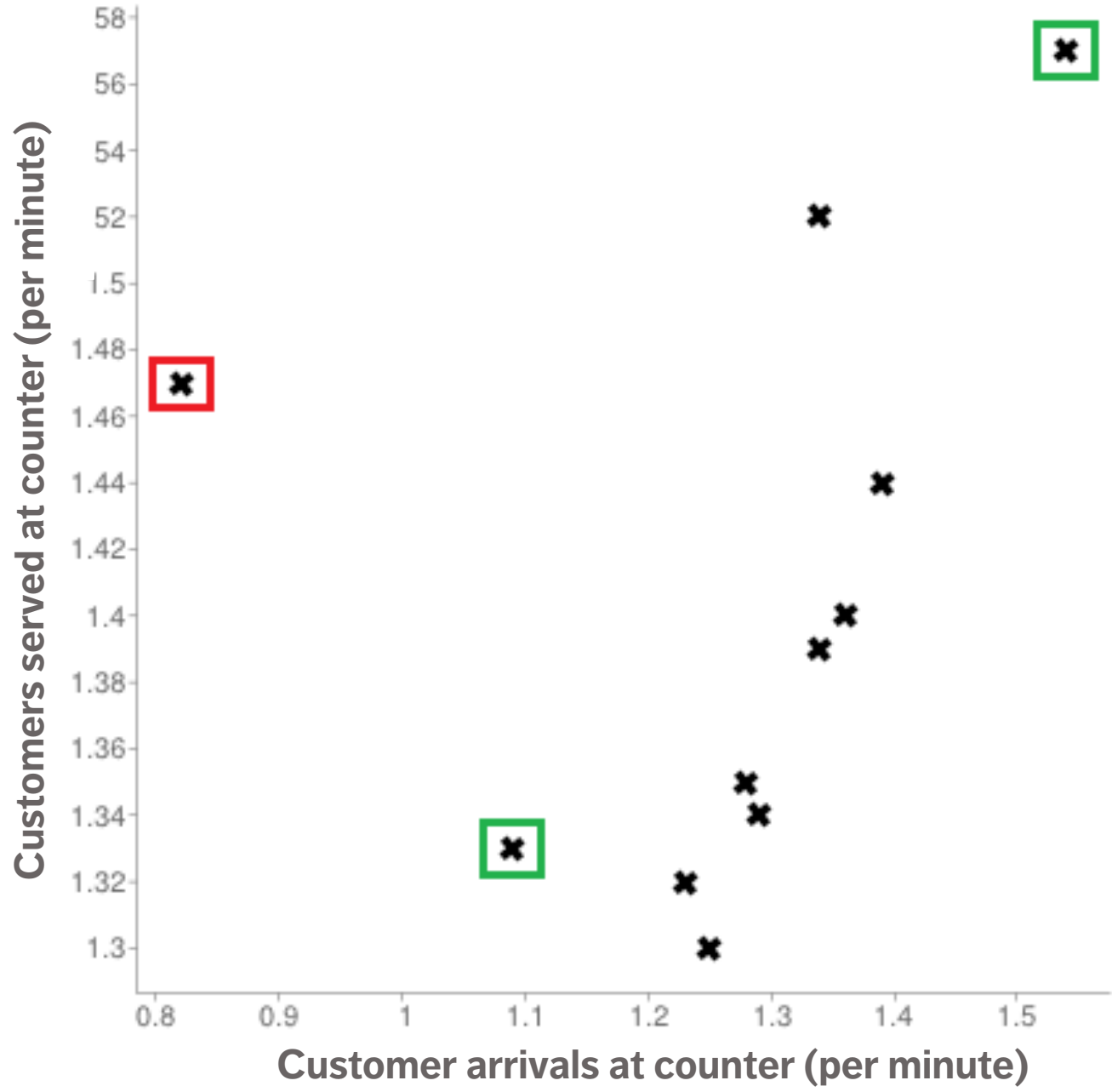
```

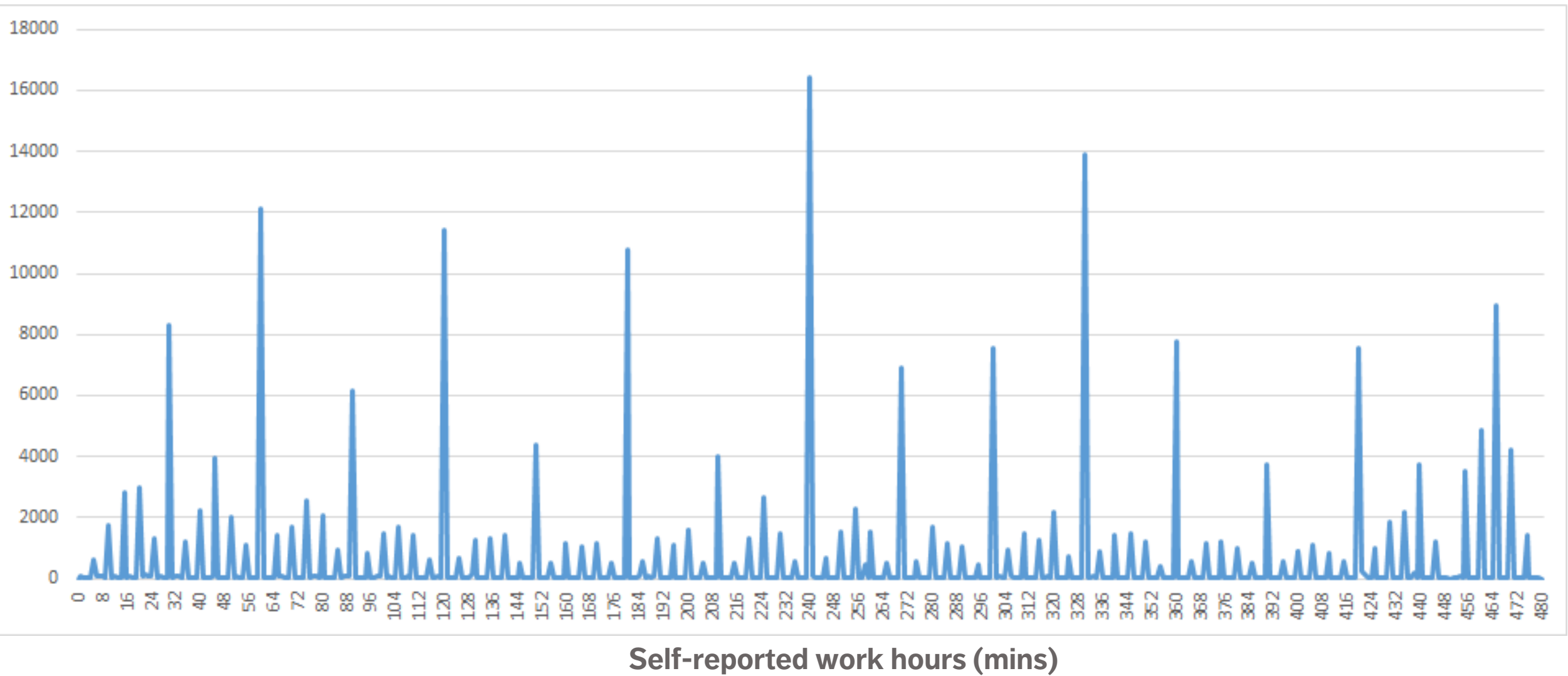
# Pre-Analysis Use

---

Data visualization can be used to set the stage for analysis:

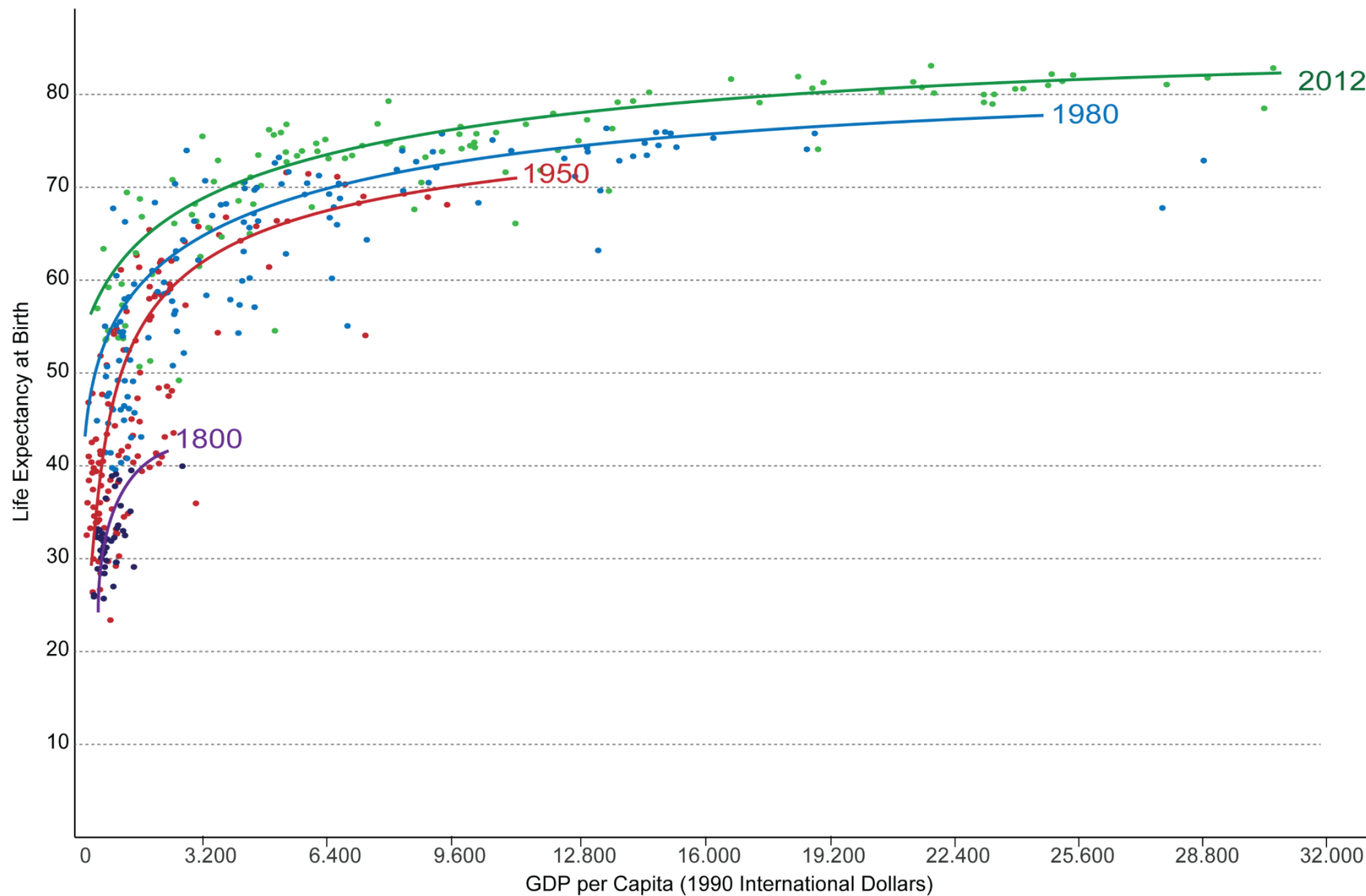
- **detecting anomalous entries**  
invalid entries, missing values, outliers
- **shaping the data transformations**  
binning, standardization, Box-Cox transformations, PCA-like transformations
- **getting a sense for the data**  
data analysis as an art form, exploratory analysis
- **identifying hidden data structure**  
clustering, associations, patterns informing the next stage of analysis





## Life Expectancy vs. GDP per Capita from 1800 to 2012 – by Max Roser

GDP per capita is measured in International Dollars. This is a currency that would buy a comparable amount of goods and services a U.S. dollar would buy in the United States in 1990. Therefore incomes are comparable across countries and across time.

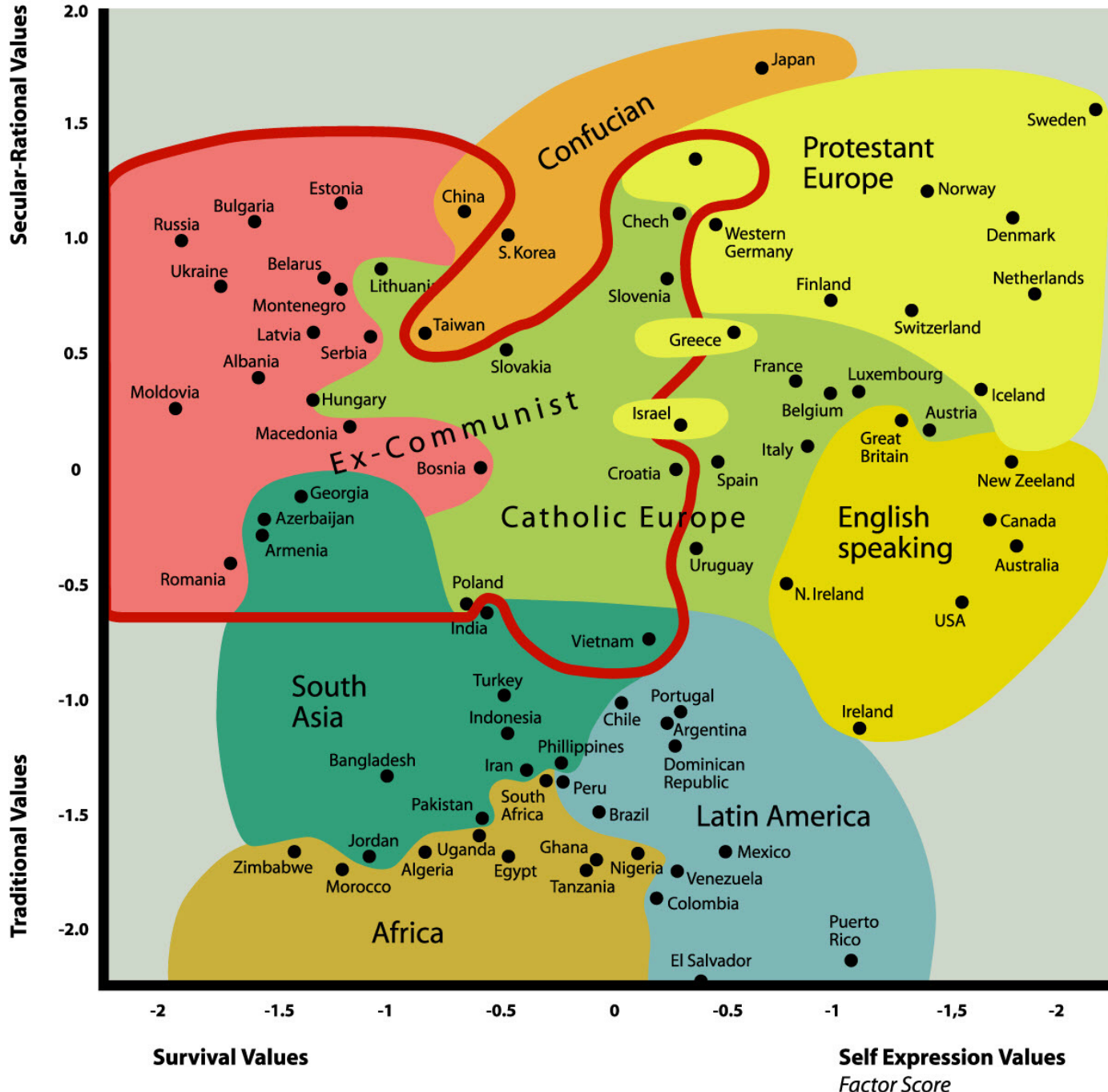


This graph displays the correlation between life expectancy and GDP per capita.

Countries with higher GDP have a higher life expectancy, in general.

The relationship seems to follow a logarithmic trend: the unit increase in life expectancy per unit increase in GDP decreases as GDP per capita increases.





**Traditional values**  
 importance of religion, parent-child ties, deference to authority and traditional family values.

**Secular-rational values**  
 less emphasis on religion, traditional family values and authority.

**Survival values**  
 emphasis on economic and physical security.

**Self-expression values**  
 high priority to environmental protection, growing tolerance of foreigners, gays and lesbians and gender equality

# Workhorse Data Exploration Charts

---

Text and Tables

Rug Charts/Number Lines

Histograms/Bar Charts

Boxplots

Line Graphs

Scatterplots

# Line Chart/Rug Chart

---

Gaps in the number line: **absence** of those numeric values in the data.

Remember: this is (possibly) different from the order that values appear in the dataset – since it is a number line, it shows where the values fall numerically.

If some values are identical, they lie on top of each other (use **jitter**?).



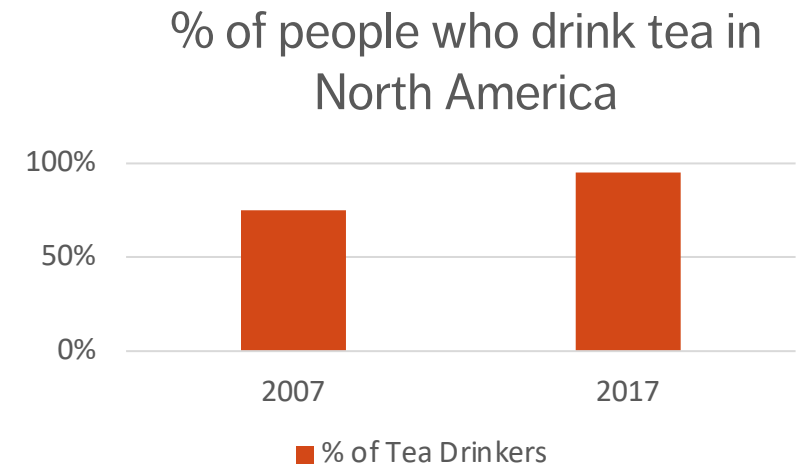
# Simple Text

---

One or two numbers to focus on.

Good at “setting the scene”.

Draws focus to an area of the report.



**95%** of the population  
drinks tea today compared to  
**75%** in 2007

# Tables

---

Tables interact with our **verbal** system, which means we **read** them:

- used to compare values
- audiences will look for their rows

Table design needs to **blend** into background

- the data should stand out, not the borders
- dense table/data: use alternating row colour

Name	Last Year	This Year
Bob	20	30
Fred	30	40
George	10	15

Name	Last Year	This Year
Bob	20	30
Fred	30	40
George	10	15

# Table Heatmaps

Leverage colour to convey magnitude

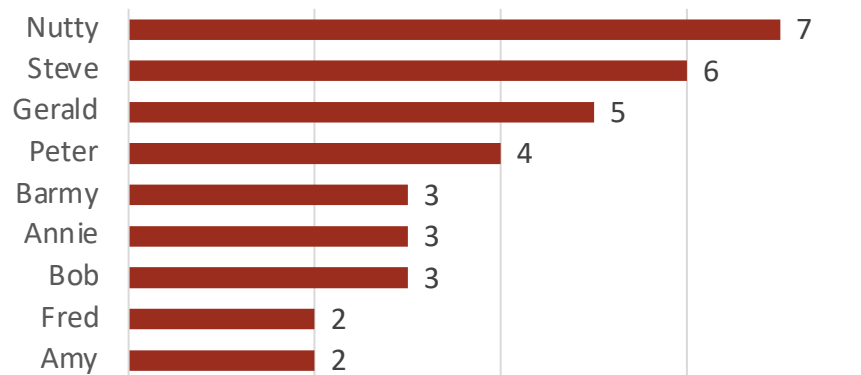
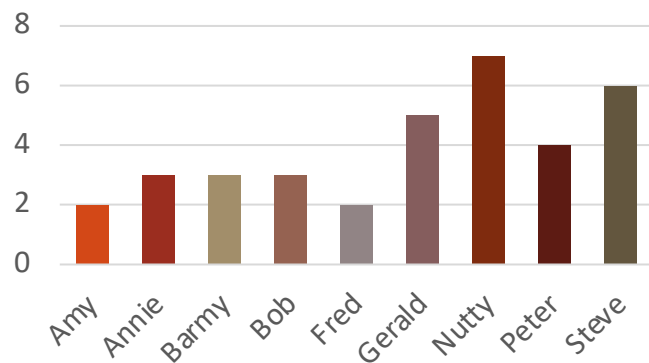
- use **single colour saturation** rather than differentiation (different colours)
- with a legend (white = low, blue = high), numbers can be removed without altering the message

	Last Year	This Year	Next Year	Optimum
George	20	20	20	20
Peter	40	35	30	25
John	10	10	5	5
Sandra	25	30	35	40

	Last Year	This Year	Next Year	Optimum
George	20	20	20	20
Peter	40	35	30	25
John	10	10	5	5
Sandra	25	30	35	40

	Last Year	This Year	Next Year	Optimum
George				
Peter				
John				
Sandra				

# Bar Charts



Very versatile and useful.

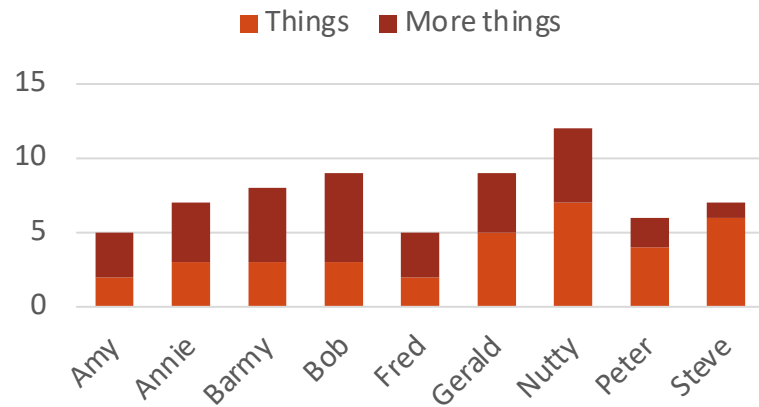
ALWAYS (?) have a zero baseline.

Use graph axis OR data labels. Axis for broad statements, data labels for more detail.

Horizontal charts are apparently **easier to read** (according to many studies).

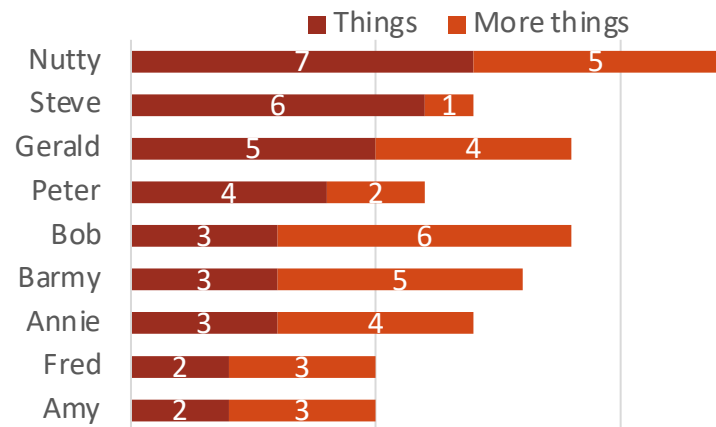
Think about the ordering of categories.

# Stacked Bar Charts



Designed for **comparing totals**, but can quickly become **overwhelming**.

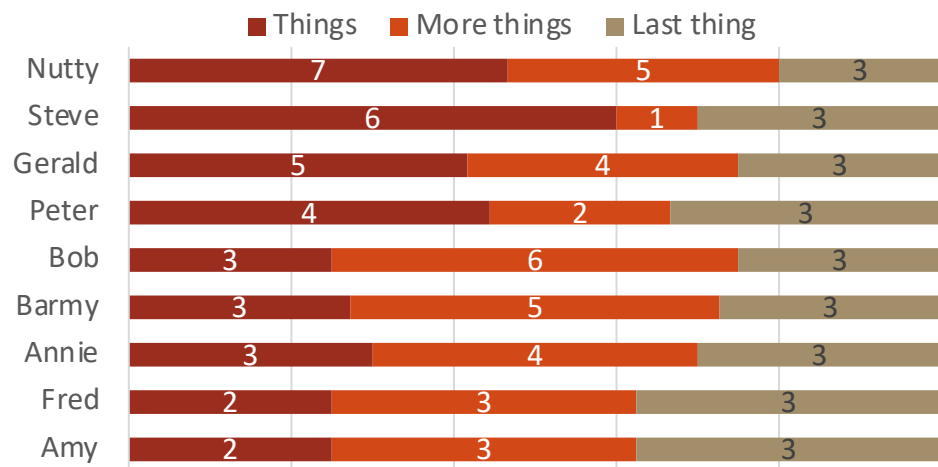
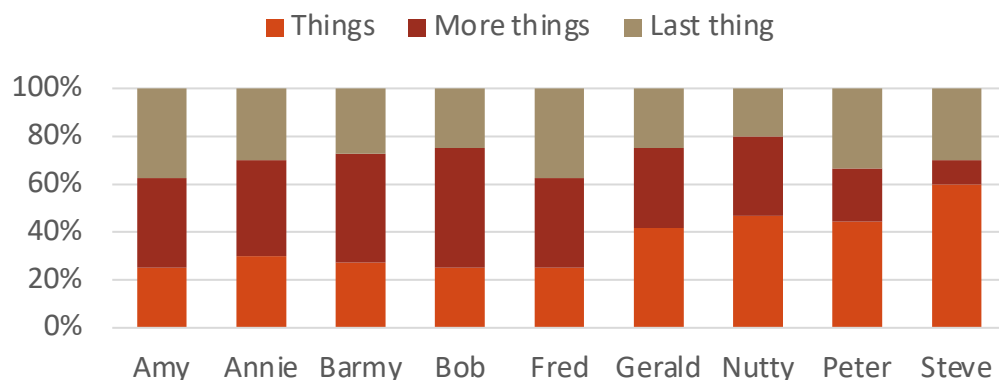
Hard to sort / order.



Filtering is complicated in Power BI (what do you click on & how the chart responds when filter is clicked on?)



# 100% Bar Charts



Work well for visualizing **proportions** of a whole on a scale from negative to positive.

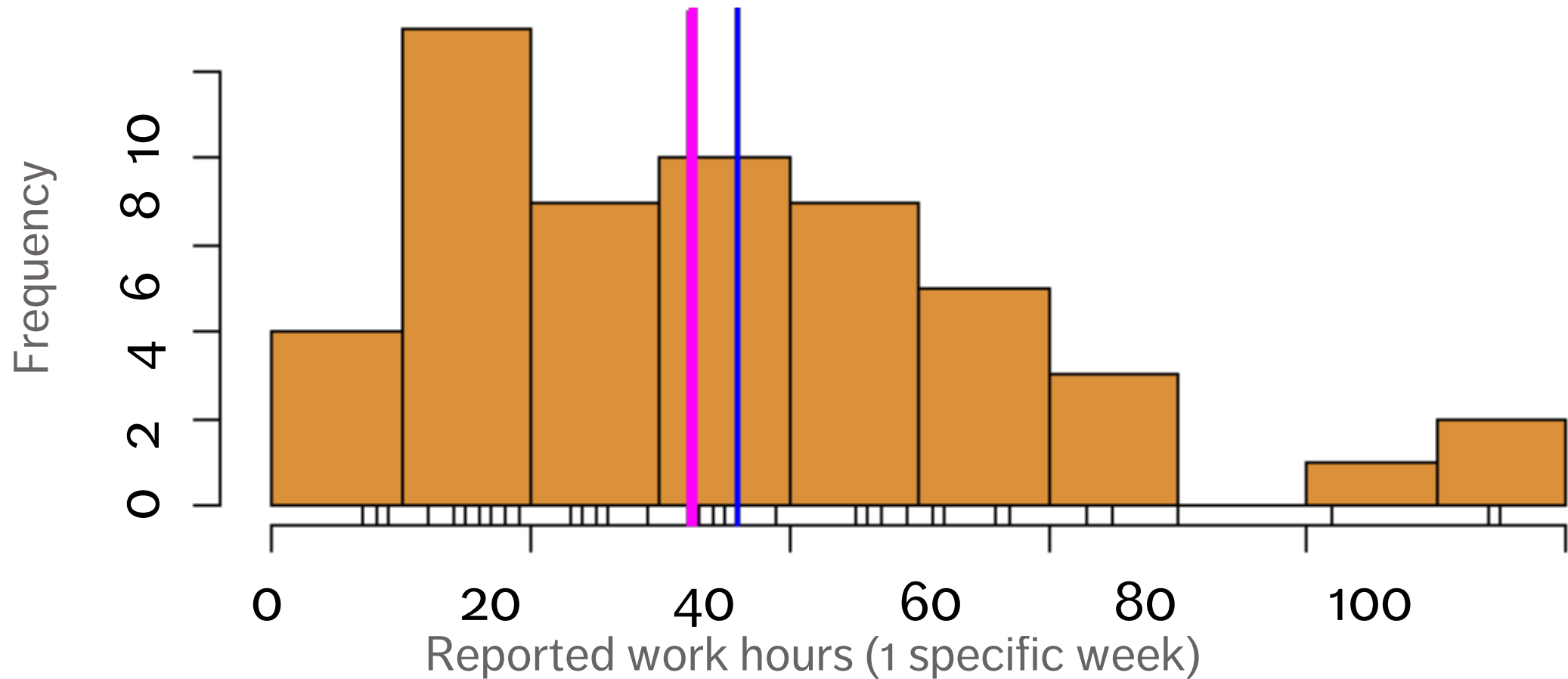
Consistent baseline on far left and right.

Easy to compare.

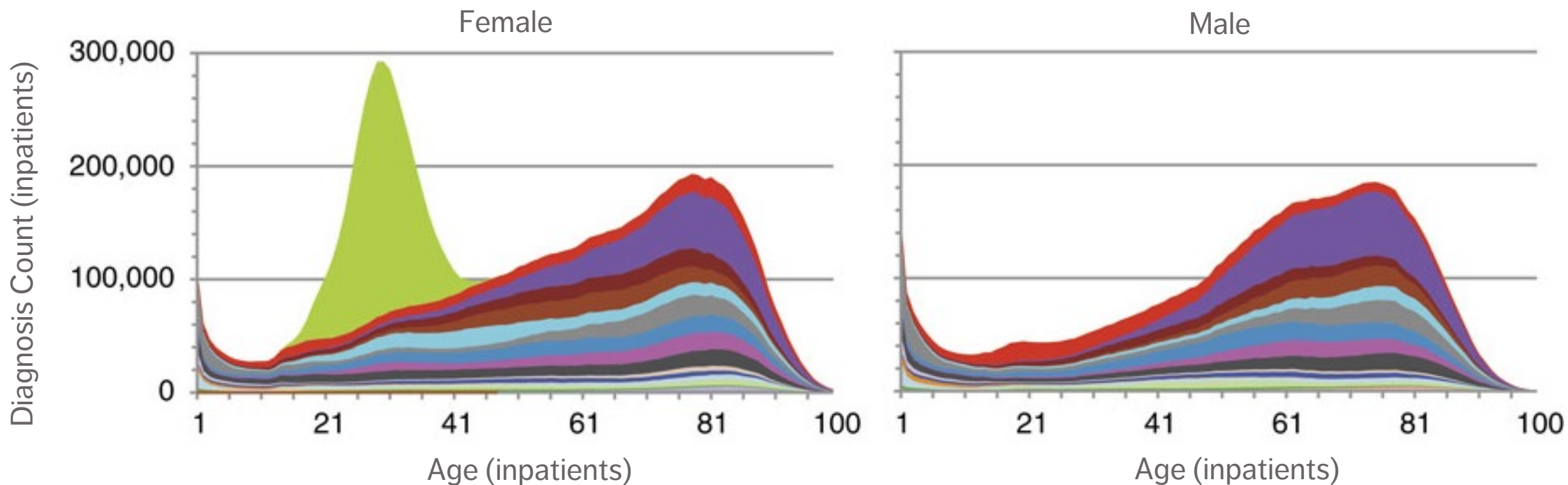
No relative measure to **magnitude** of data.

Research shows that horizontal is easier to process than vertical.

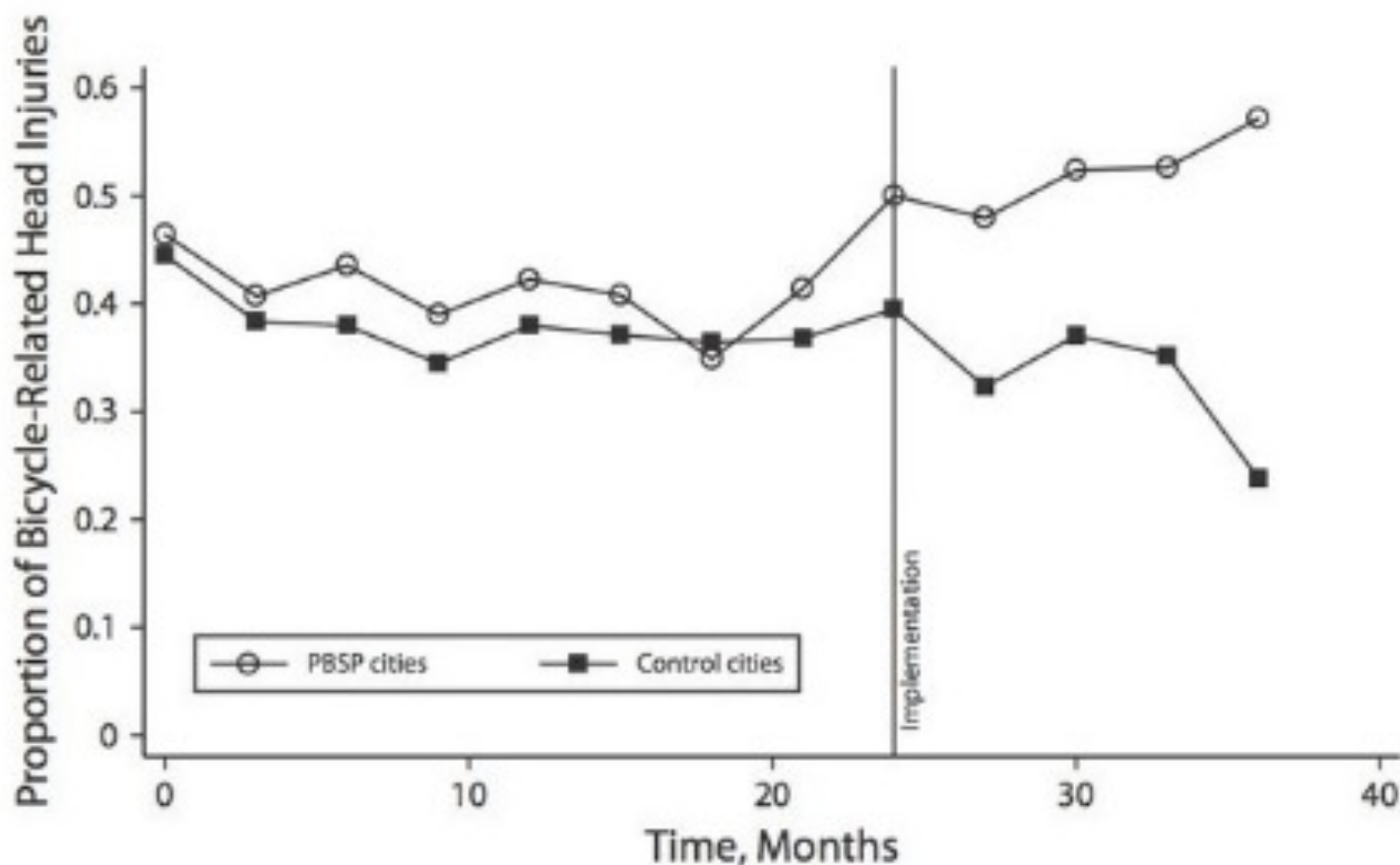
# Histogram



# Stacked Histograms

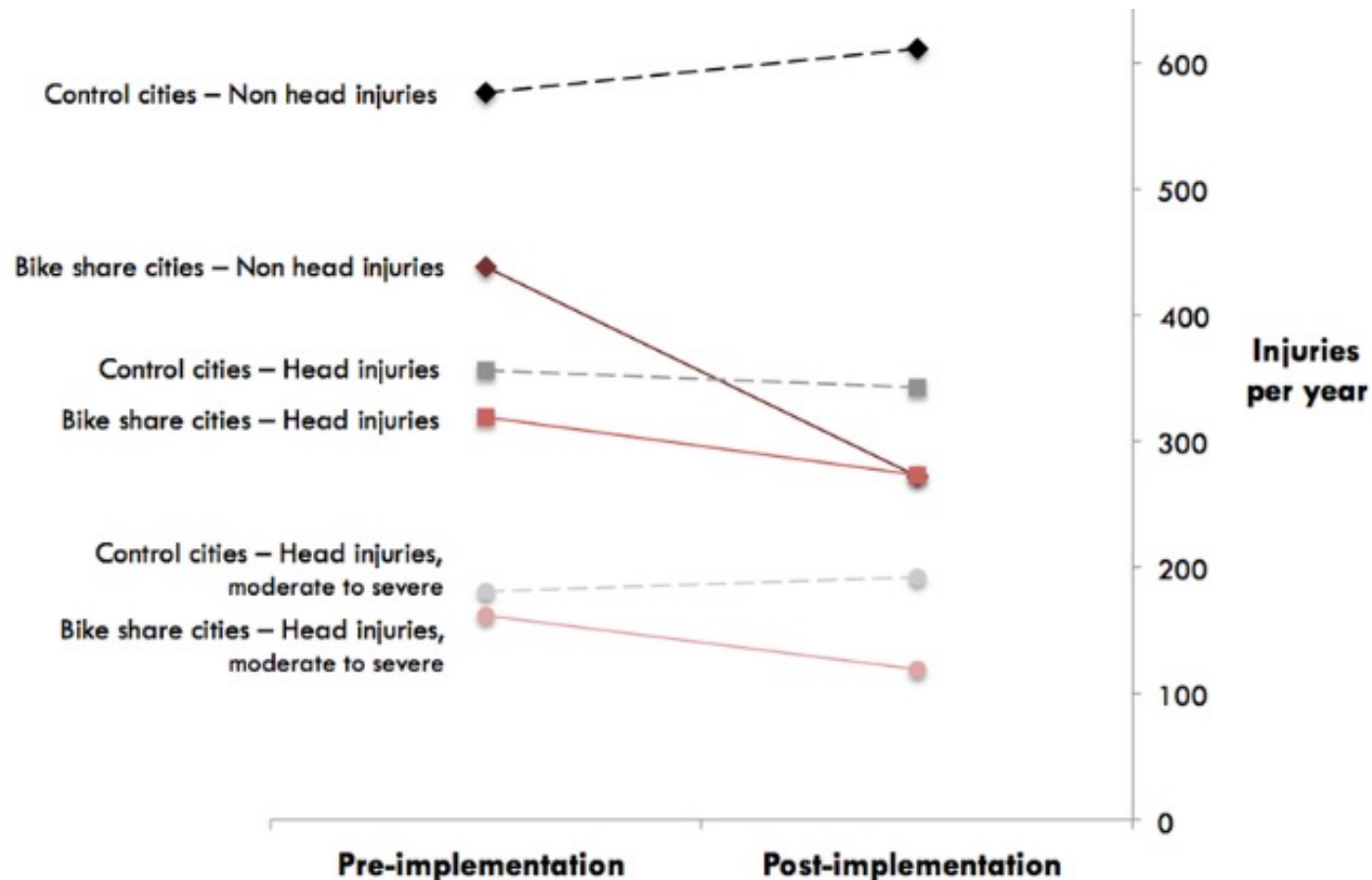


# Line Graphs



Proportion of all bicycle-related injuries that were classified as head injuries among cities with public bike share programs and control cities, centered on intervention date (vertical line); North America.

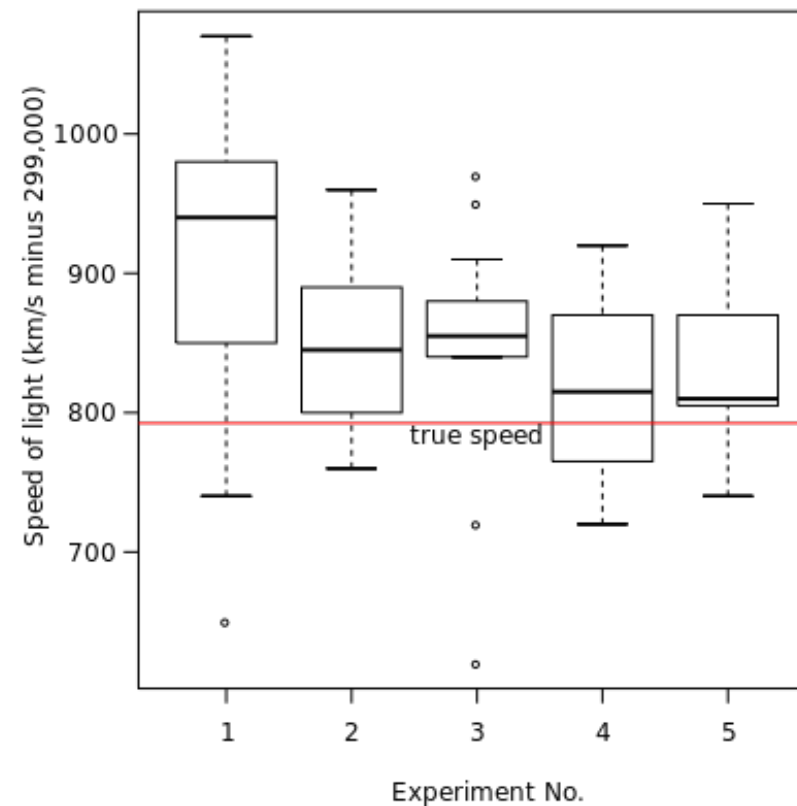
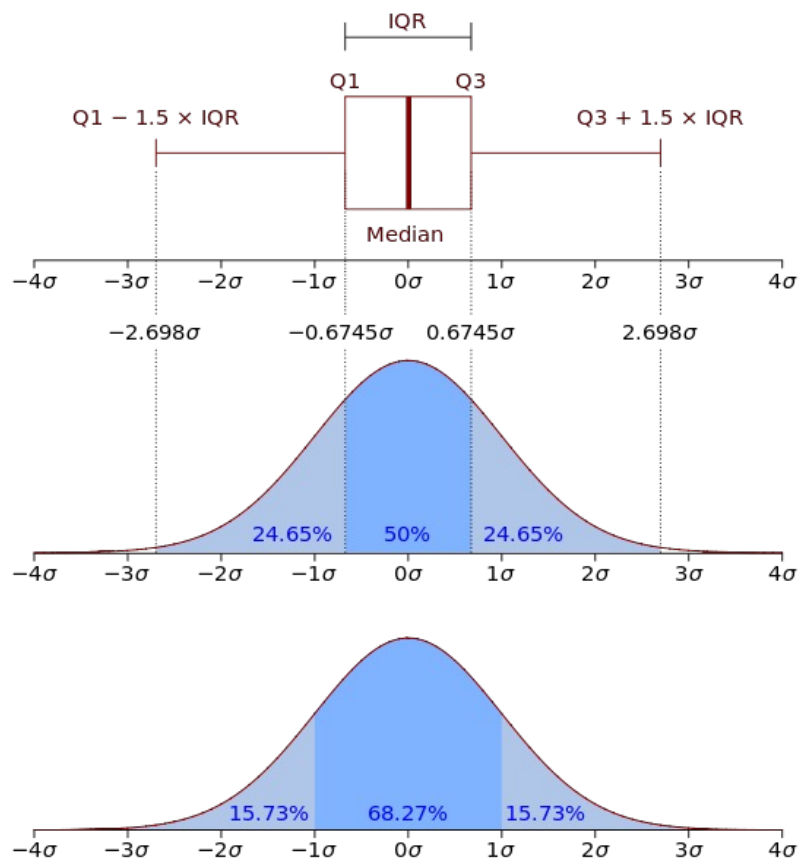
# Line Graphs



Data from new study show declines in all injuries, including head injuries after bike share system implemented.

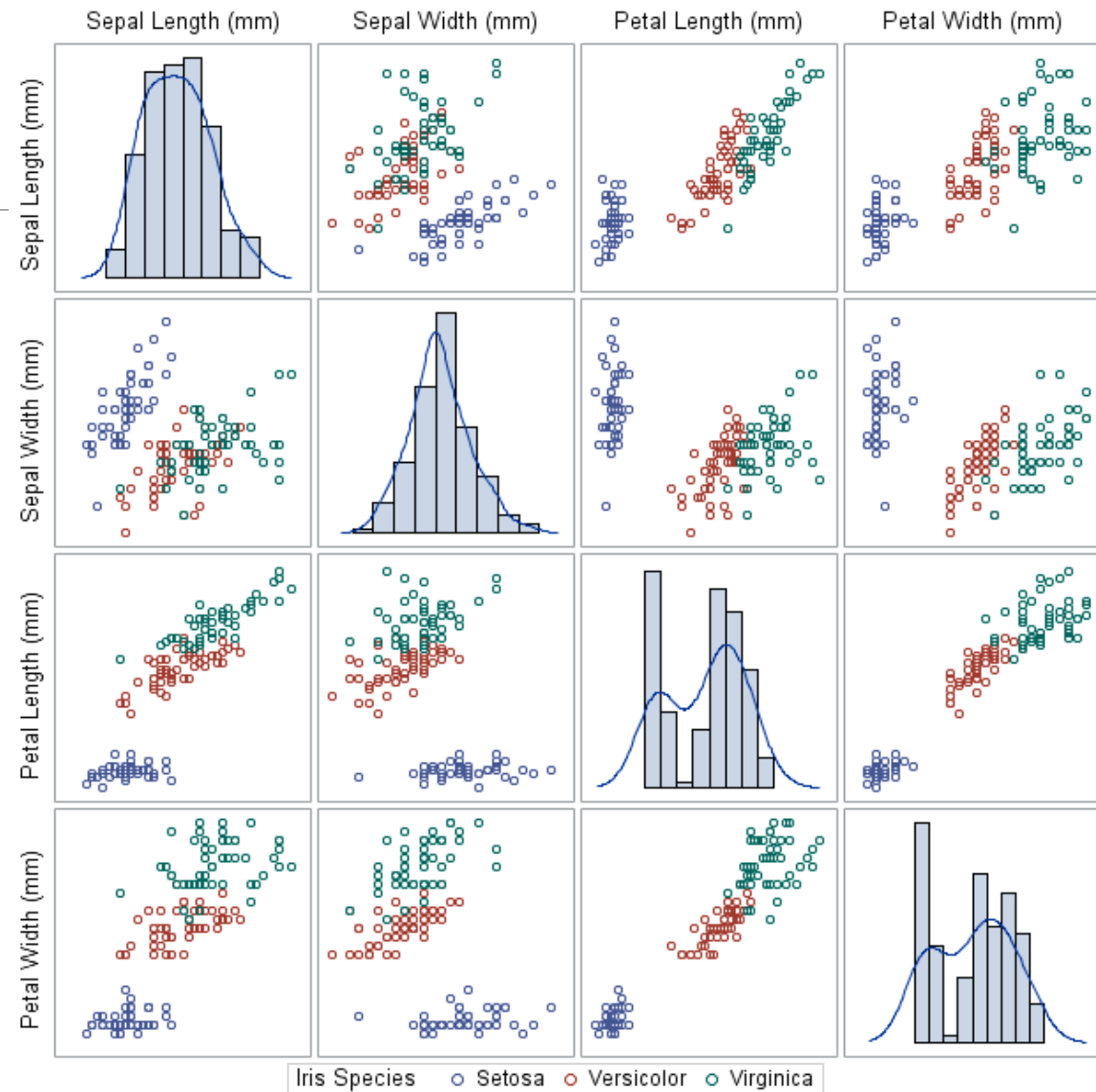
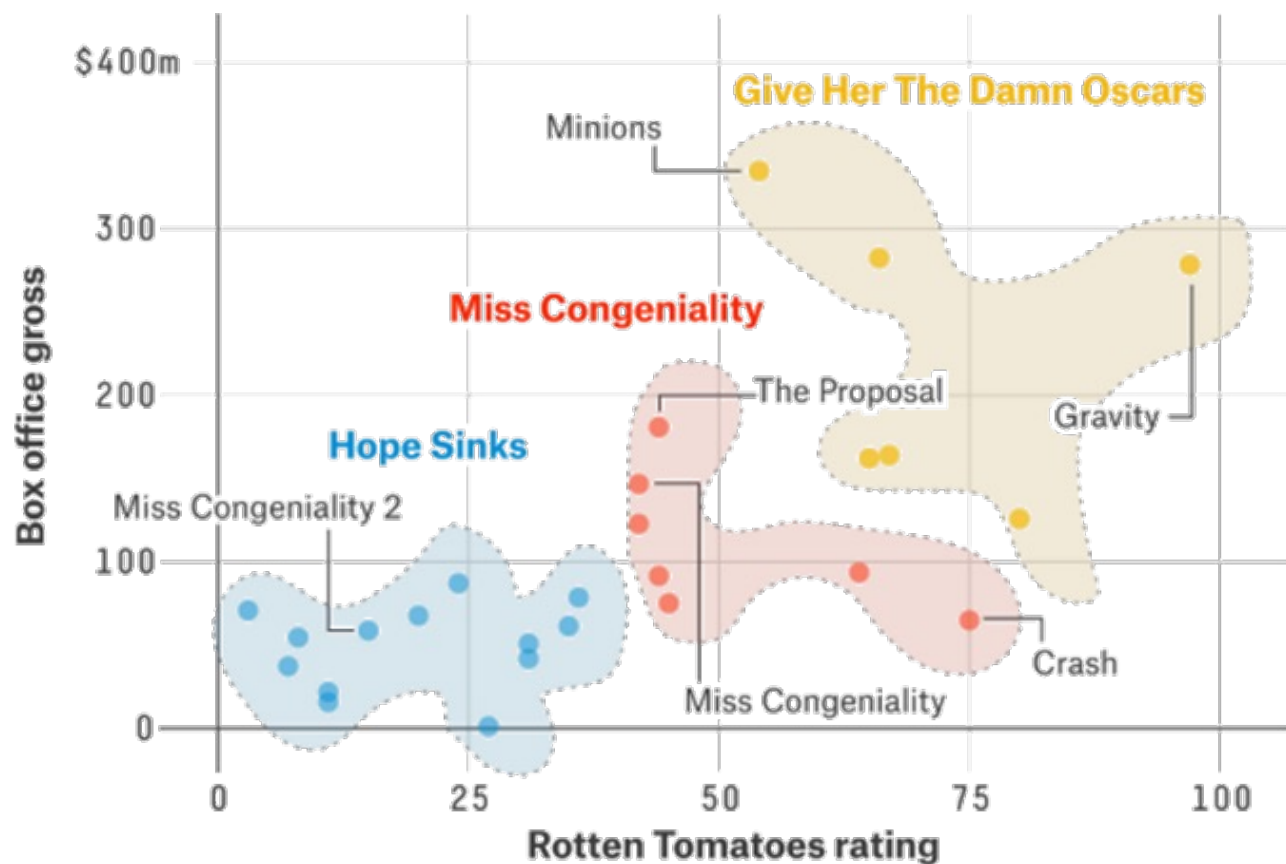
Because head injuries decline less than other injuries, they are now a larger proportion of all injuries.

# Boxplots



Experiment No.  
Michelson-Morley Experiment to Determine the Speed of Light

# Scatterplots



# Suggested Reading

Exploratory Data Analysis

*Data Understanding, Data Analysis, Data Science*  
**Data Visualization and Data Exploration**

Data and Charts

- Pre-Analysis Uses

---

*The Practice of Data Visualization*  
**Basics of Data Visualization**

Data Exploration

Workhorse Data Visualizations



# Exercises

## Exploratory Data Analysis

1. Find examples of data presentations that you consider to be particularly insightful and/or powerful. Discuss their strengths/weaknesses.
2. Find examples of data presentations that you consider to be particularly misleading and/or useless. Discuss their strengths/weaknesses.
3. How do you think new technologies (e.g. virtual or augmented reality, 3D-printing, wearable computing) will influence data presentations?