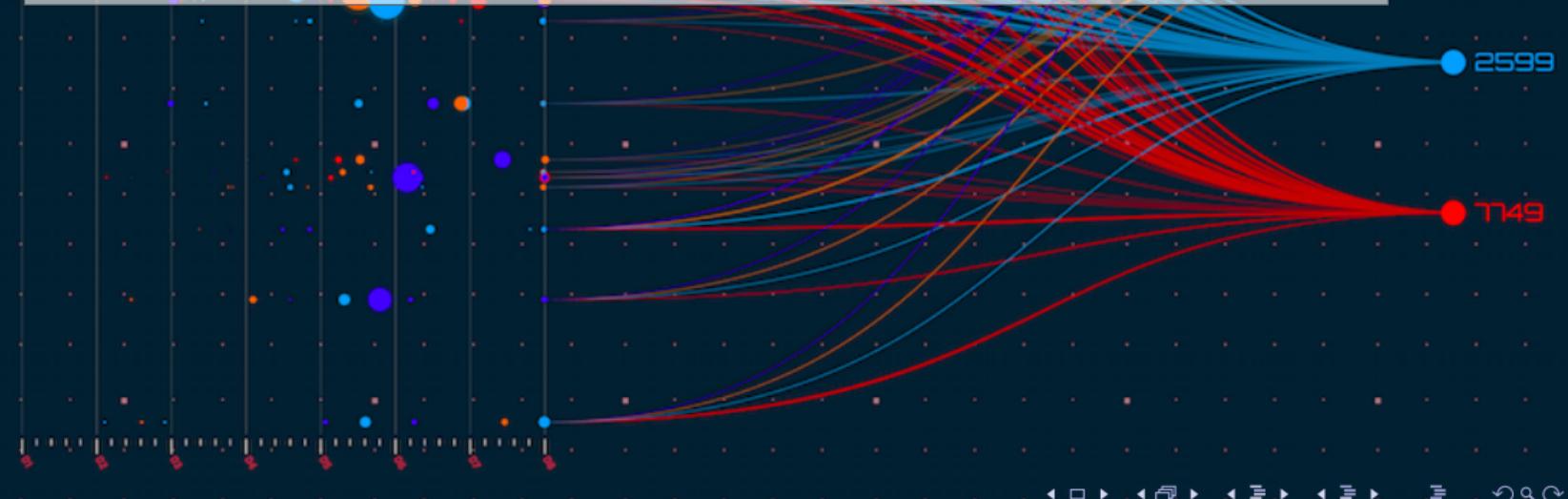
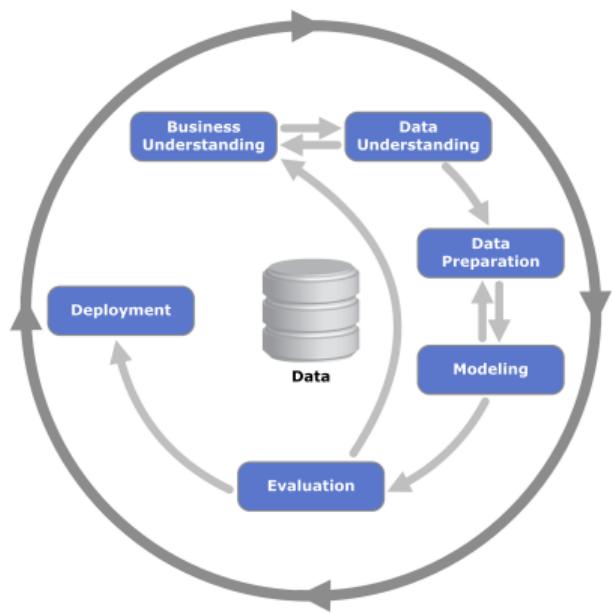


Data Visualization



Learning Outcomes

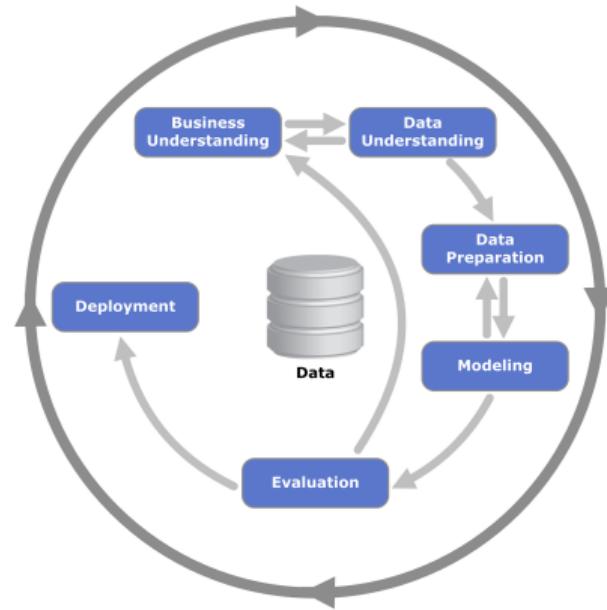


The [CRISP-DM](#) process.

- ▶ Define exploratory data analysis;
- ▶ List goals of visualizations;
- ▶ Interpret 10 commonly used visualizations;
- ▶ Offer critique of bad visualizations.

Exploratory Data Analysis

- ▶ Exploratory data analysis (EDA) is the search for patterns and trends in a given data set;
- ▶ Goals:
 - ▶ Suggest interesting hypotheses;
 - ▶ Identify data processing mistakes/outliers;
 - ▶ Find violations of statistical assumptions;
 - ▶ Identify a potential set of features for ML.



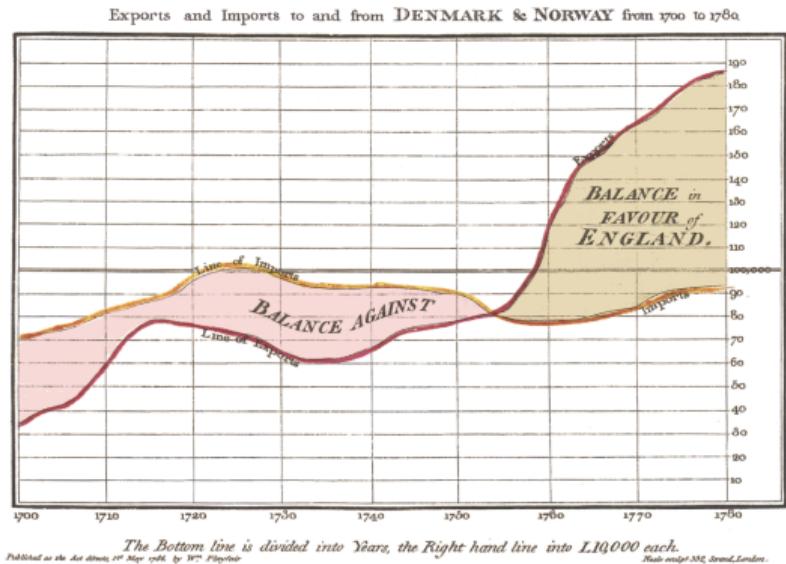
EDA falls at the end of data prep and the beginning of modeling.

Confronting New Data

- ▶ Get a feel for the data by considering...
 - ▶ What is the data generating process? Who got the data, when, and why?
 - ▶ How big is it?
 - ▶ What do the fields mean?
 - ▶ Gut checks – look for familiar or interpretable records;
 - ▶ Summary statistics, e.g. Max, Min, Quartiles, Median, etc...
- ▶ ...and also by applying visualization!
 - ▶ A key part of EDA and beyond:
 - ▶ Get a feel for what your data really looks like. Identify trends over dataset that might not be obvious from summary statistics;
 - ▶ Detect errors, e.g. outliers, insufficient cleaning, coding mistakes, bad assumptions;
 - ▶ Communicate what you found – results only become actionable after they are shared.

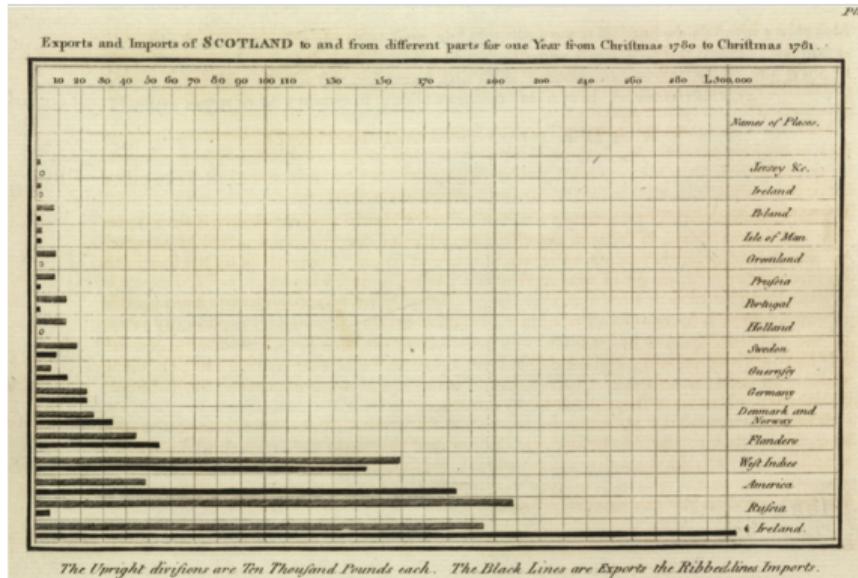
Some history...

- ▶ Many ways to extract value from data:
 - ▶ descriptive statistics;
 - ▶ hypothesis testing;
 - ▶ simulation;
 - ▶ statistical modeling;
 - ▶ machine learning/AI;
- ▶ One of the very oldest is **data visualization**:
 - ▶ William Playfair often cited as the inventor of line, bar, and pie charts, late 18th century.



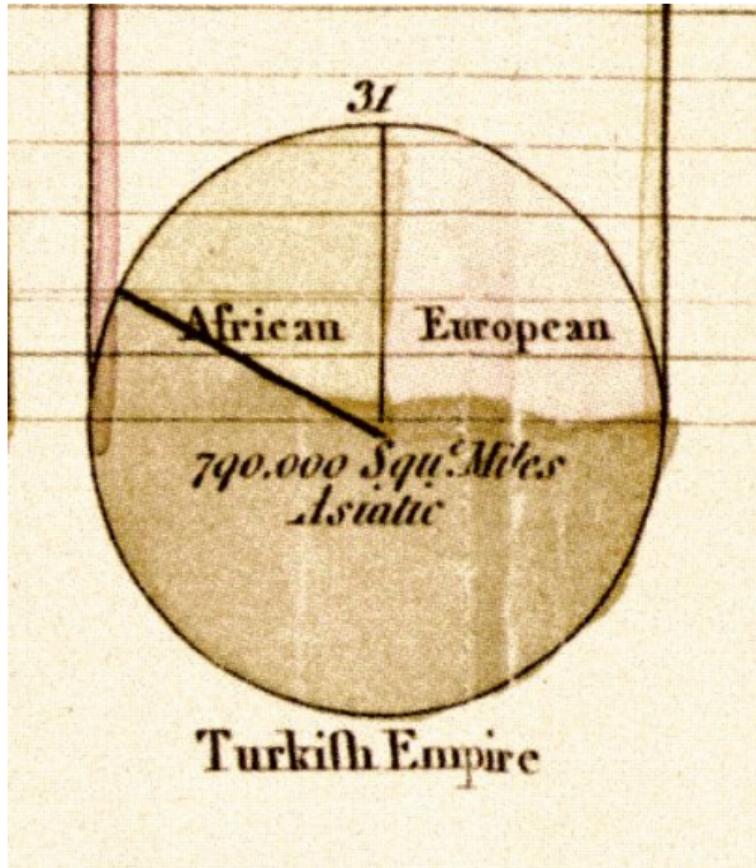
Some history...

- ▶ Many ways to extract value from data:
 - ▶ descriptive statistics;
 - ▶ hypothesis testing;
 - ▶ simulation;
 - ▶ statistical modeling;
 - ▶ machine learning/AI;
- ▶ One of the very oldest is **data visualization**:
 - ▶ William Playfair often cited as the inventor of line, bar, and pie charts, late 18th century.



Some history...

- ▶ Many ways to extract value from data:
 - ▶ descriptive statistics;
 - ▶ hypothesis testing;
 - ▶ simulation;
 - ▶ statistical modeling;
 - ▶ machine learning/AI;
- ▶ One of the very oldest is **data visualization**:
 - ▶ William Playfair often cited as the inventor of line, bar, and pie charts, late 18th century.

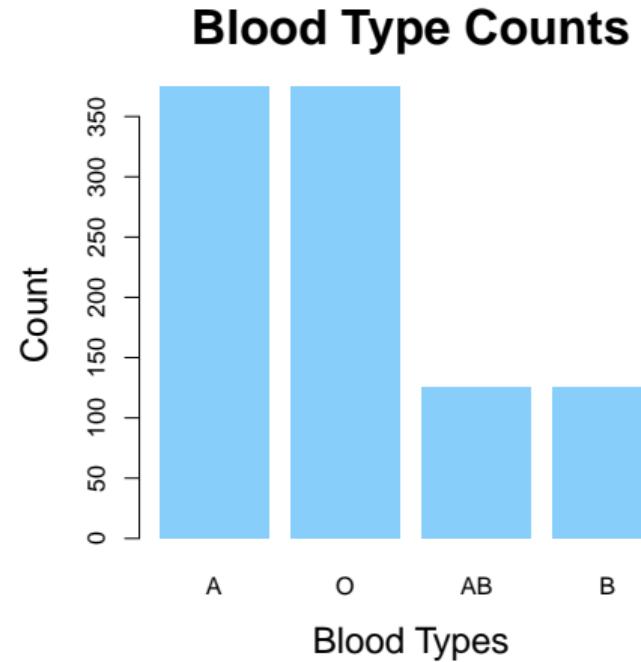


Goals:

- ▶ Characteristics of data visualization:
 - ▶ Data encoded as a visual object;
 - ▶ Goal of any data analysis is to help audience reason about the data – data viz is most reliant on human cognition for inference and meaning;
 - ▶ Form = substance;
 - ▶ Generally, necessarily low dimensional;
- ▶ Data visualization should (Tufte 1983):
 - ▶ Have message fit for audience purpose and communicate information;
 - ▶ Highlight the data to help the audience focus on the substance;
 - ▶ NOT distort the data or knowingly perpetuate a lie;
 - ▶ Put many numbers in a small, accessible space;
 - ▶ Enable the audience to grapple with large data sets coherently;
 - ▶ Encourage viewers to compare data effectively;
 - ▶ Allow investigation of both micro and macro level structure in the data;
 - ▶ Integrate closely with additional statistical analysis;
- ▶ Libraries: Matplotlib, Seaborn, Pandas, Plotly.

Some of my favorite types of graphics – one variable

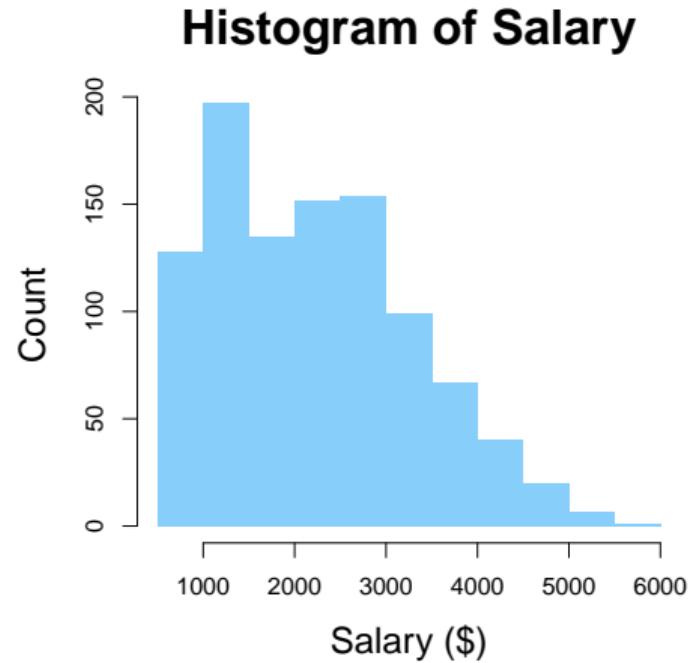
- ▶ **Bar chart;**
 - ▶ Compare amount in fixed bins (**categorical** data);
- ▶ **Histogram;**
 - ▶ Compare amount in constructed bins (**numerical** data);
- ▶ **Density plot;**
 - ▶ Smoothed amounts in constructed bins (**numerical** data);



A bar chart.

Some of my favorite types of graphics – one variable

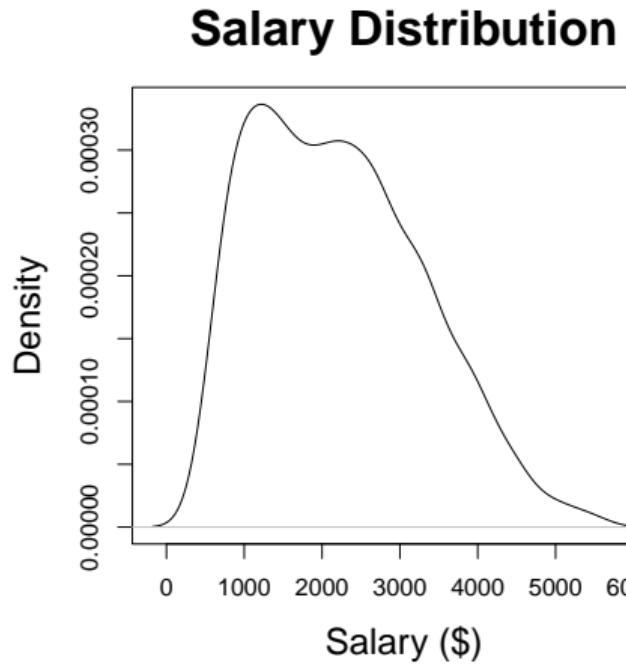
- ▶ **Bar chart;**
 - ▶ Compare amount in fixed bins (**categorical** data);
- ▶ **Histogram;**
 - ▶ Compare amount in constructed bins (**numerical** data);
- ▶ **Density plot;**
 - ▶ Smoothed amounts in constructed bins (**numerical** data);



A histogram.

Some of my favorite types of graphics – one variable

- ▶ **Bar chart;**
 - ▶ Compare amount in fixed bins (**categorical** data);
- ▶ **Histogram;**
 - ▶ Compare amount in constructed bins (**numerical** data);
- ▶ **Density plot;**
 - ▶ Smoothed amounts in constructed bins (**numerical** data);



A density.

Some of my favorite types of graphics – one variable

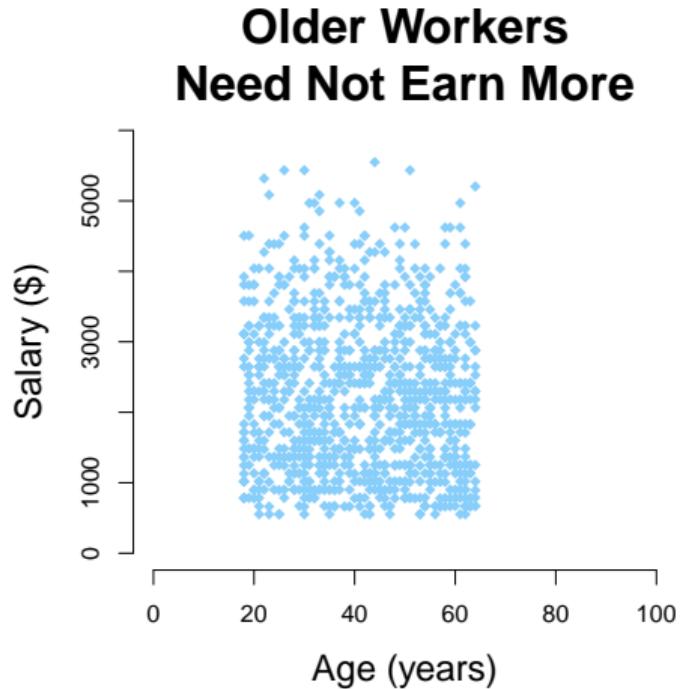
- ▶ **Bar chart;**
 - ▶ Compare amount in fixed bins (**categorical** data);
- ▶ **Histogram;**
 - ▶ Compare amount in constructed bins (**numerical** data);
- ▶ **Density plot;**
 - ▶ Smoothed amounts in constructed bins (**numerical** data);



A nicer looking density.

Some of my favorite types of graphics – two variables

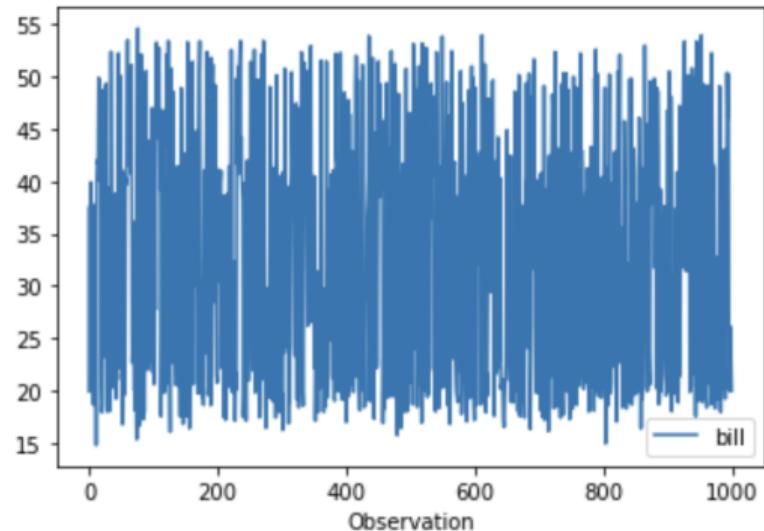
- ▶ Scatterplot/line plot/scatterplot matrix;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Joint plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Bivariate Kernel Density Plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Boxplot/violin plot;
 - ▶ Compare two variables (categorical/numerical data);
- ▶ Heatmap;
 - ▶ Compare two variables (categorical/categorical data);



A scatterplot.

Some of my favorite types of graphics – two variables

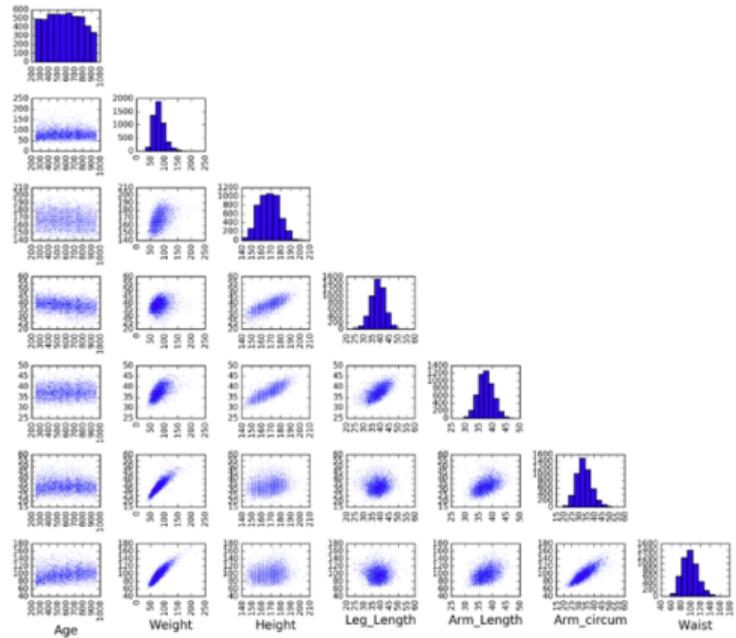
- ▶ Scatterplot/line plot/scatterplot matrix;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Joint plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Bivariate Kernel Density Plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Boxplot/violin plot;
 - ▶ Compare two variables (categorical/numerical data);
- ▶ Heatmap;
 - ▶ Compare two variables (categorical/categorical data);



A line plot.

Some of my favorite types of graphics – two variables

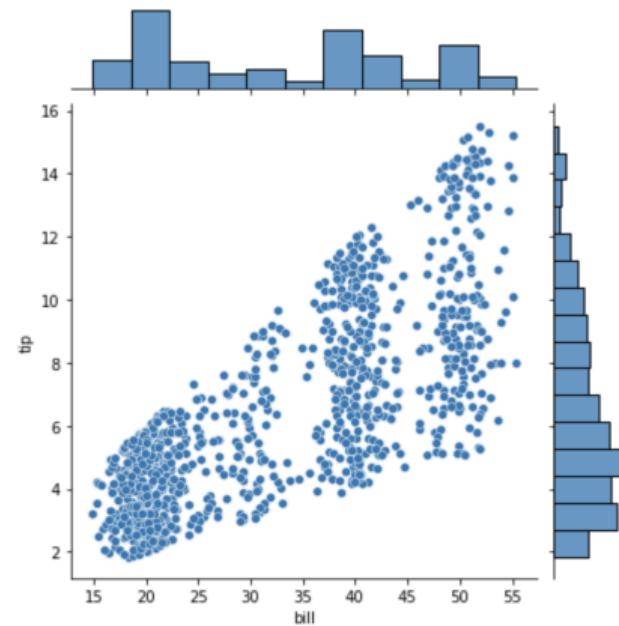
- ▶ Scatterplot/line plot/scatterplot matrix;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Joint plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Bivariate Kernel Density Plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Boxplot/violin plot;
 - ▶ Compare two variables (categorical/numerical data);
- ▶ Heatmap;
 - ▶ Compare two variables (categorical/categorical data);



A scatterplot matrix.

Some of my favorite types of graphics – two variables

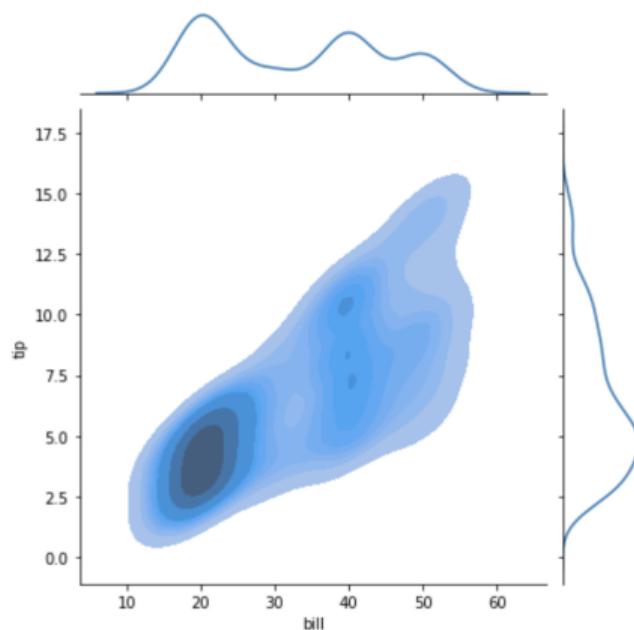
- ▶ Scatterplot/line plot/scatterplot matrix;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Joint plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Bivariate Kernel Density Plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Boxplot/violin plot;
 - ▶ Compare two variables (categorical/numerical data);
- ▶ Heatmap;
 - ▶ Compare two variables (categorical/categorical data);



A joint plot.

Some of my favorite types of graphics – two variables

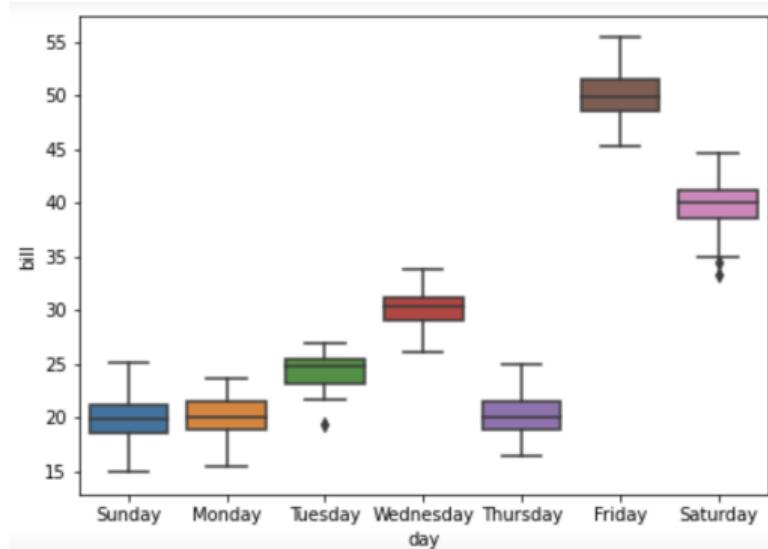
- ▶ Scatterplot/line plot/scatterplot matrix;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Joint plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Bivariate Kernel Density Plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Boxplot/violin plot;
 - ▶ Compare two variables (categorical/numerical data);
- ▶ Heatmap;
 - ▶ Compare two variables (categorical/categorical data);



A bivariate kernel density plot.

Some of my favorite types of graphics – two variables

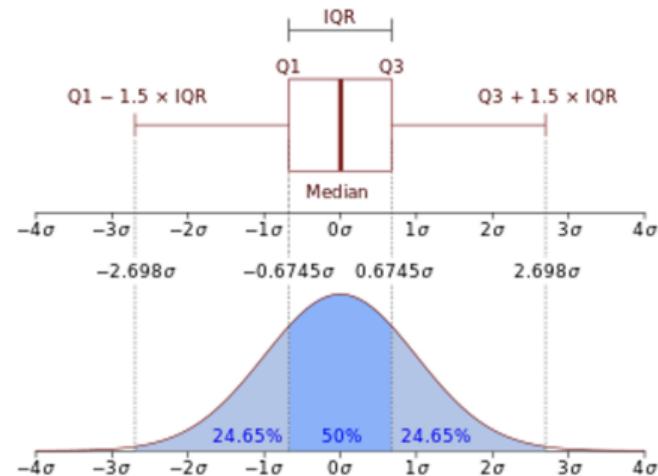
- ▶ Scatterplot/line plot/scatterplot matrix;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Joint plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Bivariate Kernel Density Plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Boxplot/violin plot;
 - ▶ Compare two variables (categorical/numerical data);
- ▶ Heatmap;
 - ▶ Compare two variables (categorical/categorical data);



A box plot.

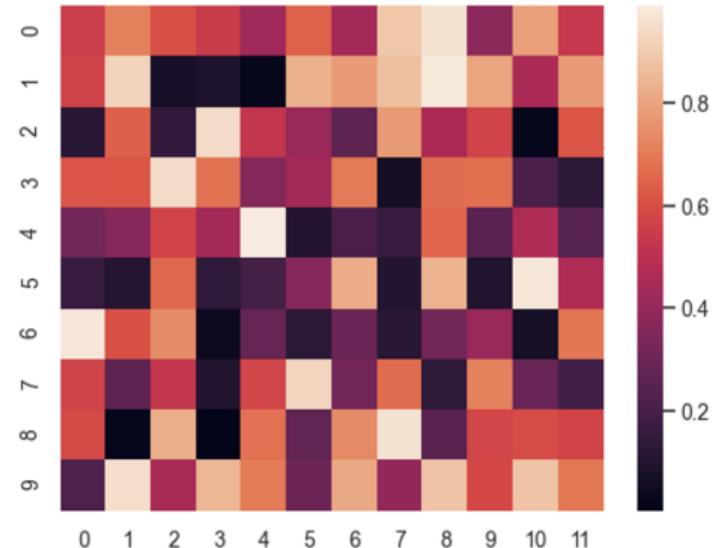
Some of my favorite types of graphics – two variables

- ▶ Scatterplot/line plot/scatterplot matrix;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Joint plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Bivariate Kernel Density Plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Boxplot/violin plot;
 - ▶ Compare two variables (categorical/numerical data);
- ▶ Heatmap;
 - ▶ Compare two variables (categorical/categorical data);



Some of my favorite types of graphics – two variables

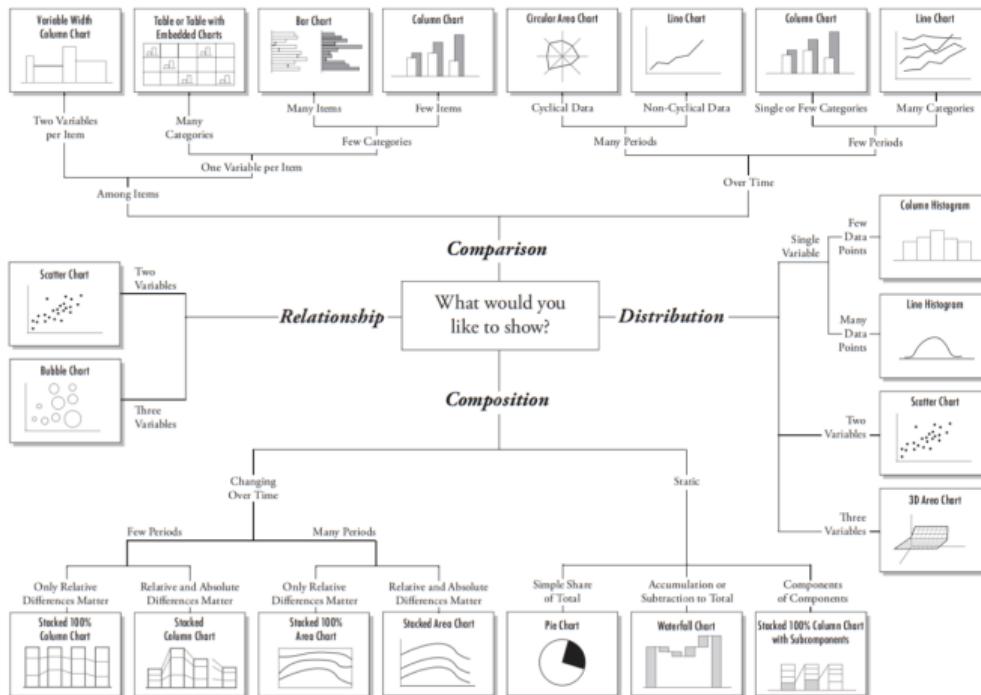
- ▶ Scatterplot/line plot/scatterplot matrix;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Joint plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Bivariate Kernel Density Plot;
 - ▶ Compare two variables (numerical/numerical data);
- ▶ Boxplot/violin plot;
 - ▶ Compare two variables (categorical/numerical data);
- ▶ Heatmap;
 - ▶ Compare two variables (categorical/categorical data);



And finally a heatmap.

What kinds of charts are there?

Chart Suggestions—A Thought-Starter



Python methods for all these...

- ▶ Bar Charts:
 - ▶ Seaborn: `.countplot()`;
 - ▶ Pandas:
`df.value_counts().plot_bar()`;
- ▶ Histograms:
 - ▶ Matplotlib: `.hist()`;
 - ▶ Seaborn: `.distplot()`;
 - ▶ Pandas: `df.plot.hist()`;
- ▶ Densities:
 - ▶ Seaborn: `.kdeplot()`;
 - ▶ Pandas: `df.plot.kde()`;
- ▶ Scatterplots:
 - ▶ Matplotlib: `.scatter()`;
 - ▶ Seaborn: `.scatterplot()`, `.pairplot()`;
 - ▶ Pandas: `df.plot.scatter()`;
- ▶ Line plot:
 - ▶ Matplotlib: `.plot()`;
 - ▶ Seaborn: `.lineplot()`;
 - ▶ Pandas: `df.plot.line()`;
- ▶ Joint plot:
 - ▶ Seaborn: `.jointplot()`;
- ▶ Bivariate Kernel Density plot:
 - ▶ Seaborn: `.jointplot()`;
- ▶ Box plot:
 - ▶ Seaborn: `.boxplot()`;
 - ▶ Pandas: `df.boxplot()`;
- ▶ Violin plot:
 - ▶ Seaborn: `.violinplot()`;
- ▶ Heatmap:
 - ▶ Seaborn: `.heatmap()`;

Some examples...

Napoleon's Progress – an exemplar?

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

Dessiné 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 largures des zones colorées à raison d'un millimètre pour dix mille hommes; ils sont de plus écrits en traits de traits des zones. Le rouge désigne les hommes qui ont été 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. Chiers, de Léger, de Fézensac, de Chambray et le journal intime de Jacob, pharmacien de l'Armée depuis le 23 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 Davout, qui avaient été détachés sur Minsk et Maliblow et qui rejoignirent Ochta et Vitebsk, avaient toujours marché avec l'armée.

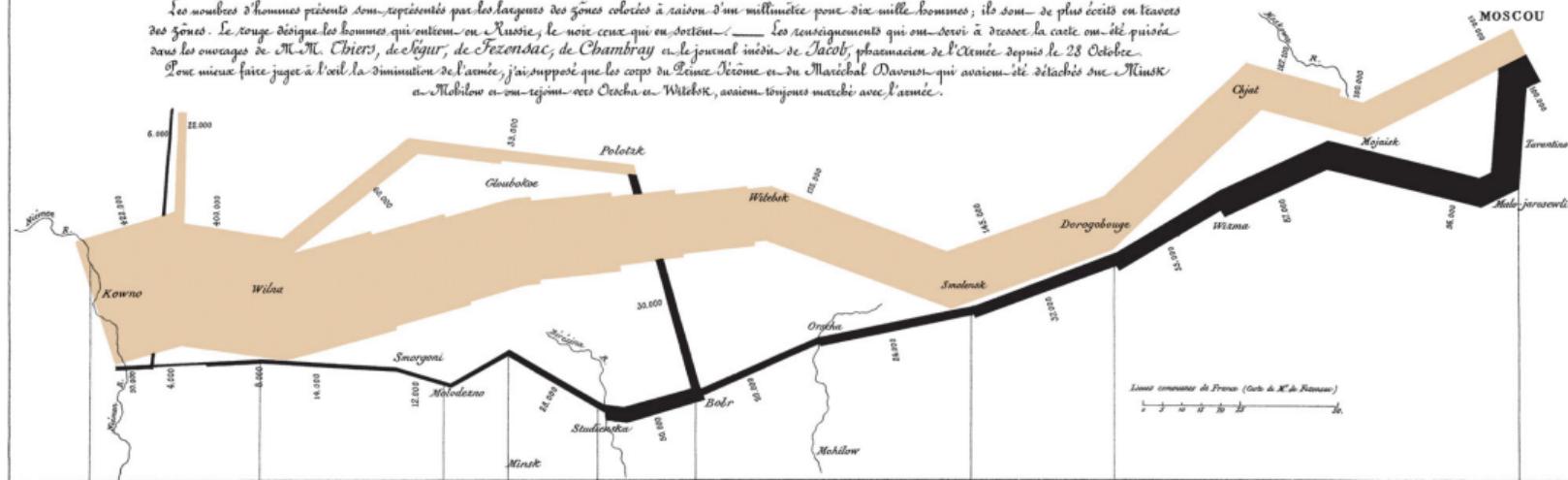
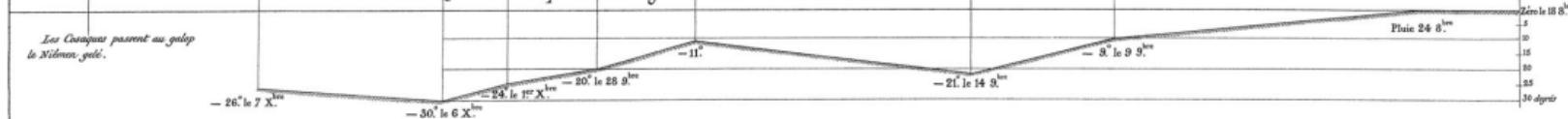


TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.



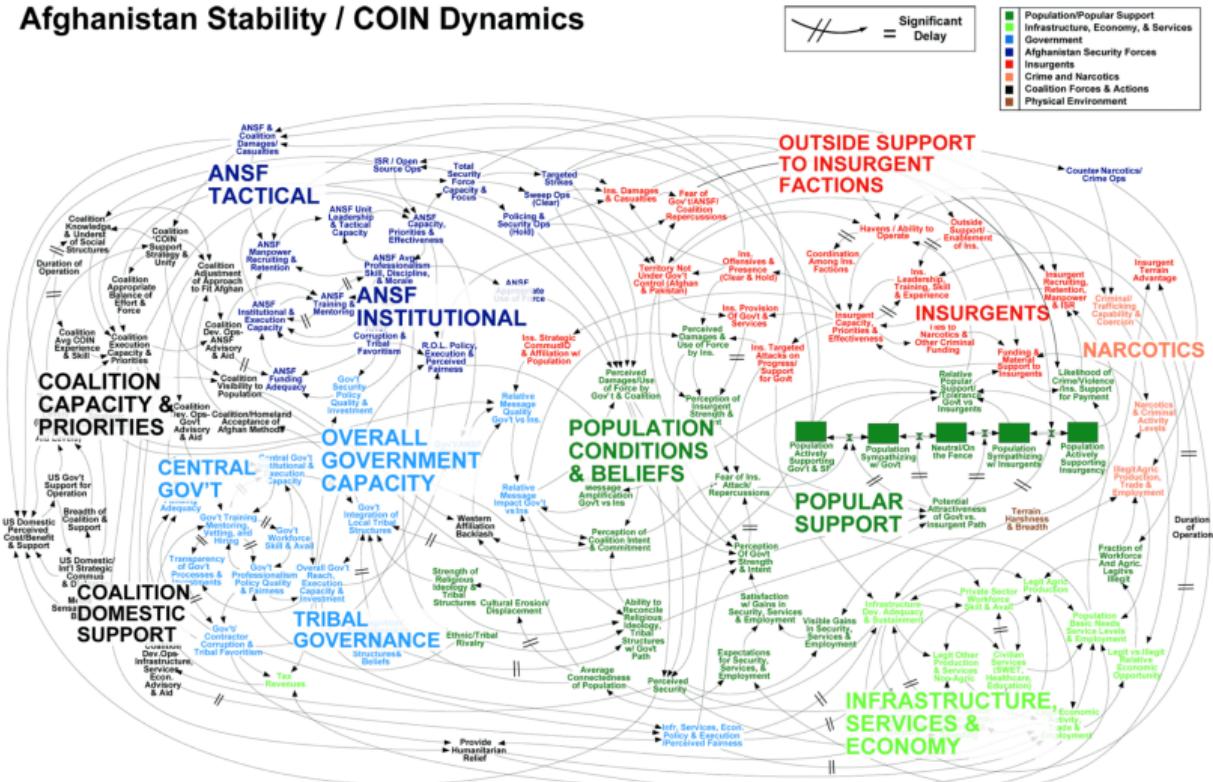
Autog. par Rognier, 8, Rue S^e Marie S^e 0^e à Paris.

Imp. Litt. Rognier et Desordain



Dynamics in Afghanistan

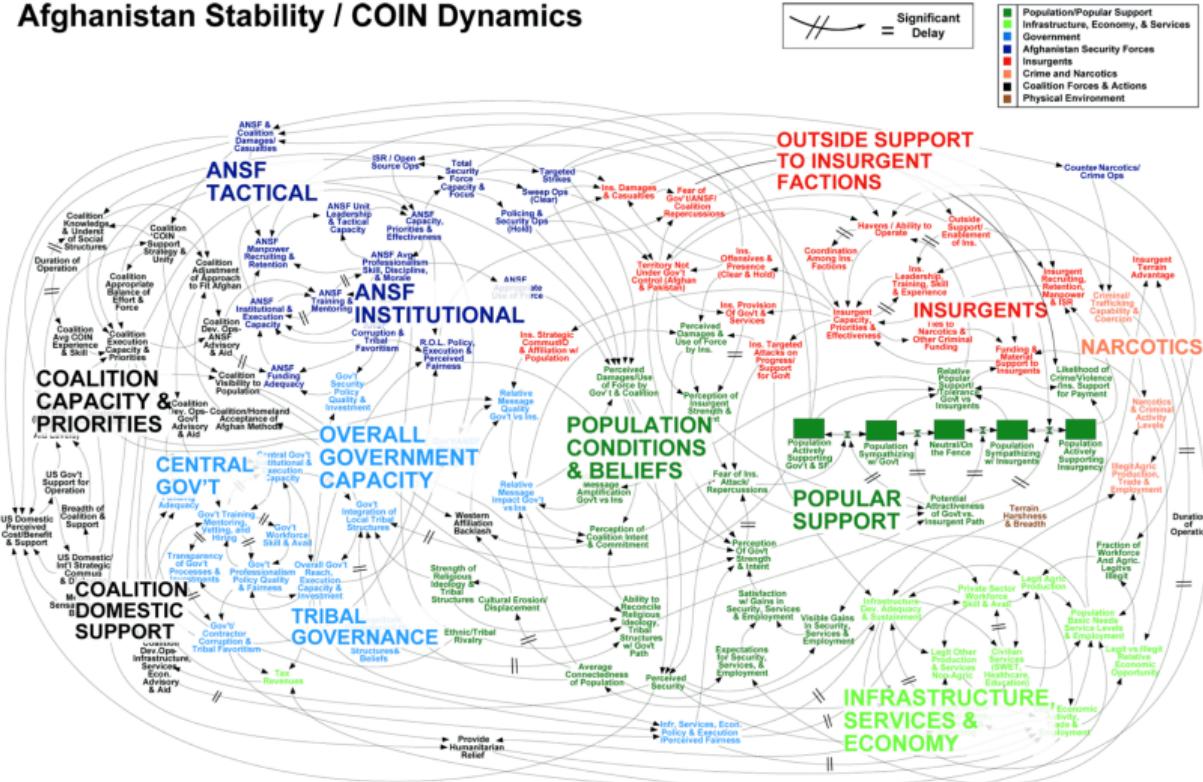
Afghanistan Stability / COIN Dynamics



Dynamics in Afghanistan – Impossible to focus.

Afghanistan Stability / COIN Dynamics

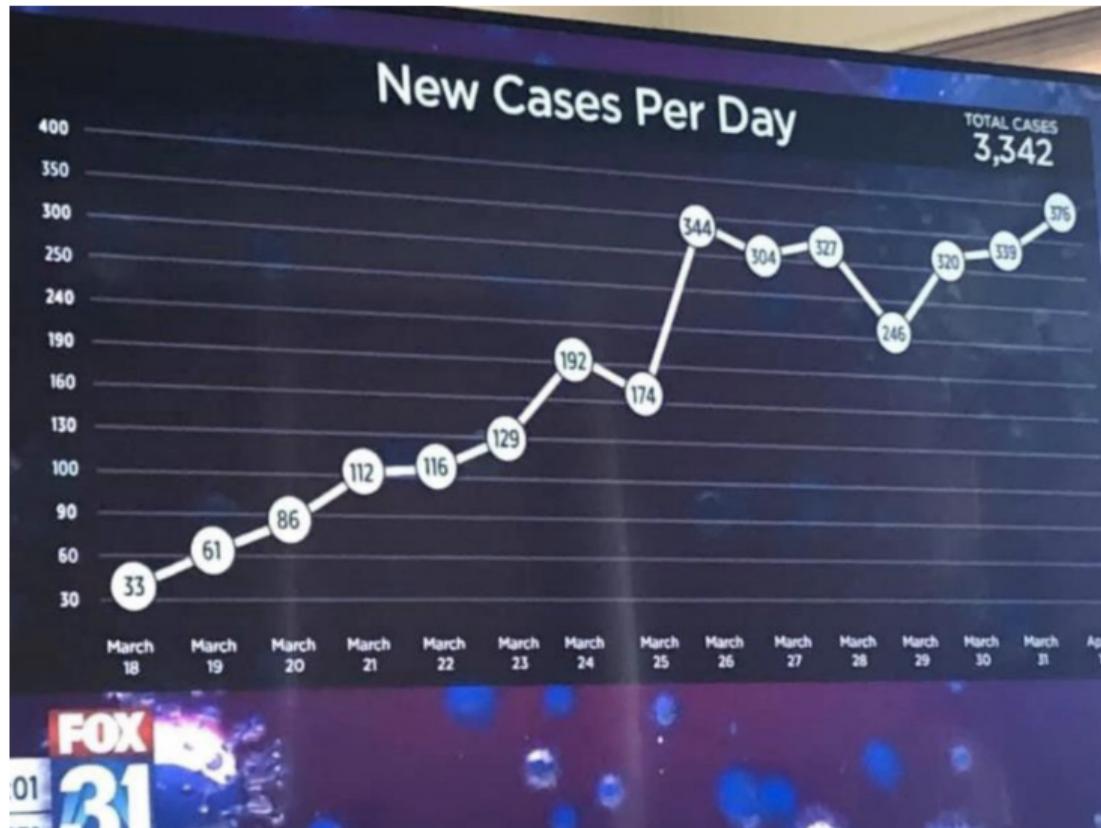
 = Significant Delay



COVID Timeline



COVID Timeline – Non-linear *y*-axis.



COVID's Reach



Horrifying new map shows no country is safe from coronavirus' deadly tentacles



Horrifying new map reveals no country safe from coronavirus' deadly tentacles
A HORRIFYING new map shows the unstoppable spread of deadly coronavirus across the globe. The incredible graphic reveals how five million Wuhan resident...
thesun.co.uk

COVID's Reach – Underlying data is unrelated to graphic message.

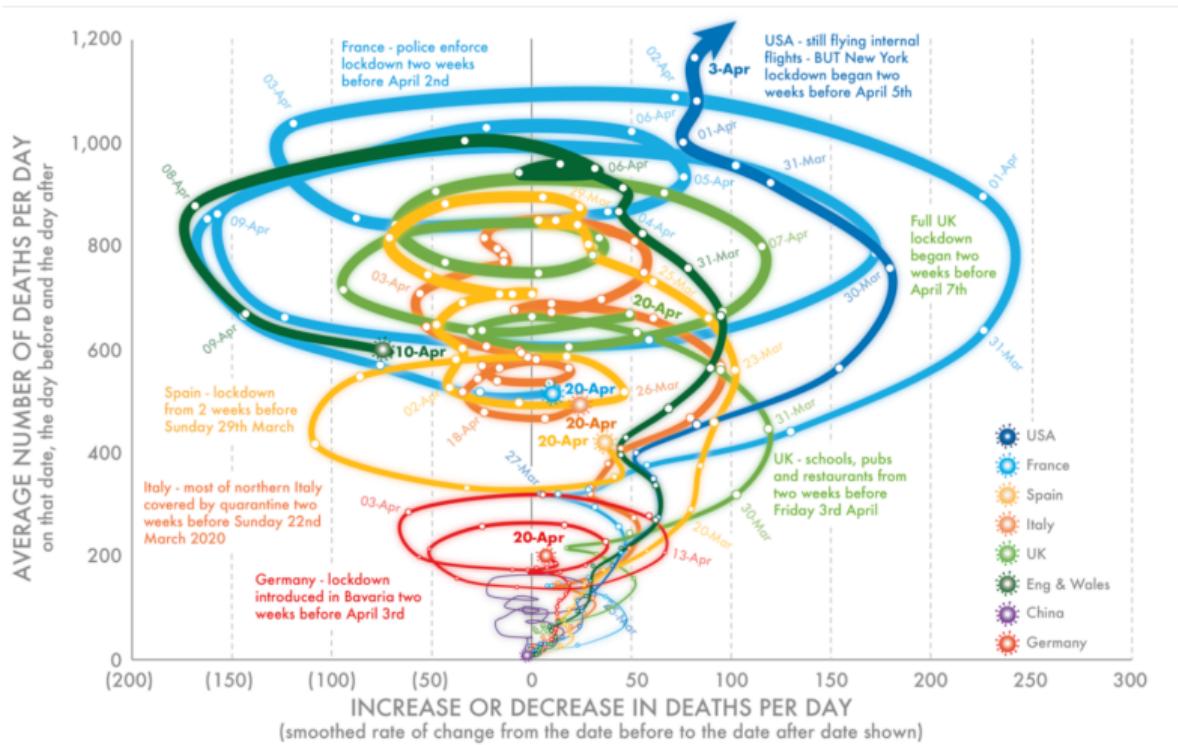


Horrifying new map shows no country is safe from coronavirus' deadly tentacles

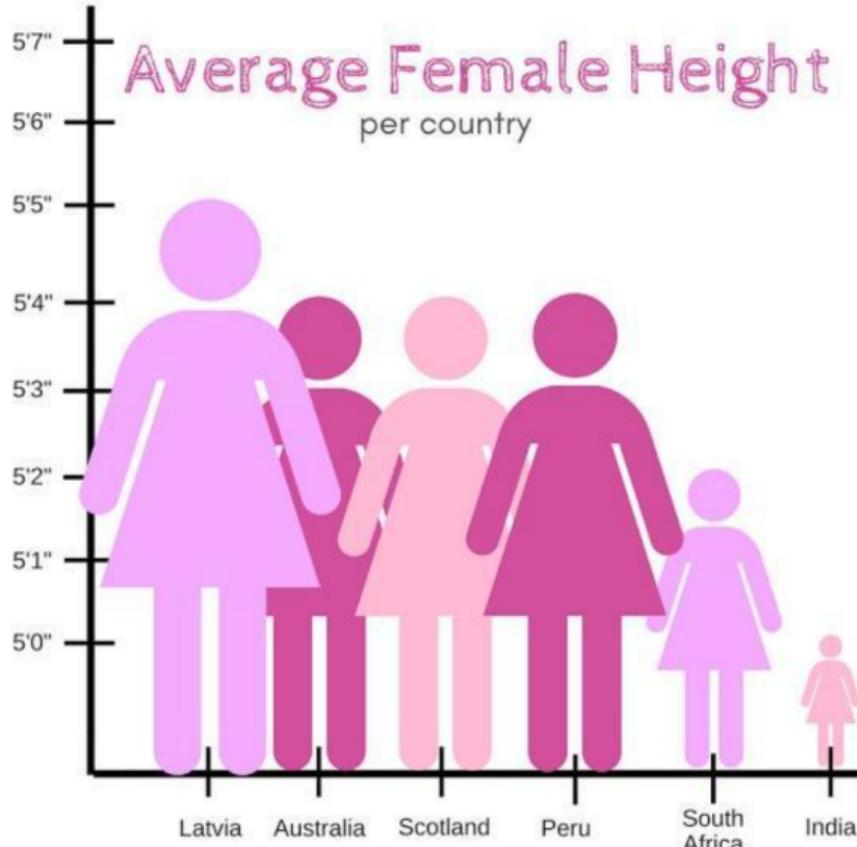


Horrifying new map reveals no country safe from coronavirus' deadly tentacles
A HORRIFYING new map shows the unstoppable spread of deadly coronavirus across the globe. The incredible graphic reveals how five million Wuhan resident...
thesun.co.uk

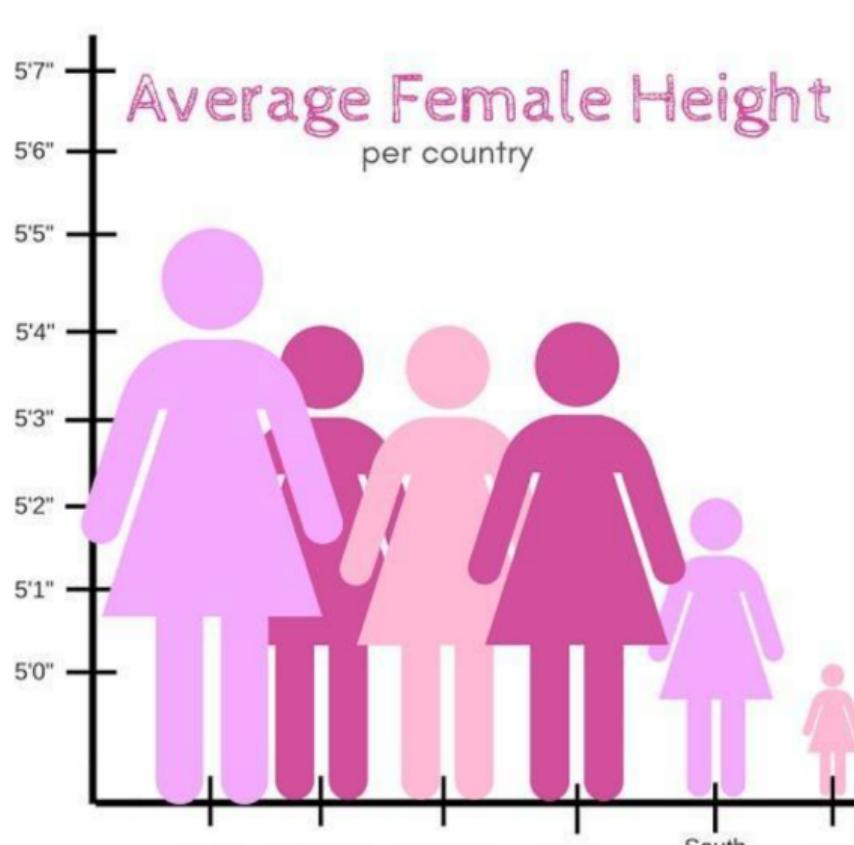
WTF?!



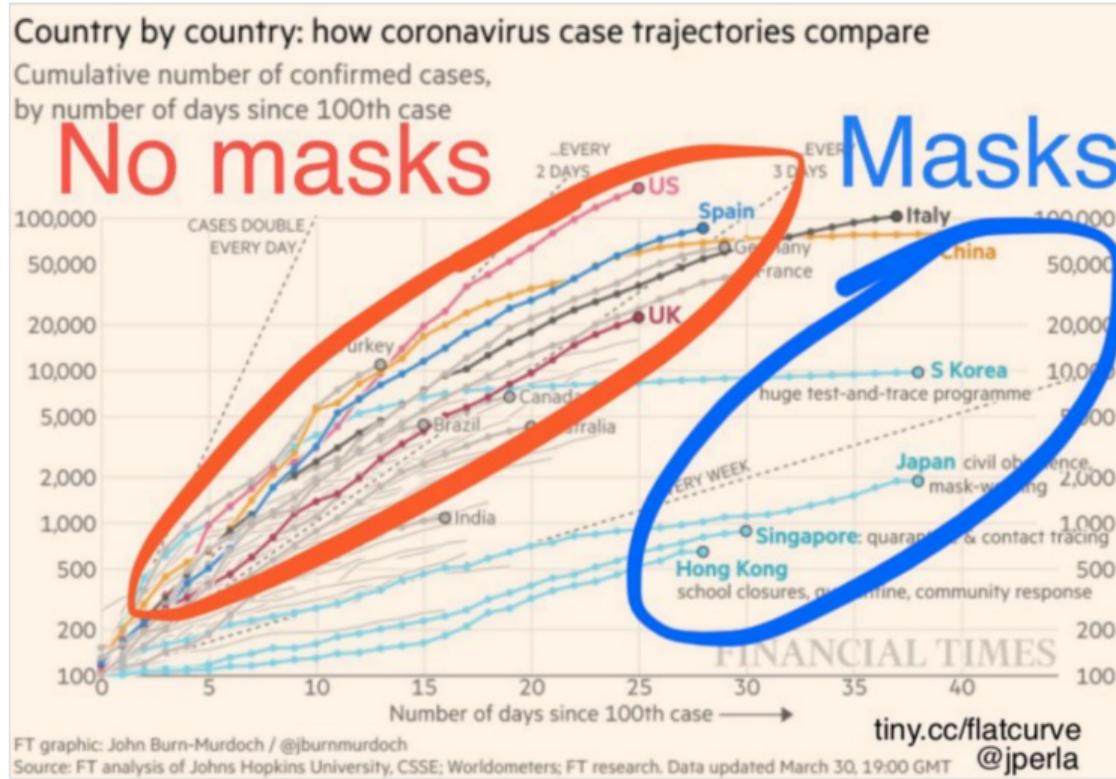
Height Differences by Country



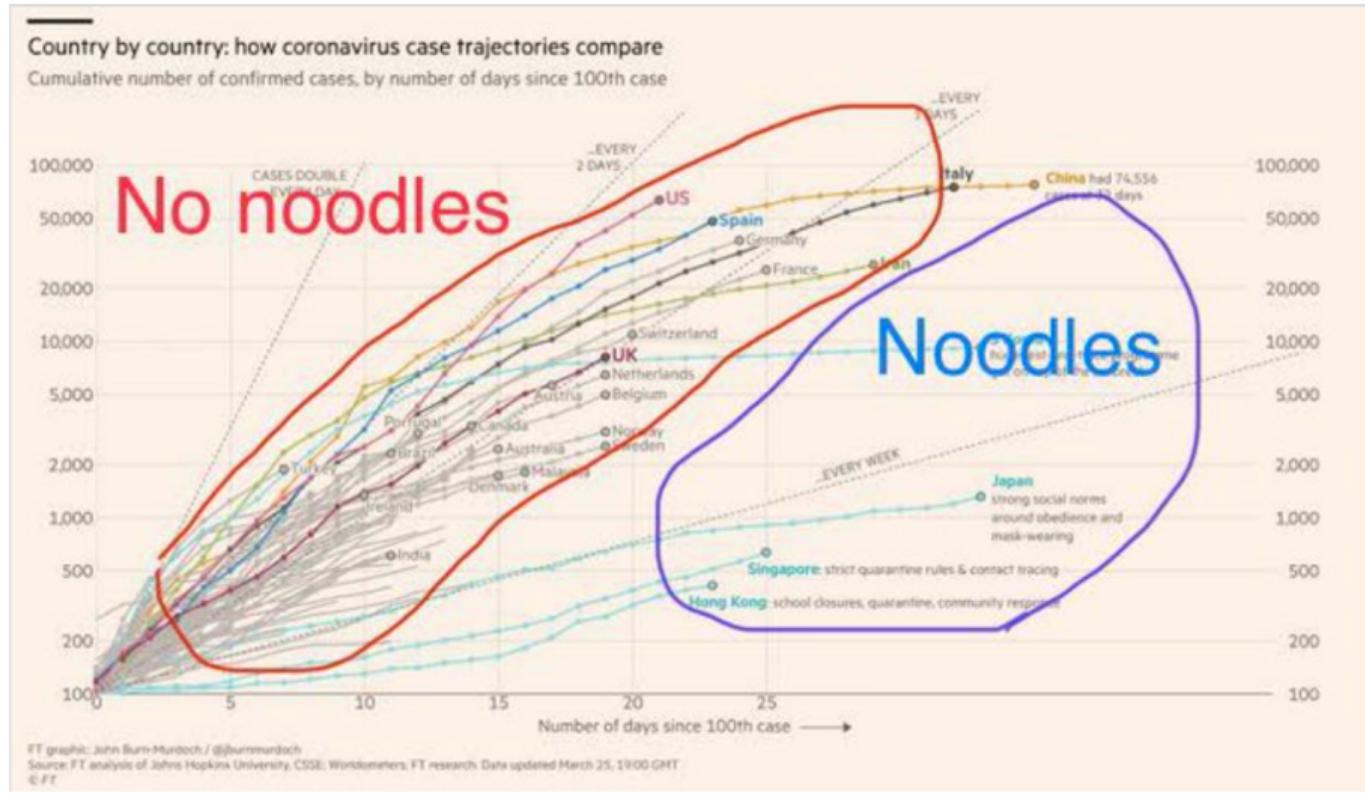
Height Differences by Country – Truncated *y*-axis, misleading area comparison.



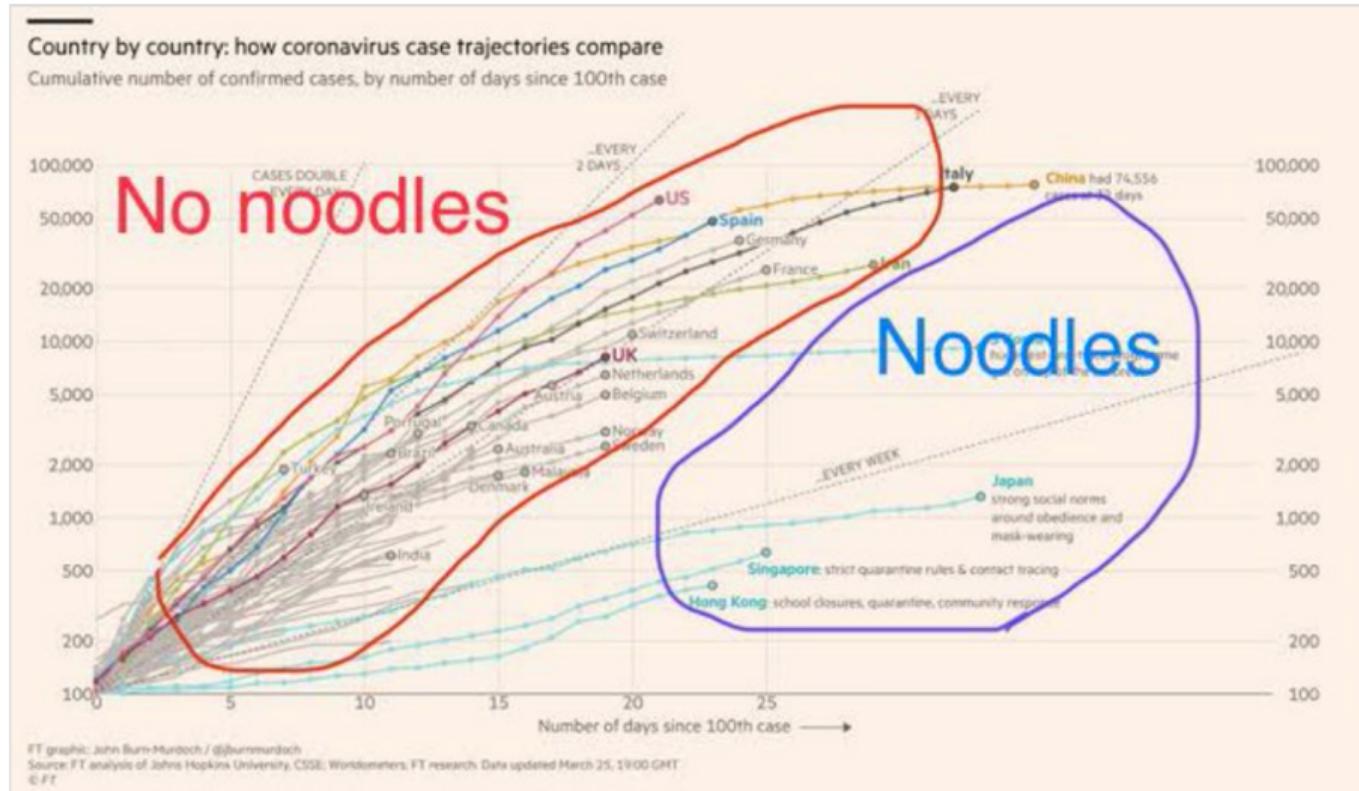
Causality?



Causality?



Causality – Falsely asserting causality.



Sino-Indian Diplomacy



Sino-Indian Diplomacy – Misleading area comparison.

