Mastering Applied Data Science

Day 2: Data Preparation and Exploration with Visualization, External Data

Acquisition



Agenda - Day 2



Introduction to Data Science

- Matplotlib, Pandas visualization
- Seaborn
- Tableau
- Supermarket Data visualization

Data Preparation and Exploration

- Feature Engineering
- Missing Value Imputation
- Creating Dummy Variables

External Data Acquisition

- Web Scraping
- API Pulls

What is Data Visualization?



Data visualization is the representation of data or information in a graph, chart, or other visual format.

It communicates relationships of the data with images.





- Words don't always paint the clearest picture. Raw data doesn't always tell the most compelling story.
- The human mind is very receptive to visual information.
 That's why data visualization is a powerful tool for communication.
- Visualization is important step in data exploration and data communication
- Gain insight into data that will guide analysis approaches
- Great technique to help explain patterns, trends, and correlations

Why Data Visualization?



A visual summary of information makes it easier to identify patterns, trends and outliers than looking through thousands of rows on a spreadsheet.

It's the way the human brain works.

Since the purpose of data analysis is to gain insights, data is much more valuable when it is visualized.

"A picture is worth a thousand words"



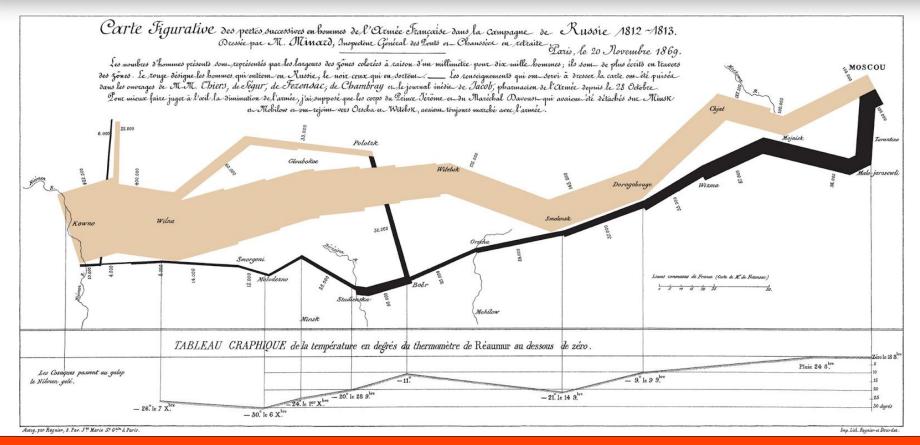
Our eyes are drawn to colors and patterns. We can quickly identify red from blue, square from circle. Our culture is visual, including everything from art and advertisements to TV and movies

Data visualization is another form of visual art that grabs our interest and keeps our eyes on the message.



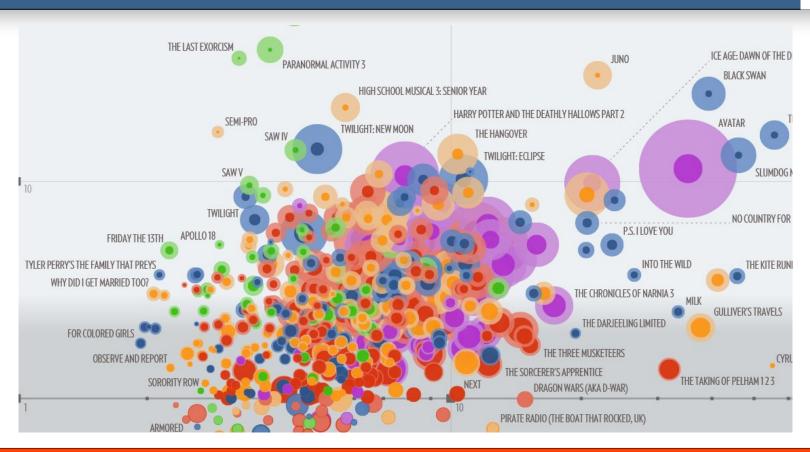






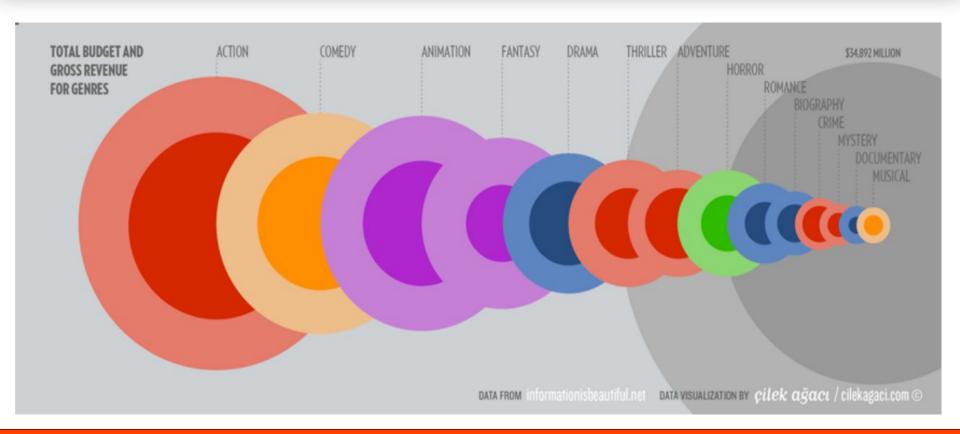
Hollywood Economics

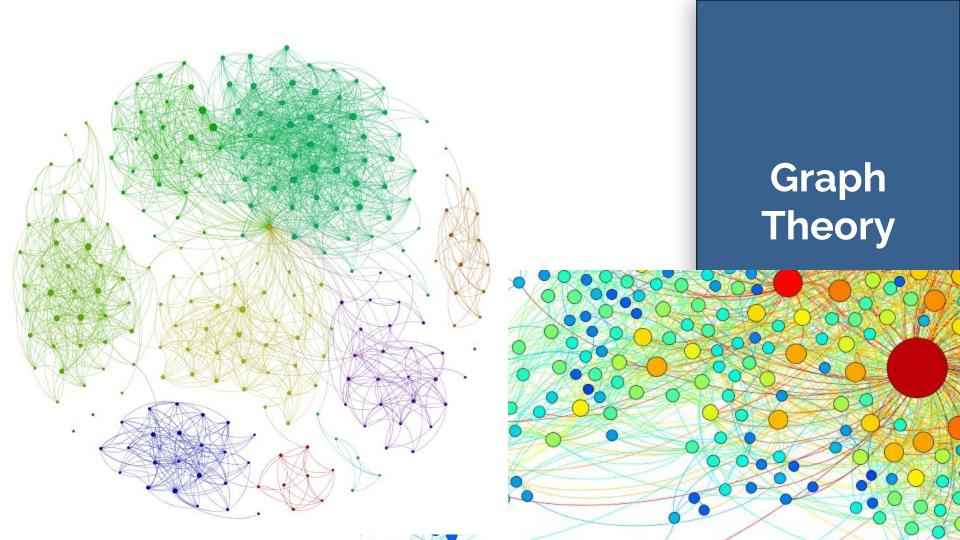




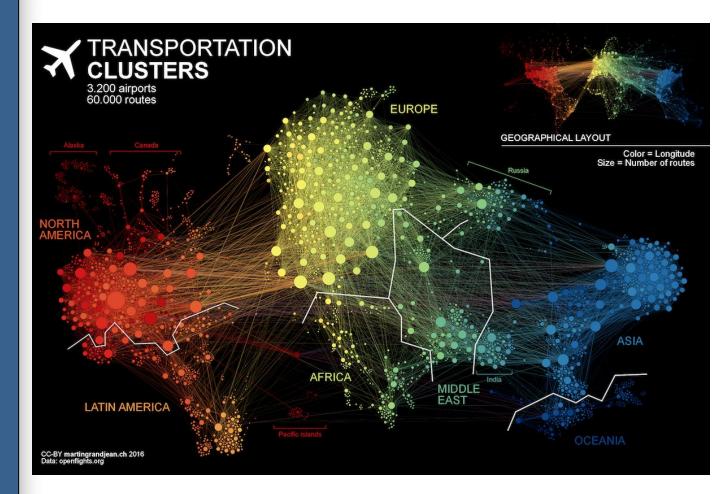
Hollywood Economics



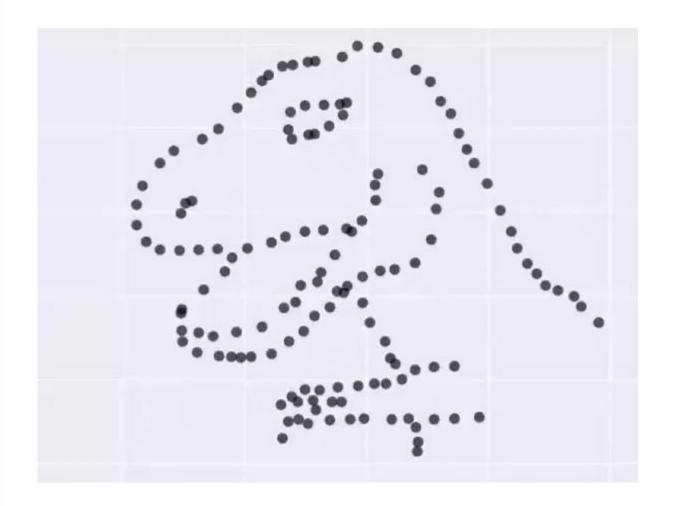




Air Traffic Control



Never trust
summary statistics
alone; always
visualize your data
Case in point - all datasets
below have the same
mean, standard deviation,
and correlation (both x/y)
to 2 decimal places, while
being drastically different.



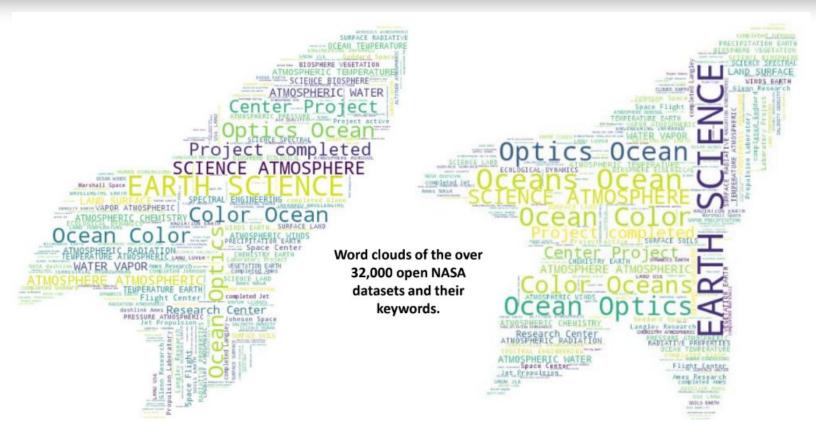
2021 This Is What Happens In An Internet Minute



What happens in an internet minute?

Word Cloud





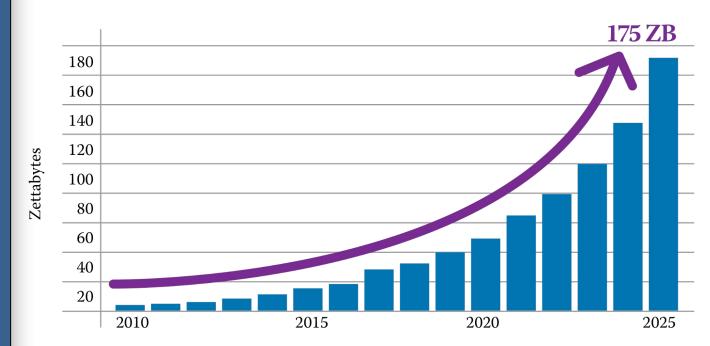
Word Cloud



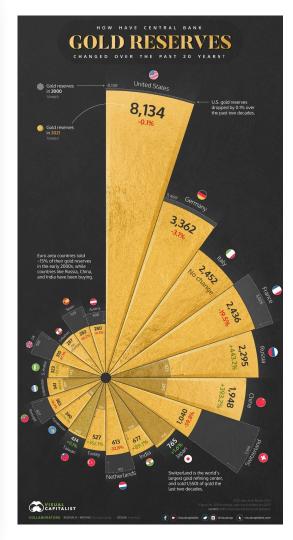




Massive Data Growth



Gold Reserves

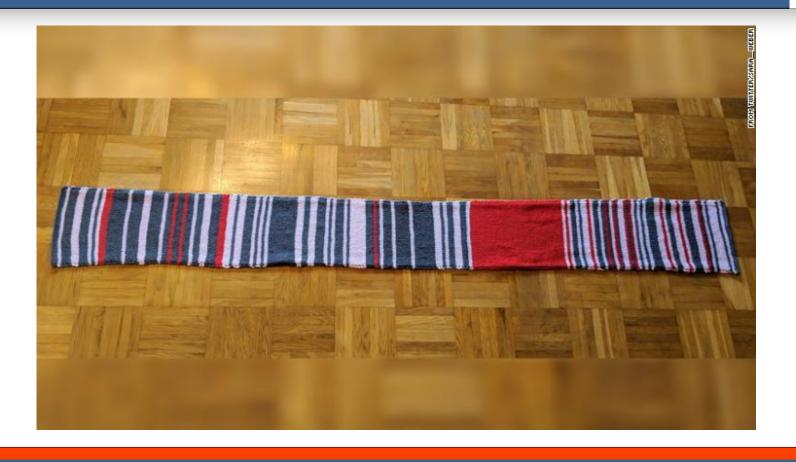


Data Science is Multidisciplinary



Computer **Domain Expertise Programming** Data Scientist Communication **Statistics** (Story-telling)

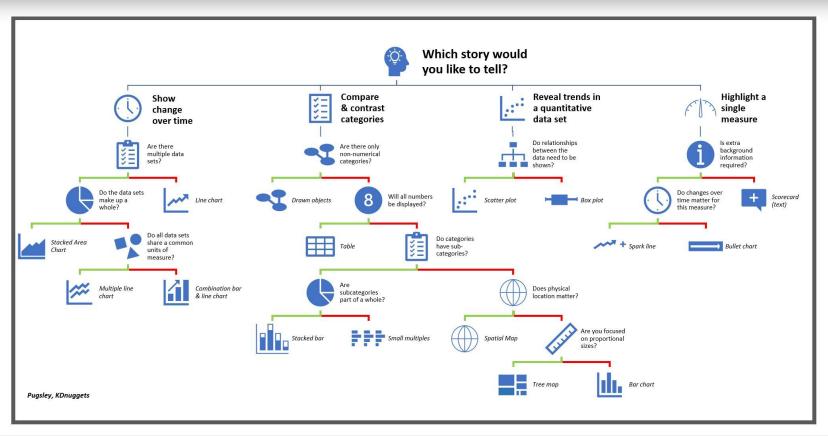




VARIABLE WIDTH TABLE WITH BAR CHART BAR CHART CIRCULAR AREA BAR CHART COLUMN CHART VERTICAL CHART VERTICAL EMBEDDED CHARTS HORIZONTAL LINE CHART LINE CHART Cyclical Data Non-Cyclical Single or Few Many Categories Two Variables Many Data Categories Few Categories Categories per Item One Variable BAR HISTOGRAM per Item Among Items Over Time Few Data SCATTER PLOT Single COMPARISON Variable LINE HISTOGRAM Variables What would you SCATTER PLOT RELATIONSHIP DISTRIBUTION like to show? BUBBLE SAZE Many Data Points: Three or more COMPOSITION Variables SCATTER PLOT Changing Static Over Time Few Periods Many Periods Relative and Only Relative Relative and Only Relative Simple Accumulation or Accumulation to Differences Absolute Differences Absolute Share of Subtraction Components total and absolute of Components Differences Matter Matter Differences Matter Total to Total difference matters STACKED 100% STACKED BAR STACKED AREA STACKED AREA PIE CHART WATERFALL STACKED 100% TREE MAP BAR CHART CHART 100% CHART CHART CHART BAR CHART WITH SUBCOMPONENTS

https://datavizcatalogue.com/





matpletlib

- Most popular graphing package in Python
- Visualization is important step in exploratory analysis and communicating findings
- We'll create five types of graphs here:
 - Line Plot
 - Scatter Plot
 - Histogram
 - Box plot

Visualization

Pandas Visualization

- Pandas is an open source high-performance, easy-to-use library providing data structures, such as dataframes, and data analysis tools like the visualization tools.
- Pandas Visualization makes it really easy to create plots out of a pandas dataframe and series. It also has a higher level API than Matplotlib and therefore we need less code for the same results.
- Here let's look at five common used graphs:
 - Scatter plot
 - Line chart
 - Histogram
 - Bar chart
 - Box plot

Pandas Visualization



Seaborn is a graphic library
 built on top of Matplotlib.

Visualization

It allows to make your charts

prettier, and facilitates some of the

common data visualisation needs

Any Questions?

Zafer Acar
Senior Data Scientist
zafer.acar@cern.ch
The Data Science Bootcamp

