

Day 2: Introduction to Data Science

Samuel Schlenker
04.11.2025, WWI 2025F





Students should ...

- Understand the fundamental definition and scope of data science as an interdisciplinary field
- Identify the three pillars of data science: domain expertise, statistics/mathematics, and computer science
- Recognize real-world applications of data science across industries (quality control, predictive maintenance, fraud detection, autonomous driving)
- Distinguish between AI, Machine Learning, Deep Learning, and Generative AI
- Understand the characteristics of Big Data (Volume, Velocity, Variety, Veracity, Value)
- Identify different data sources (open-source, private, commercial) and their accessibility
- Apply principles of effective data visualization and storytelling
- Recognize poor or misleading visualizations and understand common pitfalls
- Understand different chart types and their appropriate use cases (line charts, bar charts, scatterplots, heatmaps, etc.)
- Recognize the difference between correlation and causation
- Identify opportunities and risks associated with data science applications

Recommended Reading



Data Science
from Scratch

FIRST PRINCIPLES WITH PYTHON



Python for
Data Analysis

ADDISON WESLEY DATA & ANALYTICS SERIES 

DEEP
LEARNING
ILLUSTRATED

A Visual, Interactive Guide to Artificial Intelligence



JON KROHN

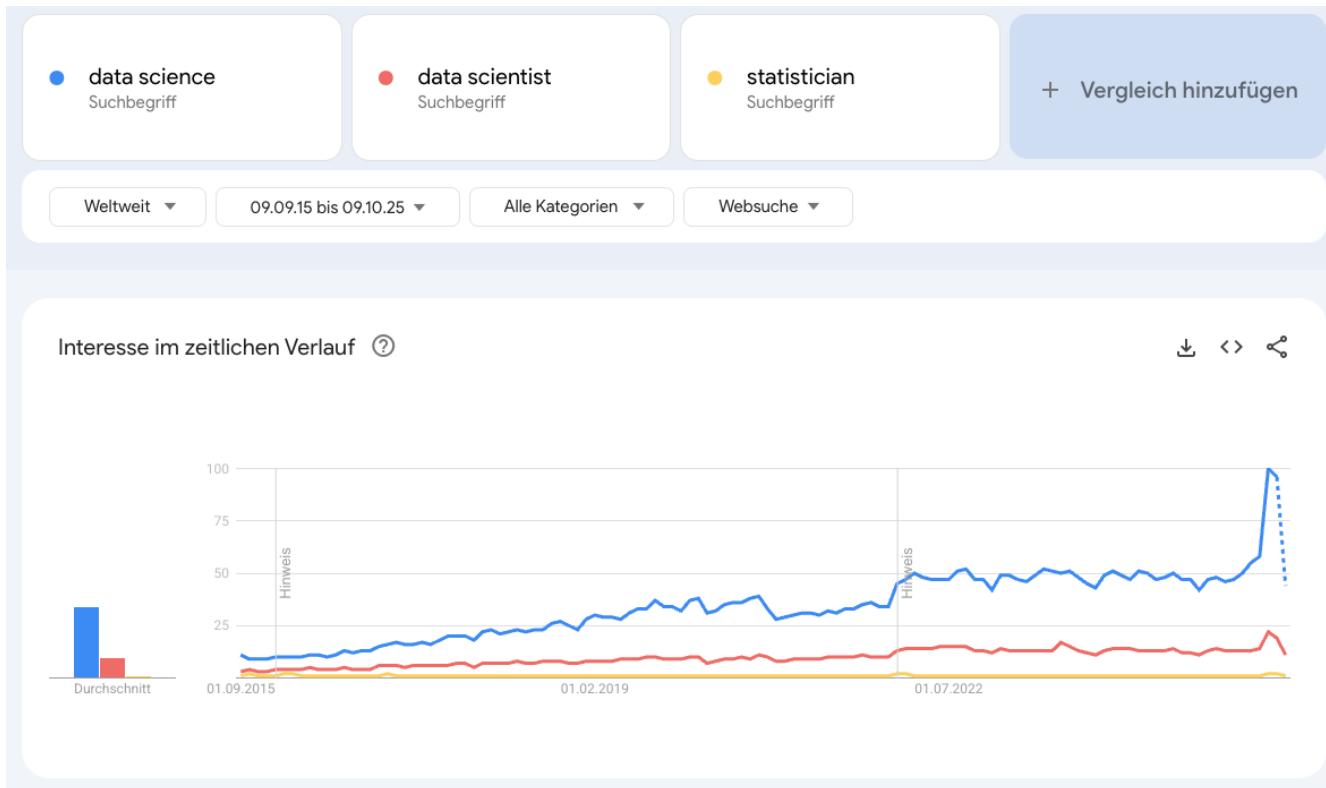
with GRANT BEYLEVELD and AGLAÉ BASSENS

Data Science from Scratch, by
Joel Grus. O'Reilly

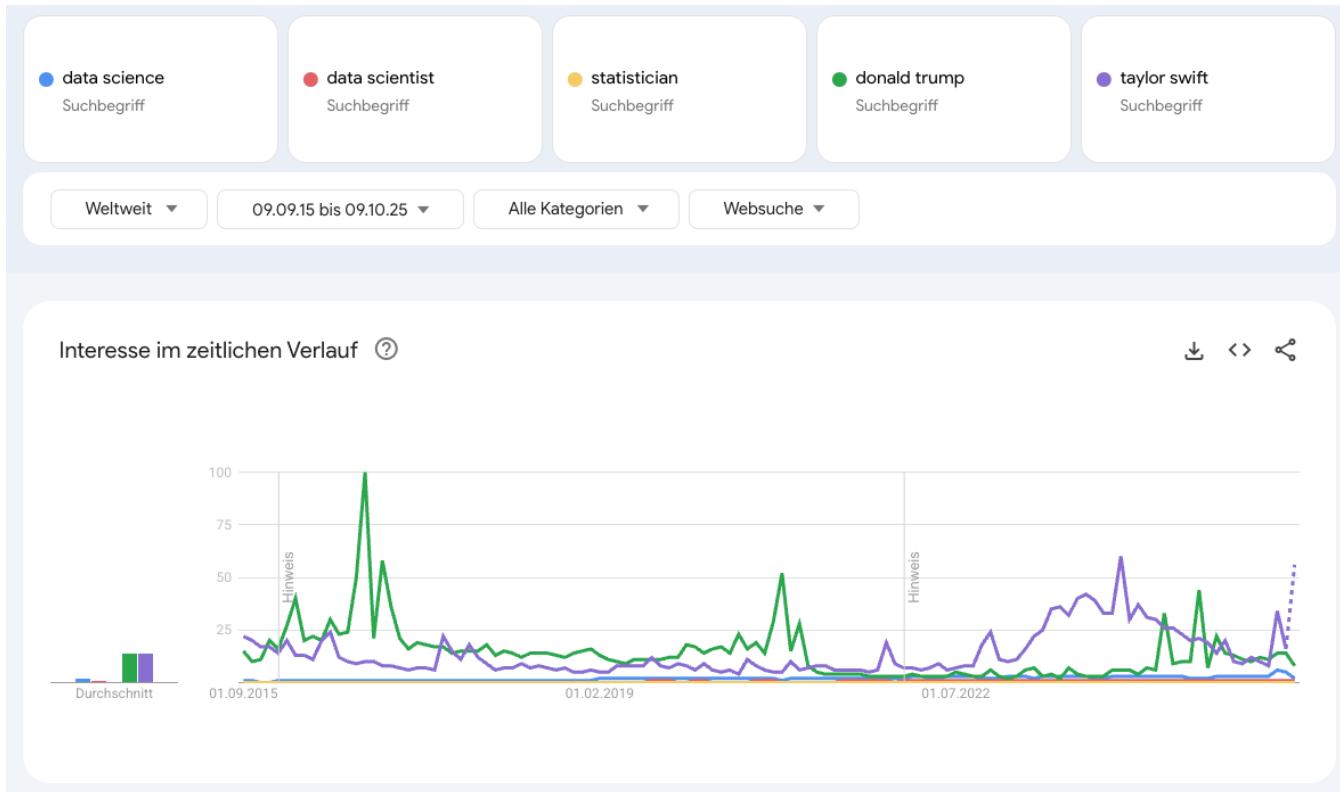
Python for Data Analysis, 2nd ed
by Wes McKinney. O'Reilly

<https://www.deeplearningillustrated.com/>

Google Trends



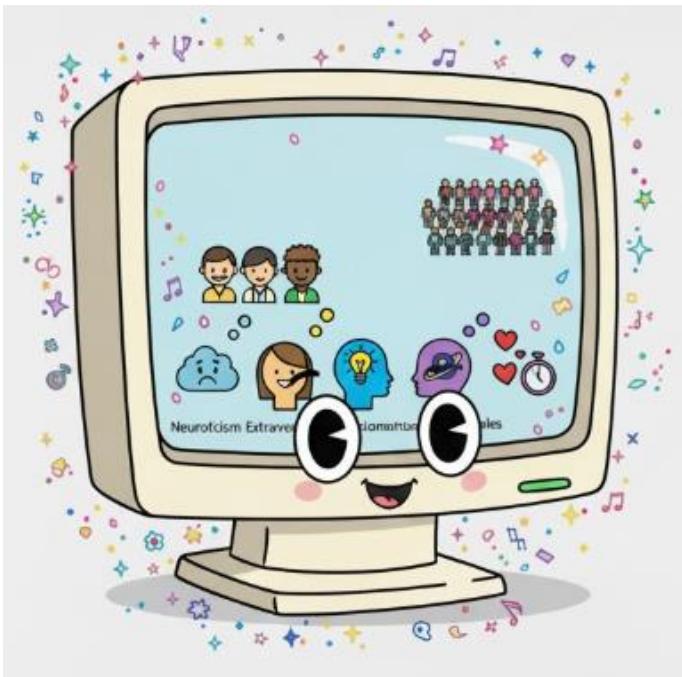
In comparison ...



Facebook knows us better than friends and family

- 2015: Study with 86,220 volunteers
- Collaboration between the University of Cambridge and Stanford University
- Questionnaire with 100 items on the Five Factor Model (FFM) of personality psychology / "Big Five"
- Neuroticism, extraversion, openness to experience, conscientiousness, and agreeableness
- Computer algorithm (linear regression) vs. assessment by individuals
 - From 10 likes: Computer is better than coworkers
 - From 70 likes: Computer is better than friends
 - From 150 likes: Computer is better than family
 - From 300 likes: Computer is better than spouse
- The average Facebook user shares 227 likes

Self-test: <https://applymagicsauce.com/demo>



As of 2011, the global size of data in healthcare was estimated to be

150 EXABYTES

[161 BILLION GIGABYTES]

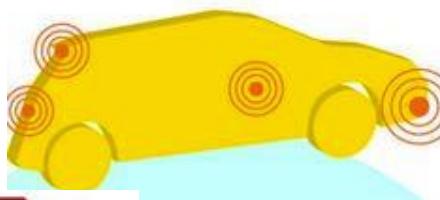
**30 BILLION
PIECES OF CONTENT**

are shared on Facebook every month



**4 BILLION+
HOURS OF VIDEO**

are watched on YouTube each month



Modern cars have close to
100 SENSORS

that monitor items such as fuel level and tire pressure

By 2016, it is projected there will be

**18.9 BILLION
NETWORK
CONNECTIONS**

– almost 2.5 connections per person on earth

400 MILLION TWEETS

are sent per day by about 200 million monthly active users

It's estimated that

2.5 QUINTILLION BYTES

[2.3 TRILLION GIGABYTES]

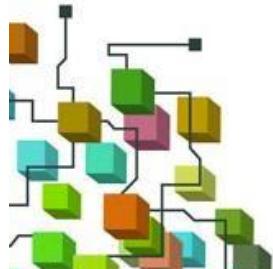
of data are created each day

Most companies in the
U.S. have at least

100 TERABYTES

[100,000 GIGABYTES]

of data stored



40 ZETTABYTES

[43 TRILLION GIGABYTES]

of data will be created by
2020, an increase of 300
times from 2005

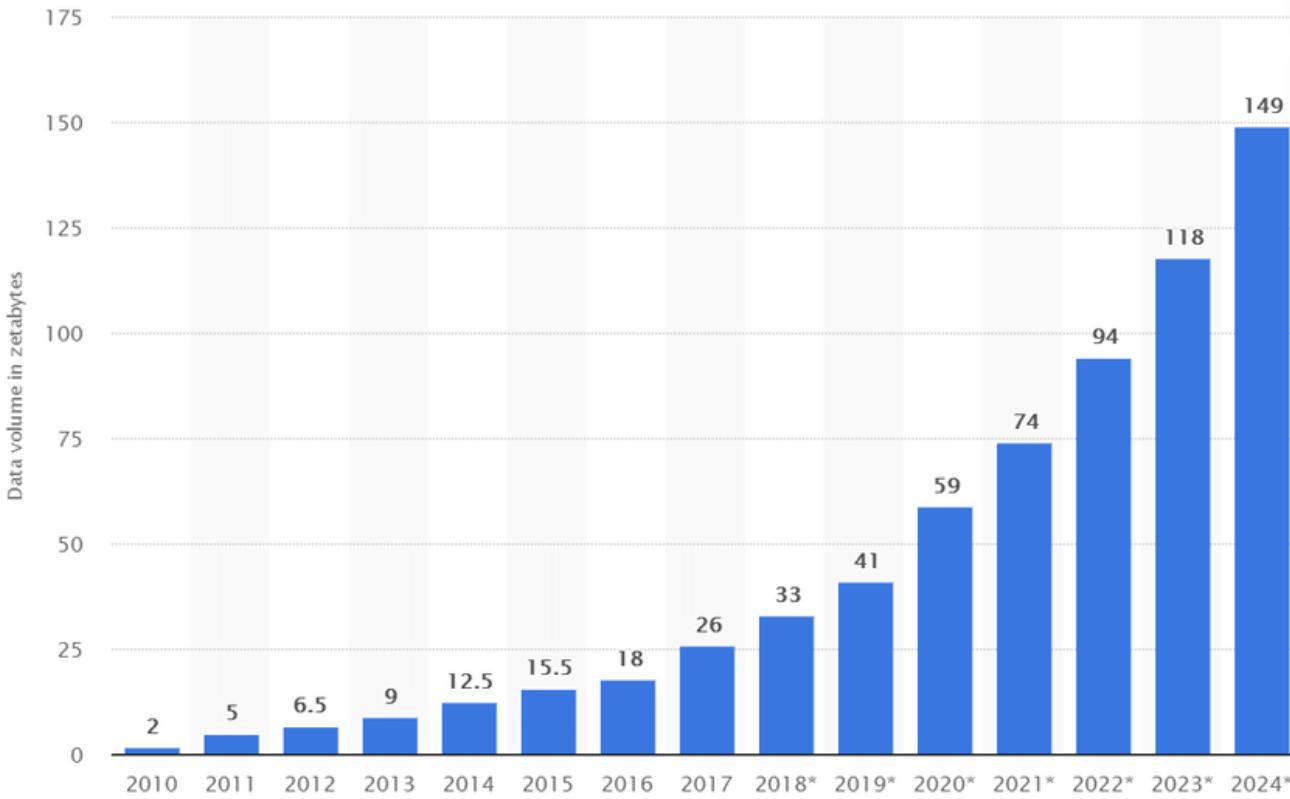
The New York Stock Exchange
captures

**1 TB OF TRADE
INFORMATION**

during each trading session



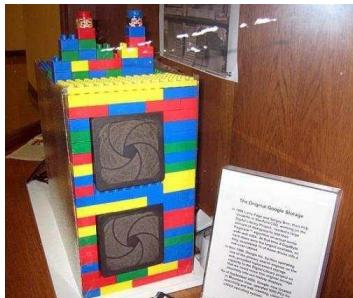
Pure Amount of data



Distributed Computing

This amount of data requires more than one computer.

Google – 1998:



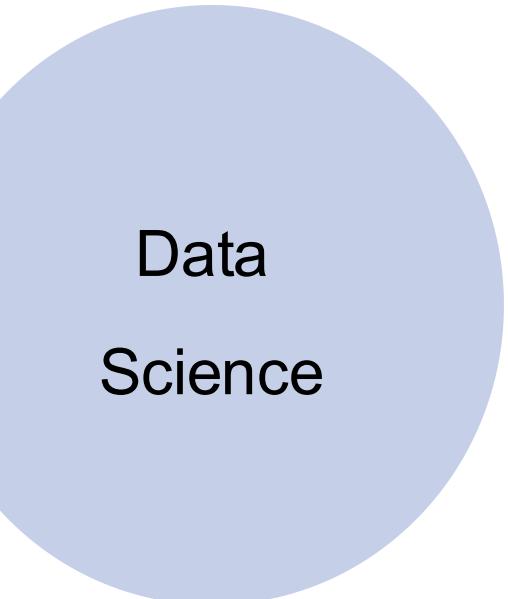
Distributed Computing

This amount of data requires more than one computer.

Google – 2014:





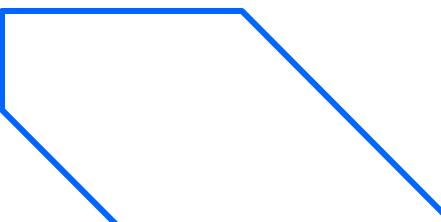


Data
Science

Needs massive data processing



Distributed computing



Why is this the case?

Data Is Driving Everything

- Modern data acquisition is inexpensive!
 - Smartphones, embedded systems, inexpensive sensors,
 - Medical devices, simulators, ...
- Data storage is inexpensive!
- Parallel (compute cluster) computation is inexpensive
 - The Cloud, clusters of computers, GPUs, tensor processors, ...
- Science only has explanatory and predictive models in a few (mostly physical sciences-related) domains
- ... So: can we use algorithms + data to understand phenomena? Build or augment models? Build detectors? Make diagnoses?



Data Is Driving Everything

“Big data”

“Data science”

“Data lakes”

“Visual analytics”

“Deep learning”

“Statistical analysis”

“Biomedical informatics”

“Business analytics”

Lots of trends in pursuit of the same goals!
Discovery, models, decision-making, ...

Also, new issues -
“Ethical algorithms”
“Reproducibility”

Data Science

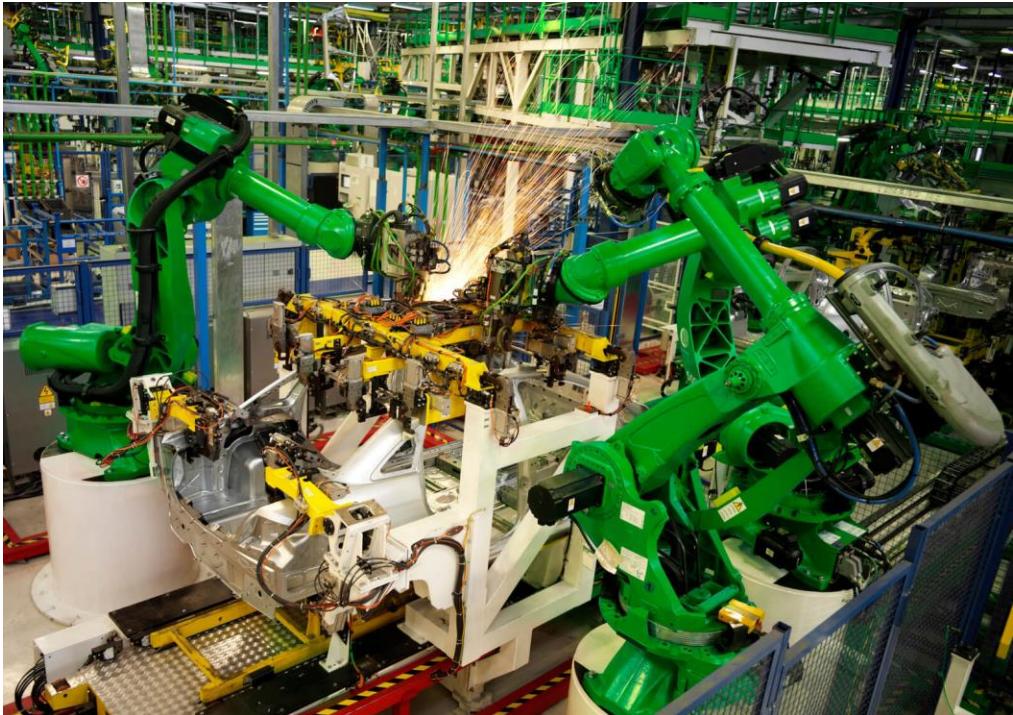


Reccomandation engines

Autonomous driving

Business intelligence

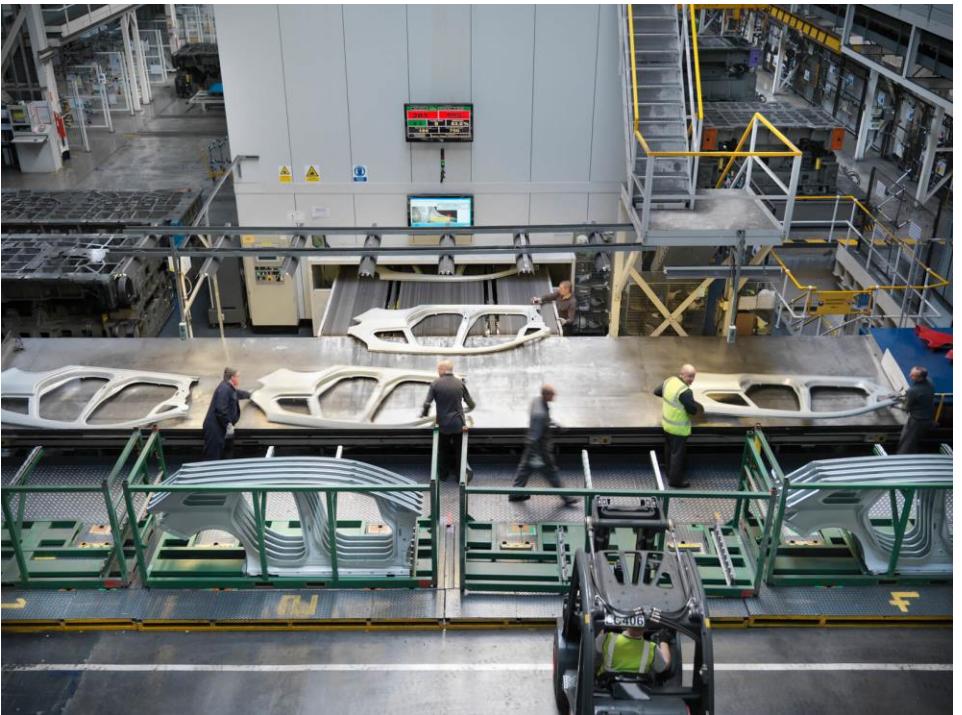
Applications of data science (Use cases)



Quality control:

- Quantity: Is there the right number of salami slices on the frozen pizza?
- Correct product: Are the right components installed in the server?
- Quality: Are there scratches in the paintwork on the car door? Is the weld seam flawless?
- Results in: Fewer complaints, Cost reduction, Higher customer satisfaction

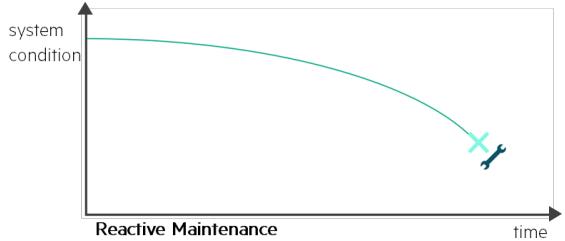
Applications of data science (Use cases)



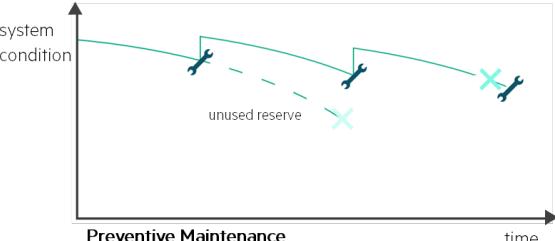
Predictive Quality: Predicting the quality of the product before the production process is complete.

- Early correction of the product, if possible, or discontinuation of production.
- Reduction of waste.
- Energy, water, and material savings by avoiding finishing steps for poor-quality products.
- Cost reduction.

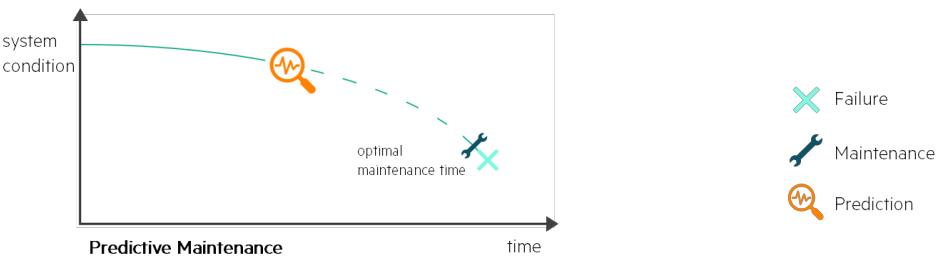
Applications of data science (Use cases)



Reactive Maintenance



Preventive Maintenance



Predictive Maintenance

- ✖ Failure
- 🔧 Maintenance
- 🔍 Prediction

Predictive Maintenance:

Predicting the failure of a production machine

- Avoiding unplanned production downtime
- Replacing/repairing the component that is about to fail at the best possible time
- Using the production machine or individual parts for as long as possible

Applications of data science (Use cases)



Fraud Detection:

Detection of credit card fraud, for example, through anomaly detection

- Early account blocking
- Loss minimization

Applications of data science (Use cases)

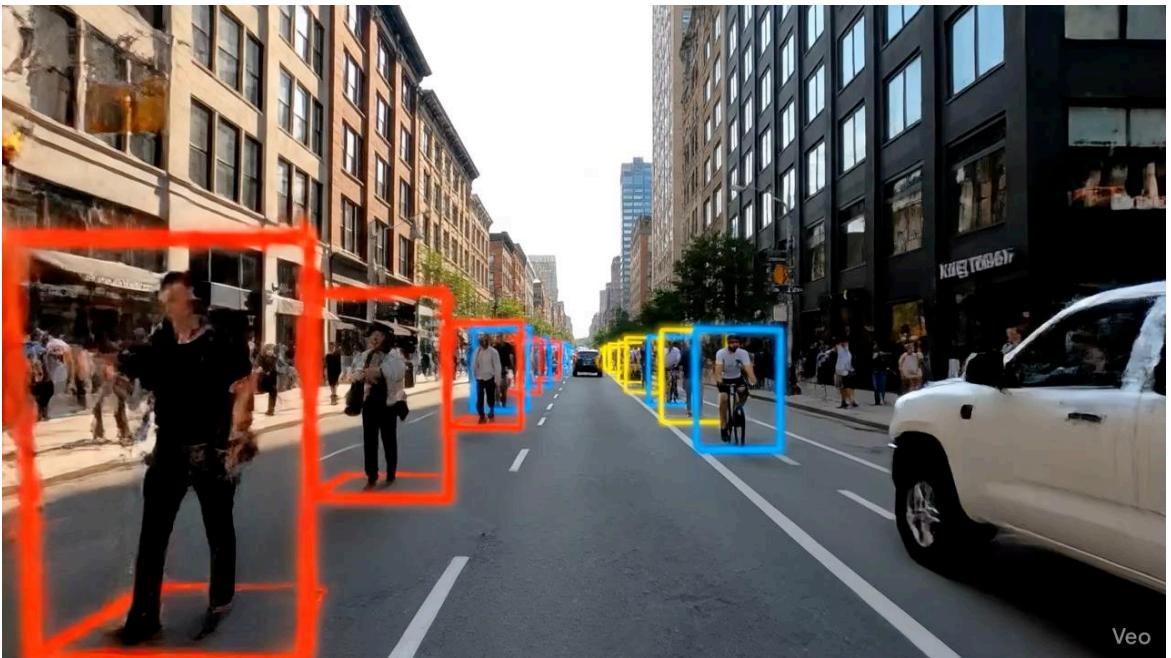


Personalized recommendations:

For example, shopping cart analysis

- Individualized shopping experience
- Higher customer loyalty
- Higher shopping cart values
- Increased sales

Applications of data science (Use cases)



Autonomous driving:

Detection of pedestrians, their direction of travel, and derivation of the appropriate action

- Traffic sign recognition
- Enables innovation
- Minimization of human error
- Reduction in the number of accidents

Applications by industry

Financial Services

- Fraud Detection
- Risk Assessment
- Algorithmic Trading
- Sentiment Analysis

Public Sector

- Urban Planning & Design
- Predictive Policing
- Citizen Engagement
- Document Summarization

Telecommunications

- Customer Support
- Network Optimization
- Predictive Maintenance
- Fraud Detection

Healthcare

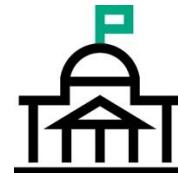
- Gene Research
- Drug Discovery
- Patient Assistant
- Medical Records

Manufacturing

- QC / Defect Detection
- Supply Chain Optimization
- Predictive Maintenance
- Generative Design



Financial Services
& Insurance



Public Sector /
Defense



Healthcare /
Life Sciences



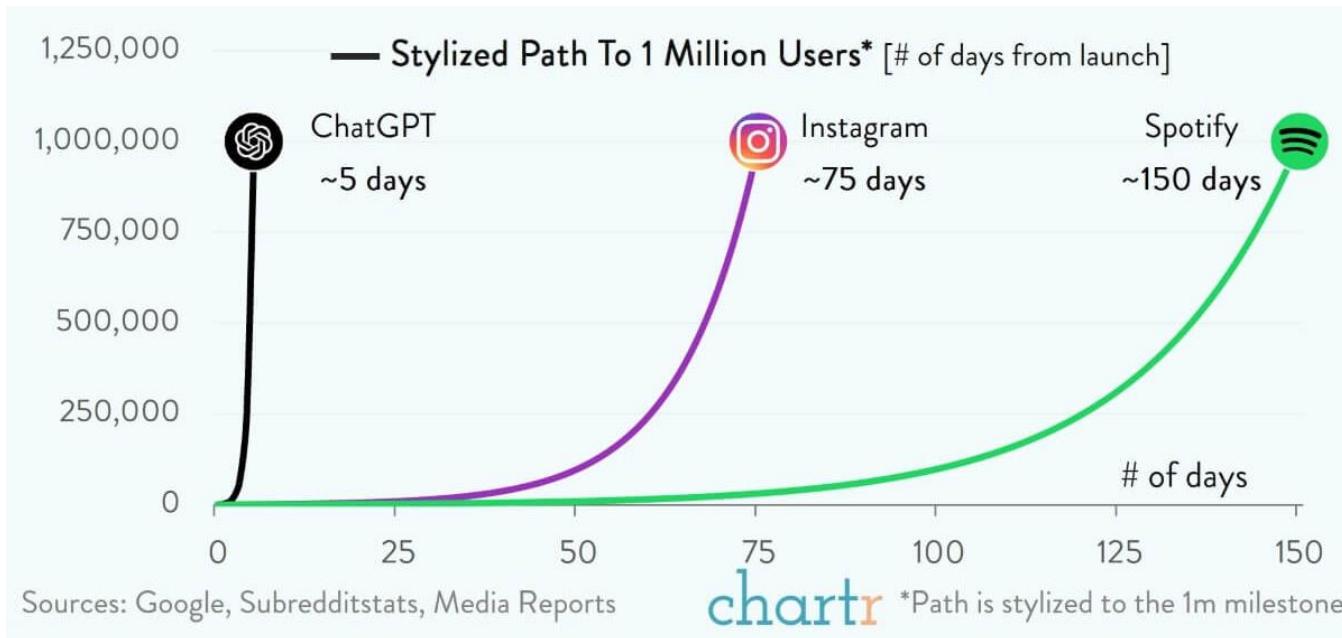
Autonomous
Driving



Manufacturing

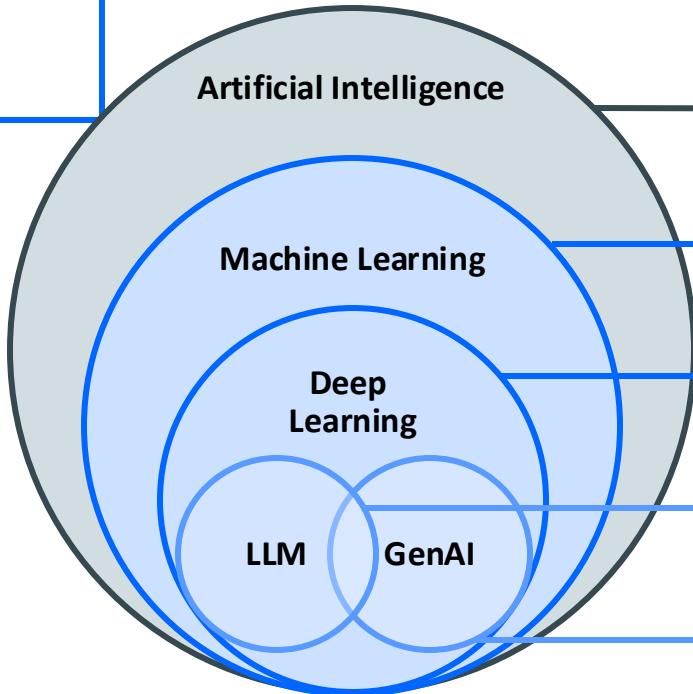
AI goes mainstream with ChatGPT

1M users in 5 days



But how does AI fit into the picture?

Navigating Artificial Intelligence: Understanding its fundamental building blocks



Artificial Intelligence (AI)

Any technology that enables machines to solve tasks in a way like humans do

Machine Learning (ML)

Algorithms that allow computers to learn from examples without being explicitly programmed (supervised & unsupervised)

Deep Learning (DL)

Using deep artificial neural networks as models, inspired by the structure and function of the human brain

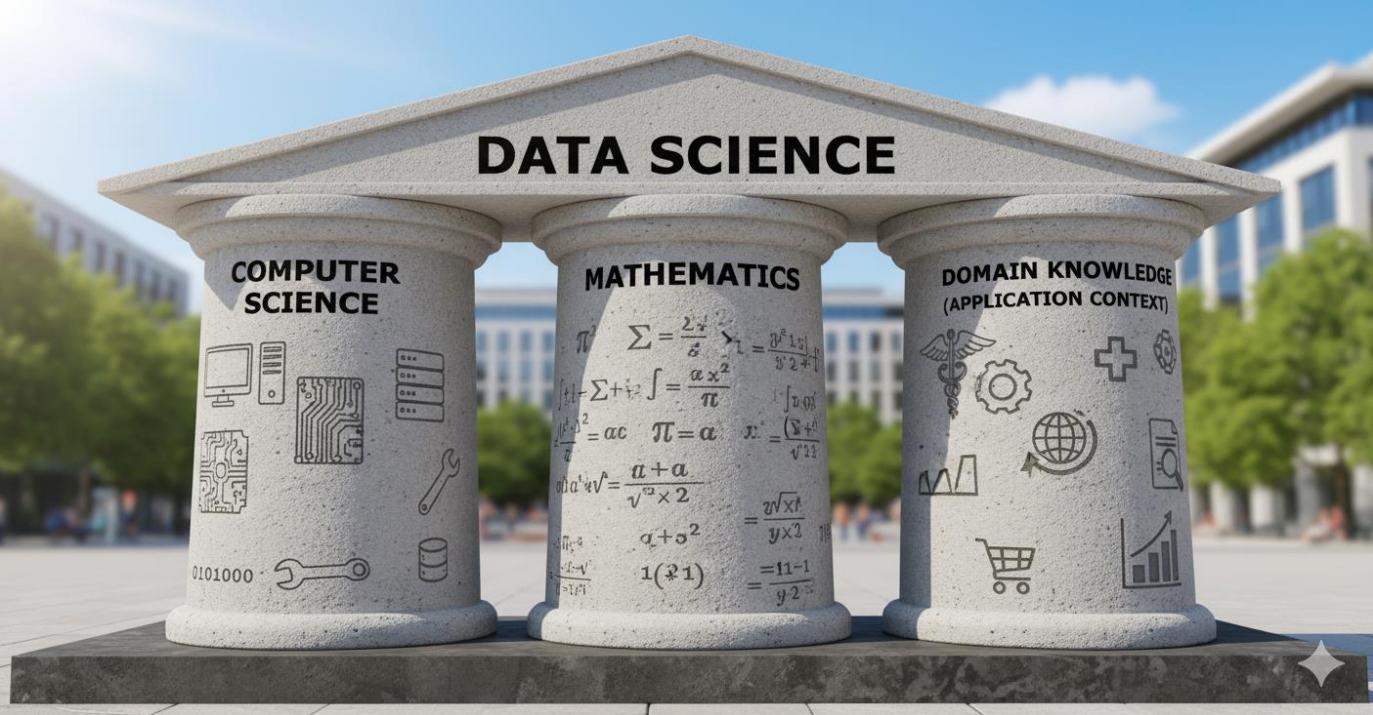
Large Language Models (LLM)

Models trained on massive datasets to understand and generate human-like text across diverse subjects

Generative Artificial Intelligence (GenAI)

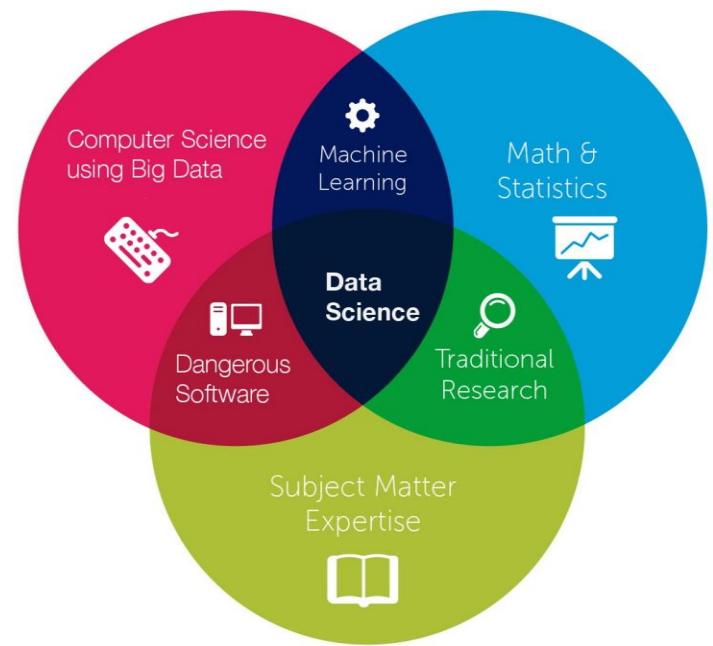
Refers to technologies that utilize machine learning models to generate human-like text, images, or other content

The three pillars of data science



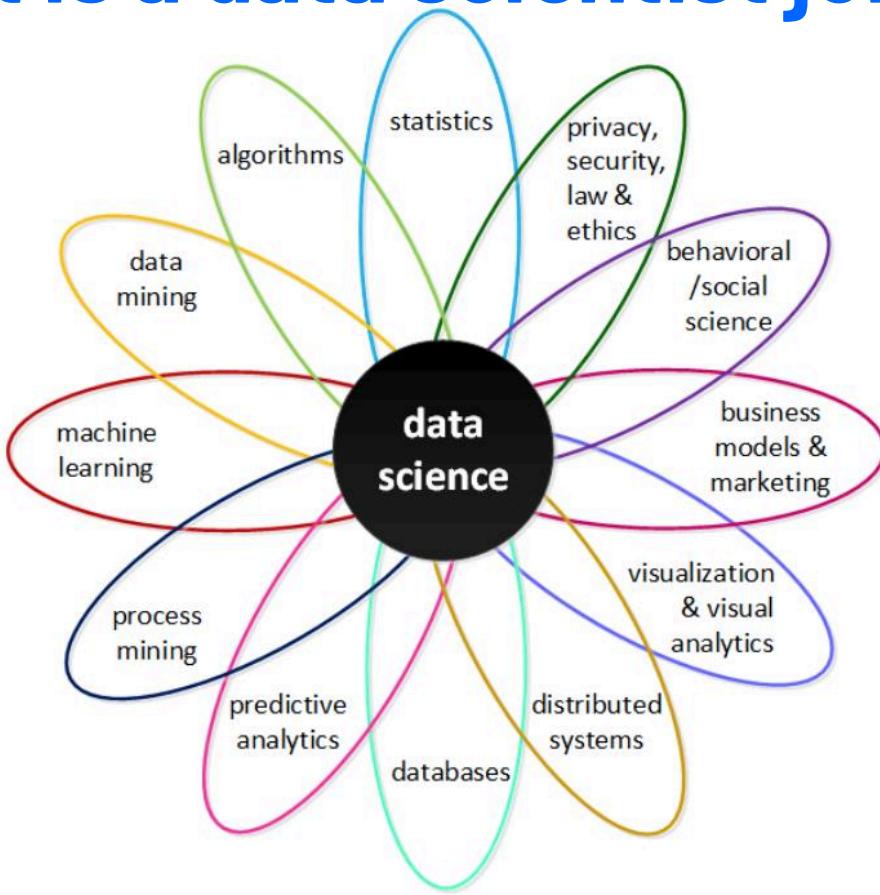
Why so many disciplines?

- Enormous amounts of data → Big Data
 - New algorithms. Must be scalable and efficient
 - New data management technologies and concepts
- High dimensionality of data
 - Data can have thousands of dimensions/attributes
 - High complexity and variability of data
- New, complex fields of application
- Data science is not (only) out-of-the-box, but domain-specific



https://miro.medium.com/max/1100/1*aXJWLmf-CYqVTrNiE2pdCw.png

So, what is a data scientist job?



And more concrete ... ?

- Data Steward
- Monitoring data quality and integrity
- Responsible for the technical accuracy of data

Tasks:

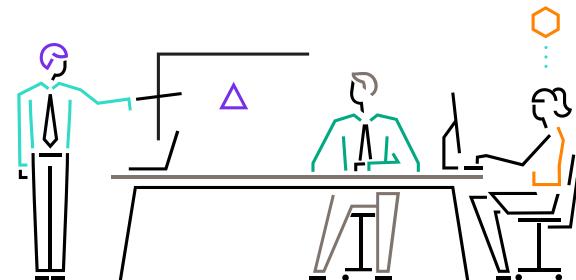
- Data owner
- Data governance
- Master of data

- Data Engineer
- Data supply
- Develop, implement, test, and maintain architecture for data storage

Tasks:

- Data Quality Improvement
- Data Transformation

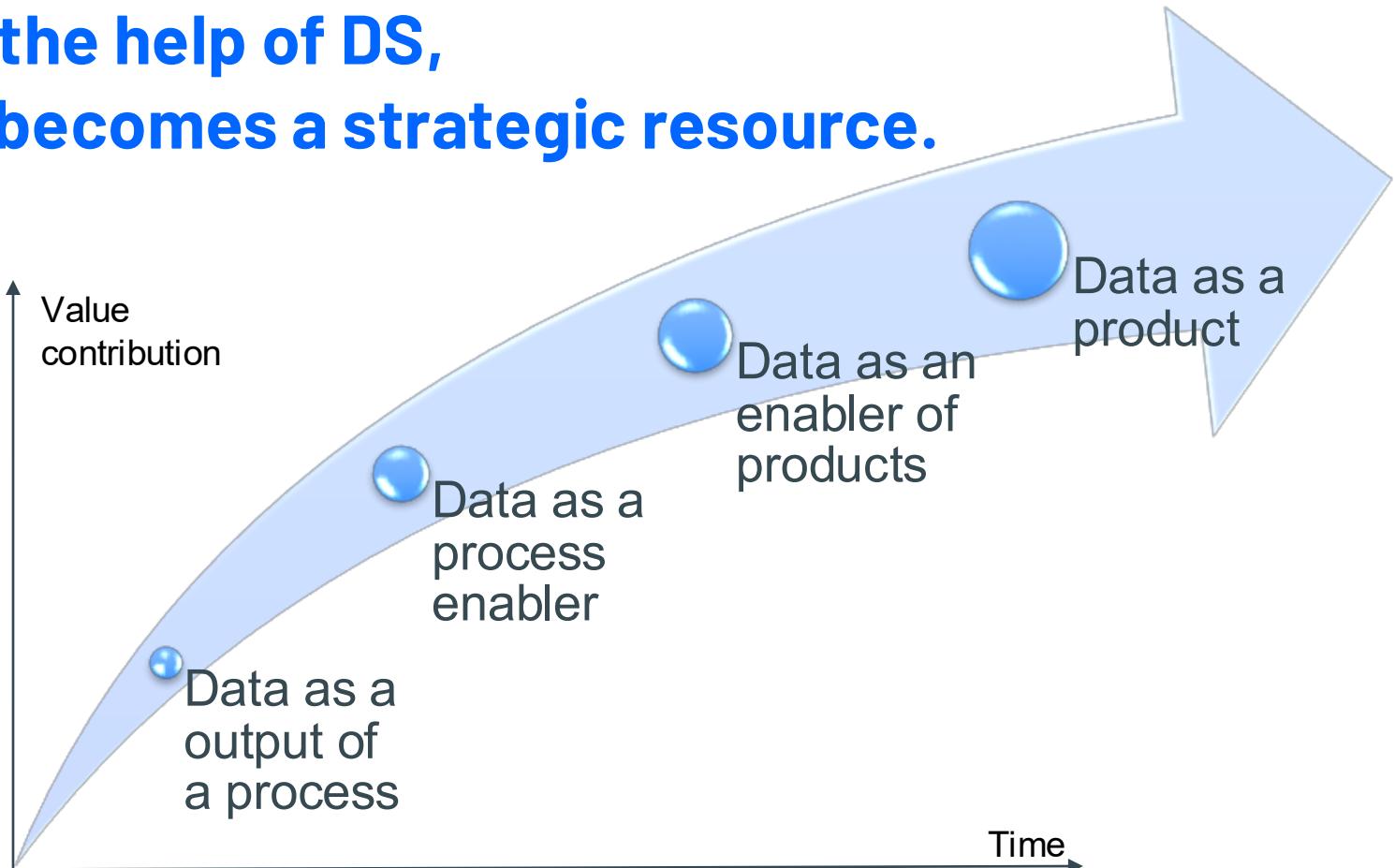
- Data Scientist
 - Provides answers to analytical questions using data
- Tasks:**
- Data Exploration & Visual Analytics
 - Model Deployment & Scoring
 - Big Data Handling & Manipulation



Big Data



With the help of DS, data becomes a strategic resource.

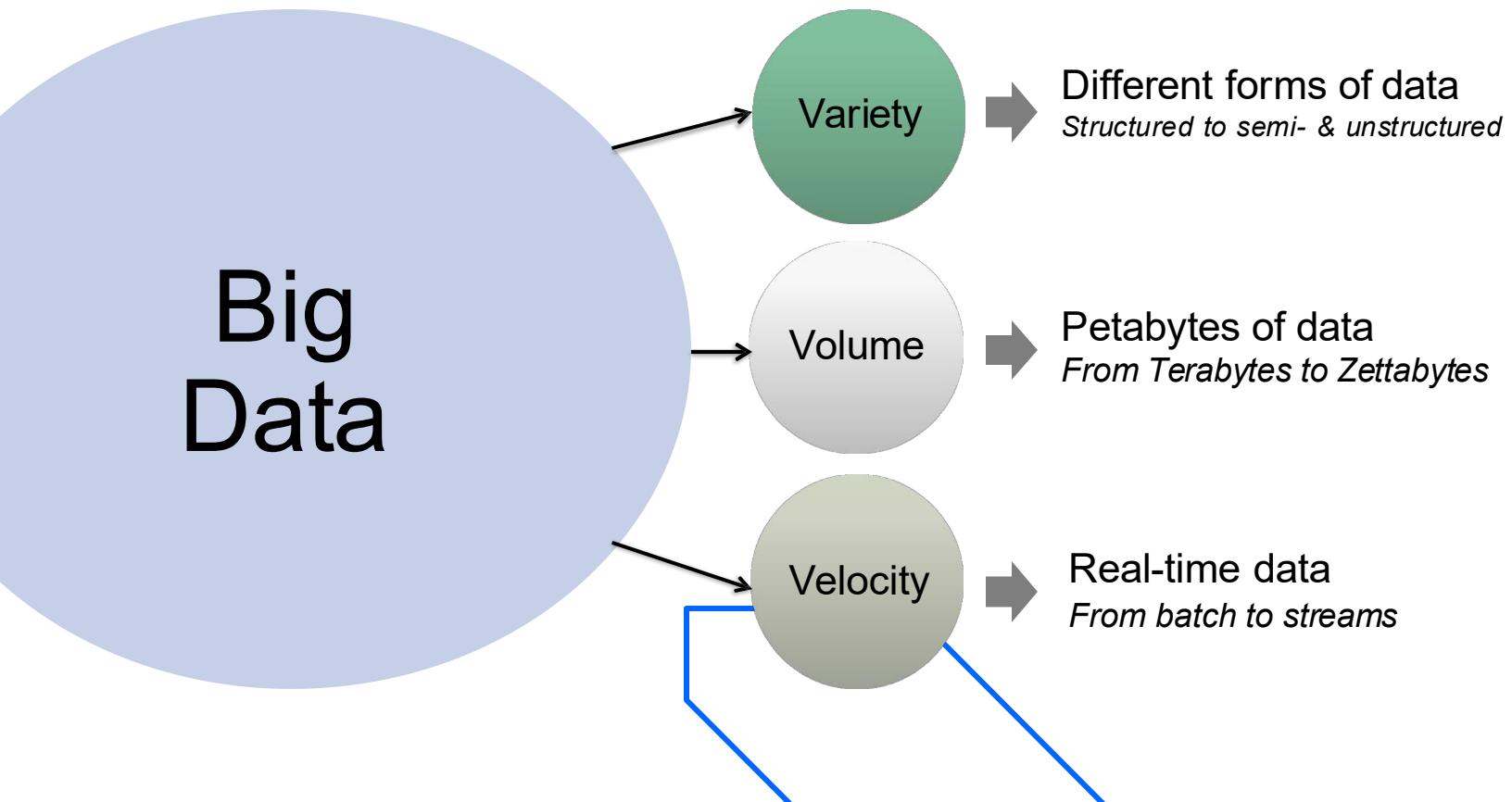


Big Data

Refers to datasets whose size is beyond the ability of typical database software tools to capture, store, manage and analyze.

Big does not stand primarily for size, but as an analogy for “overwhelming”

Big can mean “high variety”, “high volume” or “high velocity”



More Vs



Value

Turning data into
knowledge



Variability

Variation in meaning
in different contexts



Veracity

Uncertainty of
the data

- Not easy to measure
- Depend on context and intended use

Where is the data coming from?

Possible data sources

- Questionnaires, interviews, ...
- Observations & measurements (e.g. in production)
- Analysis of content (documents, web scraping, social media...)

Raw data: Data received or collected

- no variables have been manipulated
- no data removed from the record
- no summary / aggregation

Not prepared
(pre-processed)



Where is the data coming from?

Private data

Created by customers
Created during business
process execution

Commercial data

Cloud Marketplaces (e.g.,
AWS Data)
Qlik DataMarket, Statista

Open-source data

Data that is
publicly available (check
for limits
on usage)

Open-Source Data Sources examples

[Kaggle](#)

[World Health Organization](#)

[Our World in Data](#)

[Census Bureau \(U.S.\)](#)

[National Oceanic and Atmospheric Administration \(U.S.\)](#)

[UC Irvine Machine Learning Repository](#)

[Harvard Dataverse](#)

[AWS, Facebook, Google, Microsoft, ...](#)

CAN YOU GET THE DATA OUT OF SILOS?

A photograph of a rural landscape under a hazy, overcast sky. In the foreground, there's a field of tall, dry grass. In the middle ground, several large, cylindrical metal grain silos are scattered across the horizon. The silos are partially obscured by the fog. The overall atmosphere is one of isolation and separation, which serves as a metaphor for data silos.

Accounting

Marketing

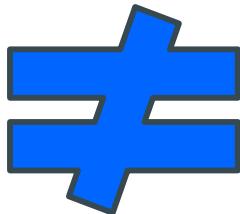
Manufacturing

After-Sales
Service

R&D

Customer
Relations

Value



Quality

While value and quality of big data may be correlated, they are conceptually different. For example, one can have high quality data about the names of all the countries in North America, but this list of names may not have much perceived value. In contrast, even relatively incomplete data about the shopping habits of people can be quite valuable to online advertisers.

Xin Luna Dong, Google, 2015

Raw data (including high quality data) itself does not hold any value, unless it is processed in analytical tasks from which humans or downstream applications can derive insight.

Gerhard Weikum, Max Planck Institute for Informatics, 2015

Data Storytelling and Visualization

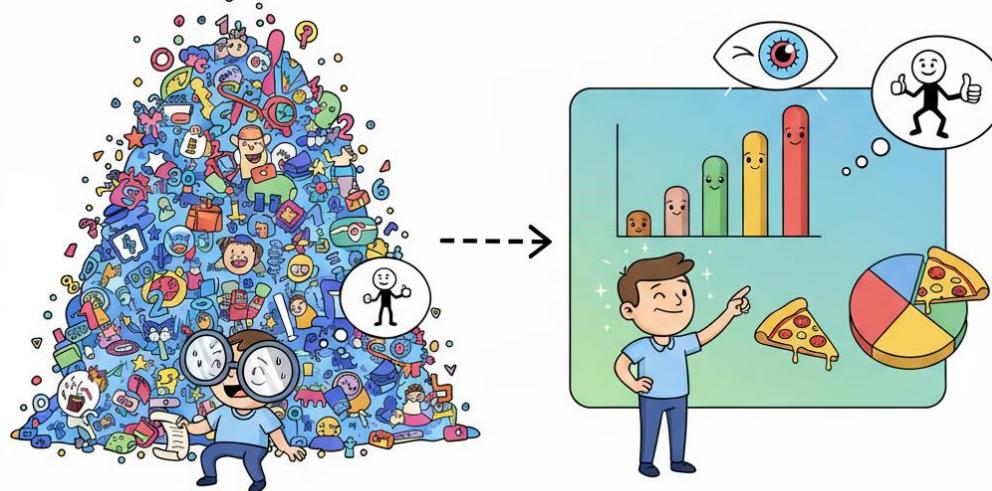


Data story telling



Why do we visualize?

- Most datasets are far too large to be examined in their raw format.
- Visual analysis uses our **pre-attentive perception** - visual cues that humans process automatically and unconsciously.
- We can perceive and interpret these types of characteristics quickly and without any special effort.
- Example: Using the length of bars to represent sales volumes is an effective choice to indicate differences in sales between categories.



Possibilities of Visual Presentation

Length

- Very good for quantitative variables
- Poorly suited for qualitative variables



Possibilities of Visual Presentation

Width

- Limited to quantitative variables
- Poorly suited for qualitative variables



Possibilities of Visual Presentation

Orientation

- Limited to quantitative variables
- Poorly suited for qualitative variables



Possibilities of Visual Presentation

Size

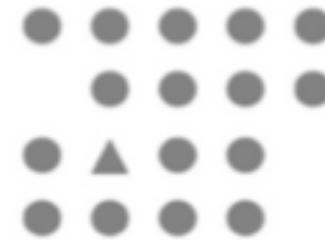
- Limited to quantitative variables
- Poorly suited for qualitative variables



Possibilities of Visual Presentation

Form

- Poor for quantitative variables
- Very suitable for qualitative variables



Possibilities of Visual Presentation

Position

- Very good for quantitative variables
- Poorly suited for qualitative variables



Possibilities of Visual Presentation

Grouping

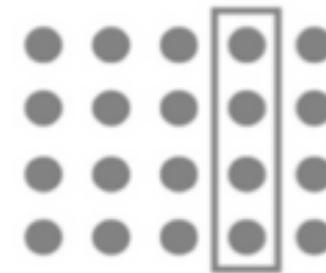
- Limited to quantitative variables
- Very suitable for qualitative variables



Possibilities of Visual Presentation

Containment

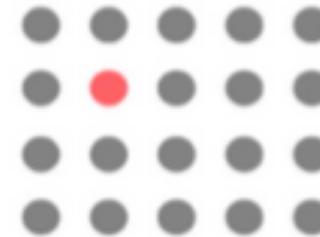
- Poor for quantitative variables
- Limited suitable for qualitative variables



Possibilities of Visual Presentation

Hue

- Poor for quantitative variables
- Very suitable for qualitative variables



Possibilities of Visual Presentation

Colour intensity

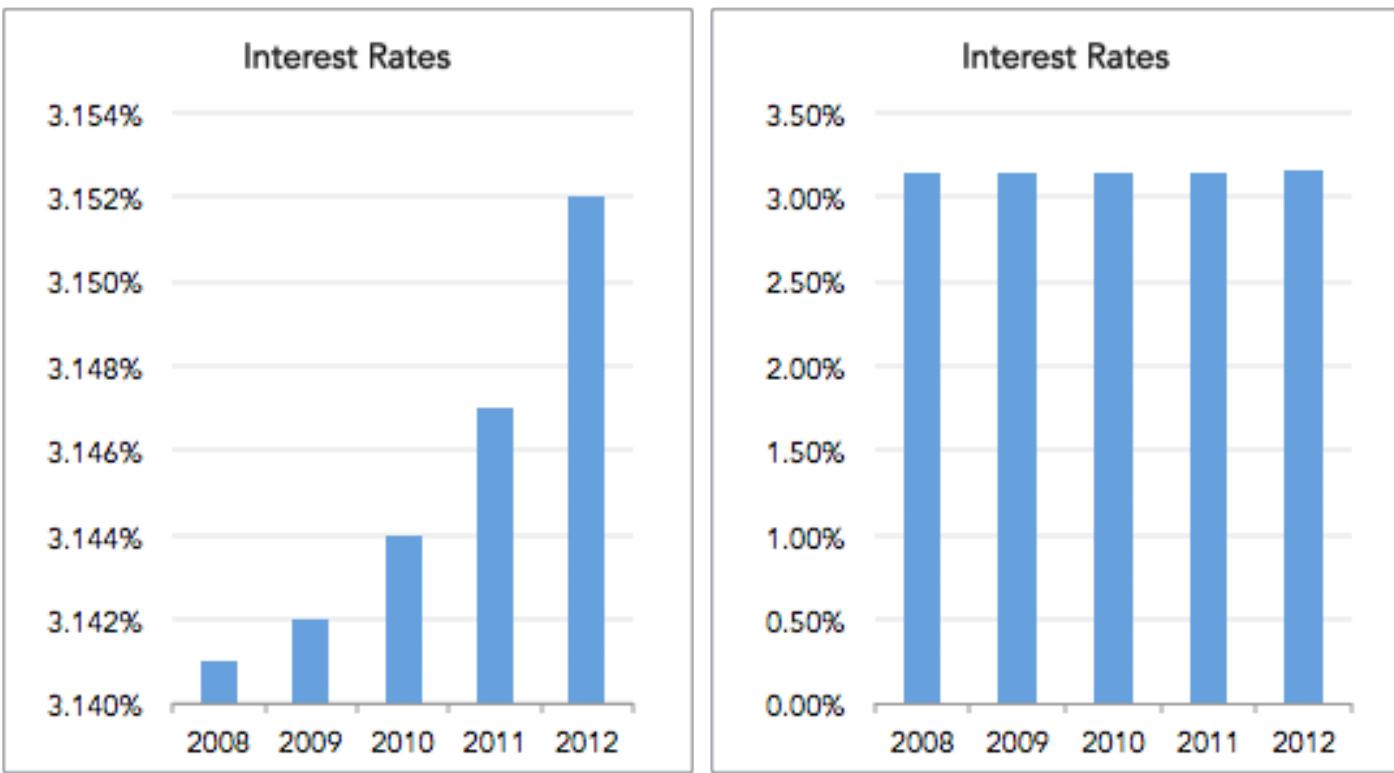
- Limited to quantitative variables
- Poorly suited for qualitative variables



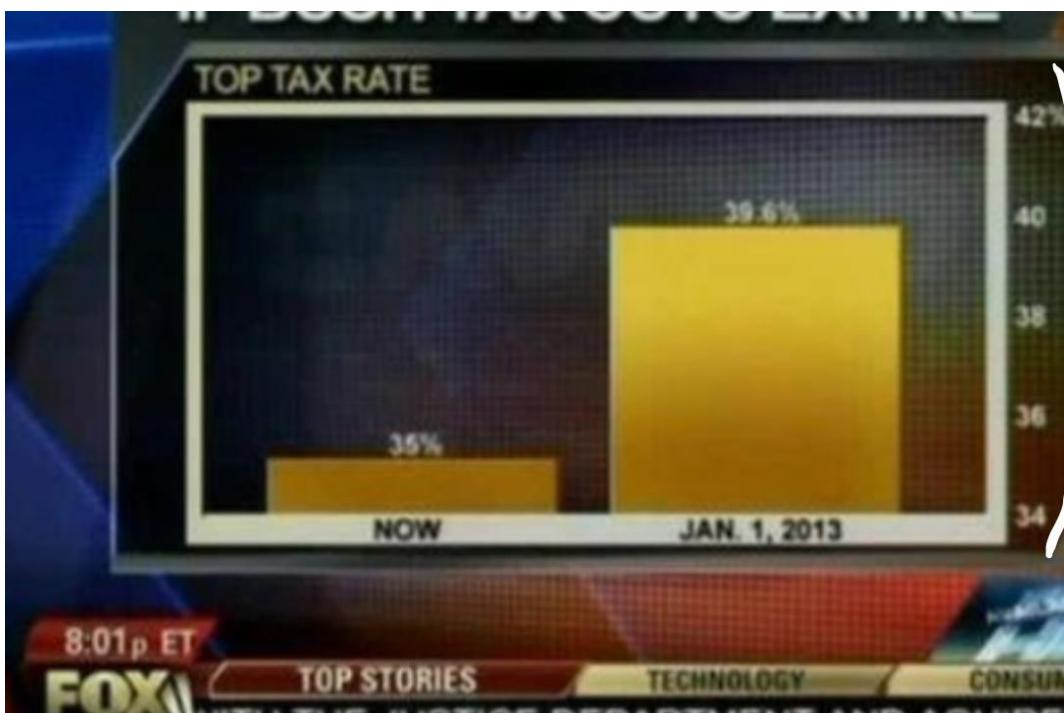
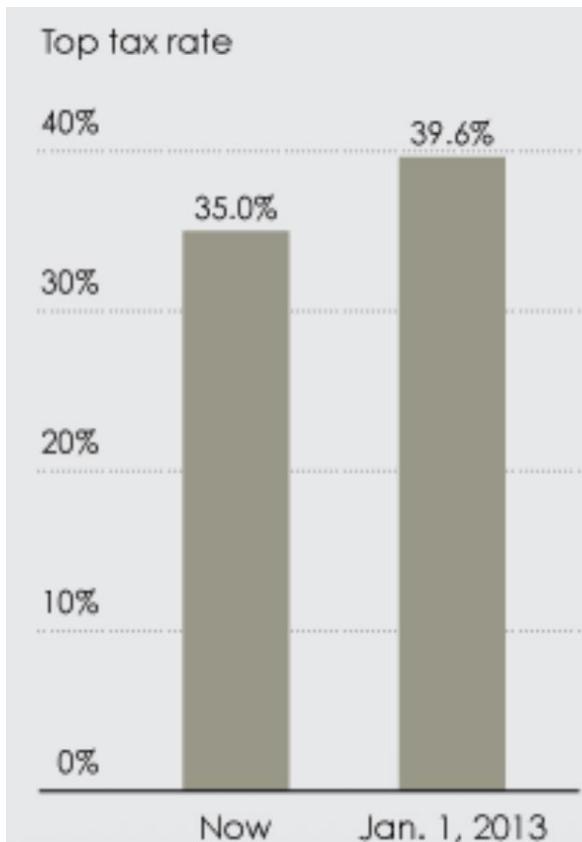
Poor or misleading visualizations



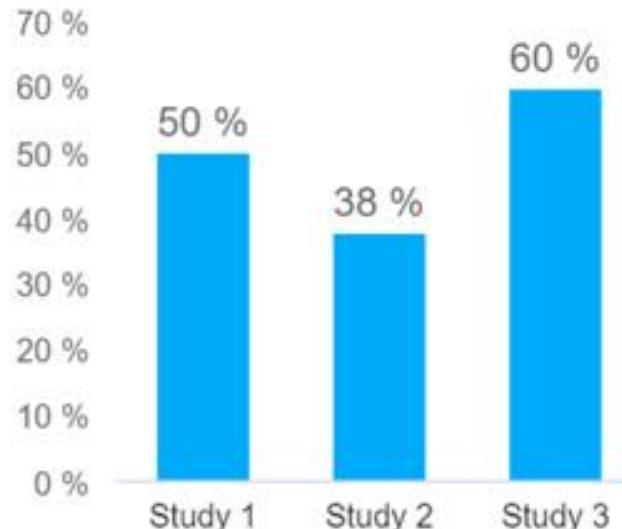
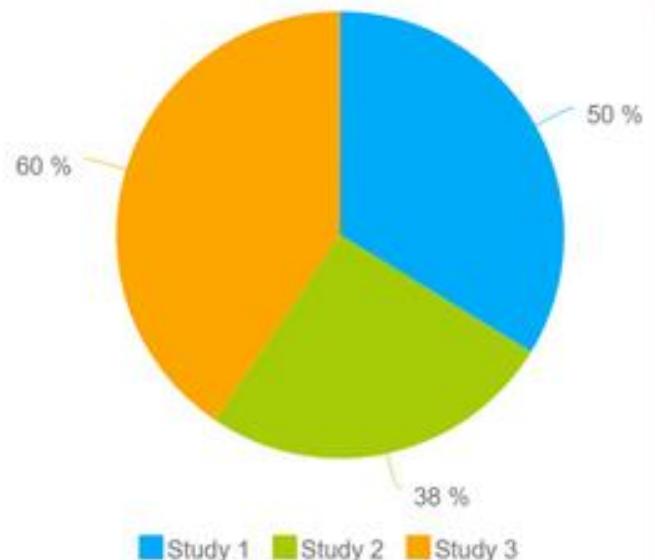
Poor or misleading visualizations



If Bush's tax rate reduction expires ...

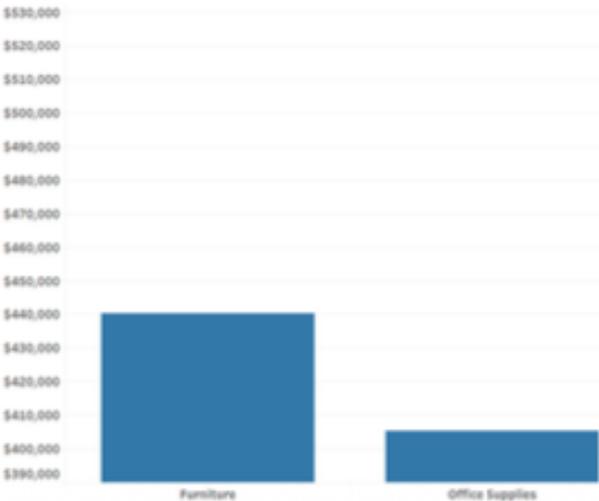


Poor or misleading visualizations

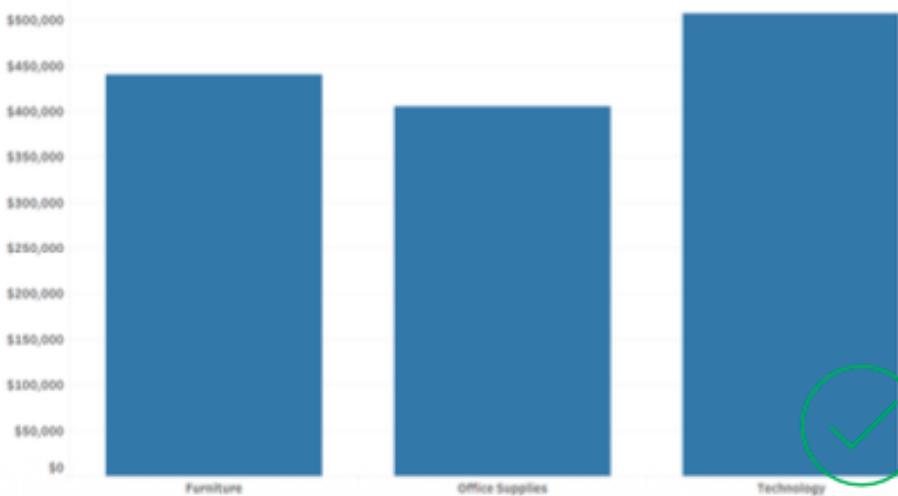


Poor or misleading visualizations

Shipping cost

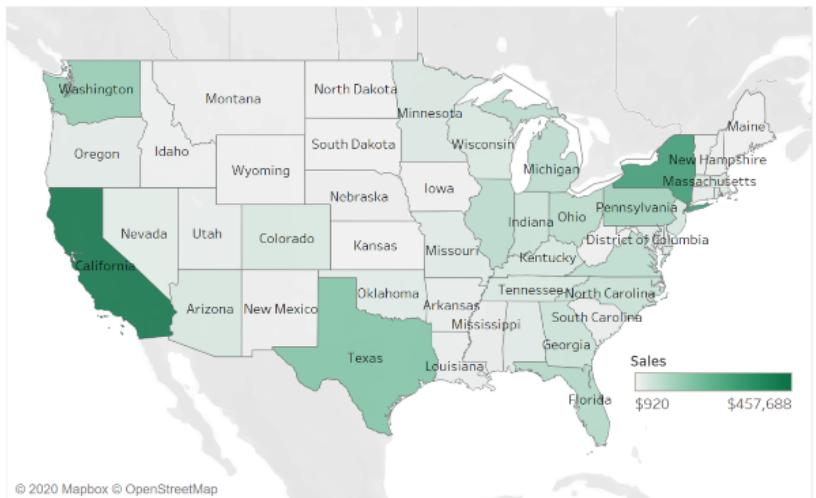


Shipping cost

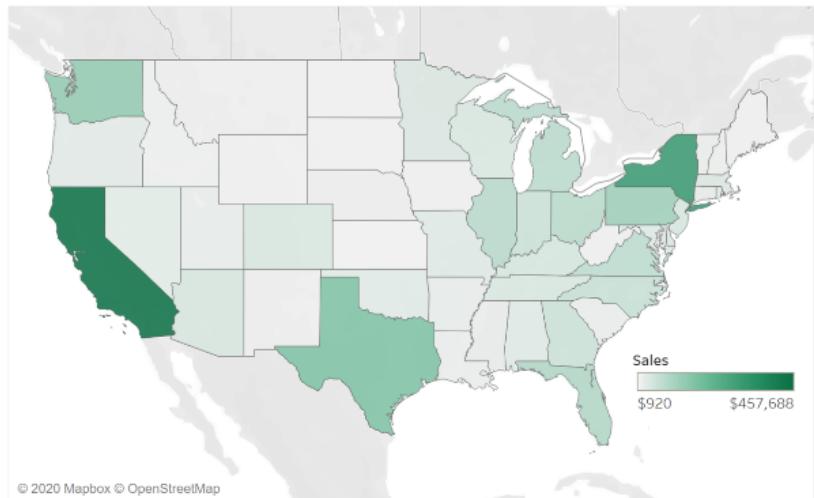


Poor or misleading visualizations

total sales map



total sales map



! **ineffective**

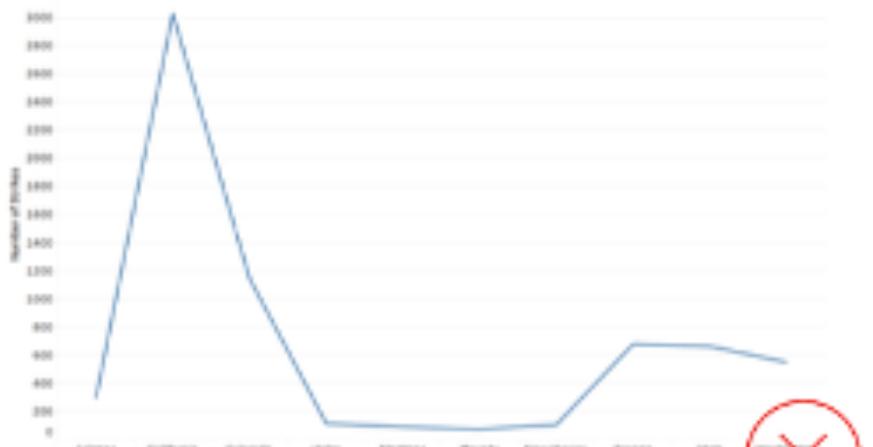
✓ **effective**



Poor or misleading visualizations

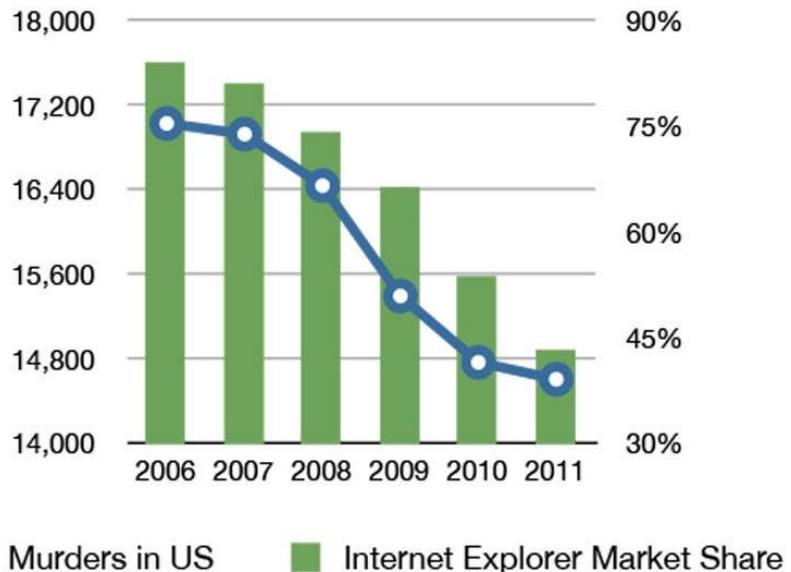


Poor or misleading visualizations



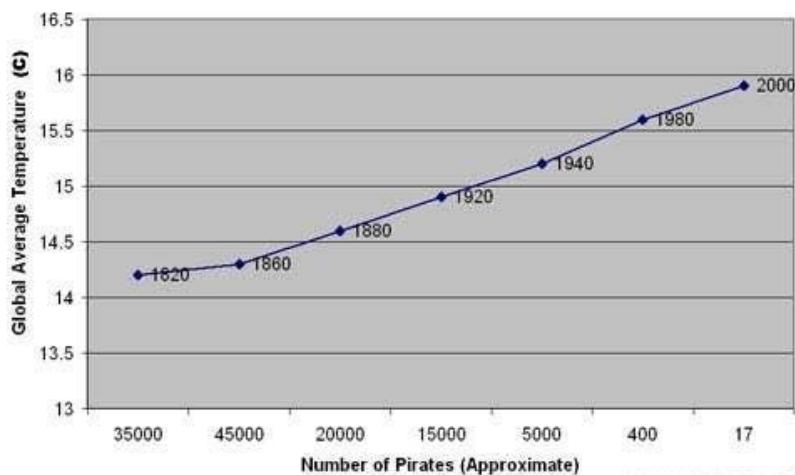
Correlation = causality?

Internet Explorer vs Murder Rate



<https://www.buzzfeednews.com/article/kjh2110/the-10-most-bizarre-correlations>

Global Average Temperature Vs. Number of Pirates



<https://www.tylervigen.com/spurious-correlations>

CHART TYPES



Line Chart

- Displays the change over time for a measure
- What kind of questions does this chart answer?
 - How has this variable changed over the past period?
 - When did this variable change?
 - How quickly did this variable change?
 - What are the trends? Can future trends be derived?
- Example: Stock market prices of the last five years



Bar Chart (horizontal | vertical)

- Comparison of data of different categories (dimensions)
- What kind of questions does this chart answer?
 - Which of these categories shows the highest/lowest value?
 - Are there any extraordinary categories?
 - What is the gap (deviation) between the lowest and highest values of different categories?
- Example: Sales per department



Bullet Diagram (horizontal | vertical)

- Modification of a bar chart that shows the performance of a primary measure of achievement of key figures.
 - What kind of questions does this chart answer?
 - As with the bar chart
 - Additional: Comparison per bar against a key figure
- Example: What is the actual turnover compared to the expected turnover?



Histogram

- Representation of the distribution of values
- What kind of questions does this chart answer?
 - Are events grouped around a certain probability?
 - Which group shows the highest values?
 - Which area covers the most observations?
- Example: Students' performance in an exam



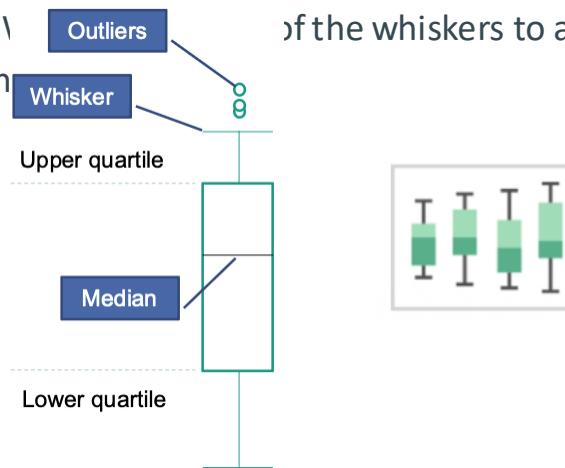
Boxplot

- Representation of the distribution within categories (dimensions)
- What kind of questions does this chart answer?
 - In which range are the values of most of the data in a category located?
 - Are there outliers in the data?
 - What is the median of values in a category?
- Example: Distribution of discounts in different product groups



Boxplot

- Inside the box is the middle 50 percent of the data
- Whiskers (antennas) describe boundaries outside of which we speak of outliers
 - No uniform definition
 - Definition according to John Tukey: Outliers are points beyond $1.5 \times (\text{IQR})$ of the whiskers to a maximum of 1.5 times the interquartile distance



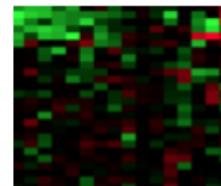
Scatterplot

- Displays relationships between two numeric variables
- What kind of questions does this chart answer?
 - Are there any patterns when looking at the data points?
 - Are there correlation relationships between the variables?
 - Are there exits in the data?
- Example: Relationship between daily calorie intake and a person's body weight.



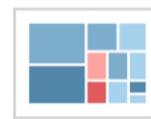
Heat Map

- Comparisons between two variables
- What kind of questions does this chart answer?
 - Is there a relationship between two variables?
 - Are certain areas particularly prominent?
 - Do two variables correlate?
- Example: Which nations won medals at the Olympics (divided into gold, silver,
● bronze)?



Tile Chart (Tree Map)

- Displays the proportion of an overall distribution.
- What kind of questions does this chart answer?
 - How much does this value contribute to the total?
 - How does the distribution of a variable change over time?
- Example: What proportion of total sales do an item's sales by sales region provide?
- Alternative chart types: pie charts, area charts, stacked bar charts



Pie Chart

- Displays the proportion of an overall distribution.
- What kind of questions does this chart answer?
 - How much does this value contribute to the total?
 - How does the distribution of a variable change over time?
- Example: What proportion of the votes did a party receive in the last elections?
- Alternative chart types: tree map, area charts, stacked bar charts



Map Charts

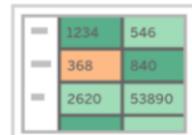
- Map charts represent positions and geographic patterns in the data.
- Variants: filled cards, point distribution cards, symbol cards, density cards...
- What types of questions can this diagram answer?
 - Which place has the largest/smallest values?
 - Which region has the largest/smallest values?
 - How do you represent deviations geographically?
- Example: How is COVID-19 spreading worldwide?



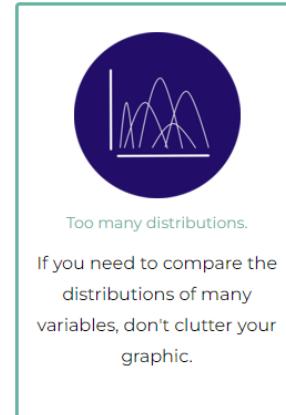
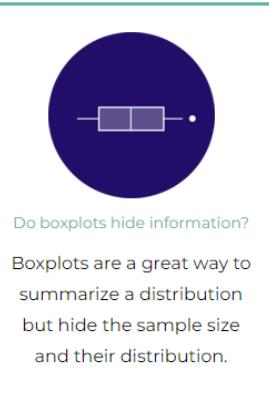
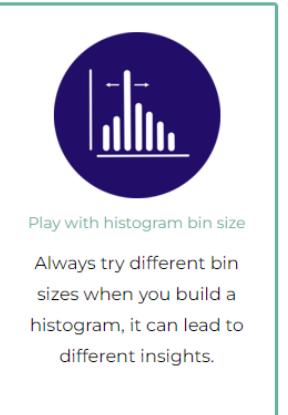
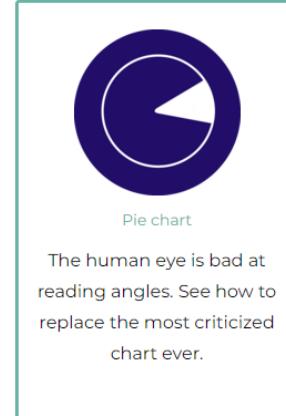
Highlight Table

- Data table with color coding
- What kind of questions does this chart answer?
 - Shows all the details of the data points
 - Colors can be assigned to specific dimensions and categories or highlight key figures (quantiles, max/min values, ...)
 - Apply color markers from diagrams and thus display their details.

	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022
	ADAC	EMEA																		
Accessories	7%	9%	11%	6%	20%	21%	25%	25%	26%	26%	25%	25%	25%	25%	25%	25%	25%	25%	25%	25%
Appliances	18%	8%	17%	12%	7%	13%	10%	14%	12%	10%	14%	17%	15%	14%	15%	16%	14%	15%	15%	14%
Art	22%	12%	14%	9%	20%	23%	25%	24%	22%	27%	29%	28%	23%	24%	24%	24%	24%	24%	24%	24%
Books	13%	16%	17%	19%	18%	19%	18%	27%	23%	10%	14%	17%	15%	20%	21%	21%	21%	21%	21%	21%
Booksellers	12%	13%	14%	13%	12%	13%	14%	9%	8%	8%	8%	8%	12%	12%	12%	12%	12%	12%	12%	12%
Chairs	1%	12%	13%	10%	6%	8%	7%	5%	12%	10%	9%	9%	9%	9%	9%	9%	9%	9%	9%	9%
Clothes	18%	18%	18%	17%	8%	11%	18%	18%	13%	13%	13%	13%	13%	13%	13%	13%	13%	13%	13%	13%
Envelopes	11%	12%	13%	7%	20%	27%	18%	27%	12%	17%	18%	19%	19%	19%	19%	19%	19%	19%	19%	19%
Fasteners	8%	9%	11%	5%	18%	22%	15%	22%	12%	14%	14%	14%	14%	14%	14%	14%	14%	14%	14%	14%
Furniture	17%	13%	17%	17%	10%	25%	17%	12%	2%	9%	1%	2%	2%	2%	2%	2%	2%	2%	2%	2%
Labels	13%	12%	10%	12%	20%	23%	27%	27%	29%	1%	23%	28%	28%	28%	28%	28%	28%	28%	28%	28%
Machine	11%	12%	9%	10%	9%	9%	7%	9%	1%	20%	20%	20%	1%	21%	21%	21%	21%	21%	21%	
Paper	11%	7%	14%	12%	17%	14%	15%	21%	21%	14%	13%	13%	13%	13%	13%	13%	13%	13%	13%	13%
Phones	21%	15%	18%	16%	20%	21%	22%	22%	7%	12%	12%	12%	12%	12%	12%	12%	12%	12%	12%	12%
Storage	8%	12%	13%	12%	6%	7%	10%	9%	14%	14%	13%	13%	13%	13%	13%	13%	13%	13%	13%	13%
Supplies	7%	8%	6%	4%	24%	13%	14%	10%	25%	18%	28%	28%	14%	16%	14%	16%	14%	16%	14%	16%
Tables	12%	2%	6%	1%	2%	12%	4%	23%	12%	4%	12%	12%	12%	12%	12%	12%	12%	12%	12%	12%



Tipps and Guidelines



Presenting results

Promises

- Enhanced Analyses and Predictions: Leveraging large volumes of data can enable deeper insights and more accurate forecasts than ever before.
- Uncertainty at the Micro Level, Accuracy at the Macro Level: Big Data allows for precise identification of patterns and trends at the macro level, despite uncertainty at the individual level.
- Quality is Crucial: Despite the vast amounts of data, ensuring data quality is essential for drawing valid conclusions.

Risks

- Incorrect Modeling and Conclusions: The complexity and volume of Big Data can lead to errors in data modeling, resulting in false or misleading outcomes.
- Privacy and Individual Rights: Processing large datasets poses risks to privacy and can conflict with individual rights.

Recommendations for Mitigating Risks

- Rigorous Data Verification: Implement strict data quality verification processes to ensure the integrity of analyses.
- Transparent Modeling: Promote transparency in your modeling processes to build trust and avoid misinterpretations.
- Preserving Privacy: Adhere to privacy regulations and practices to protect individuals' rights and foster trust in Big Data initiatives.
- Reason about results: be able to explain; why have predictions been made?
- Reason about input: data set; metrics

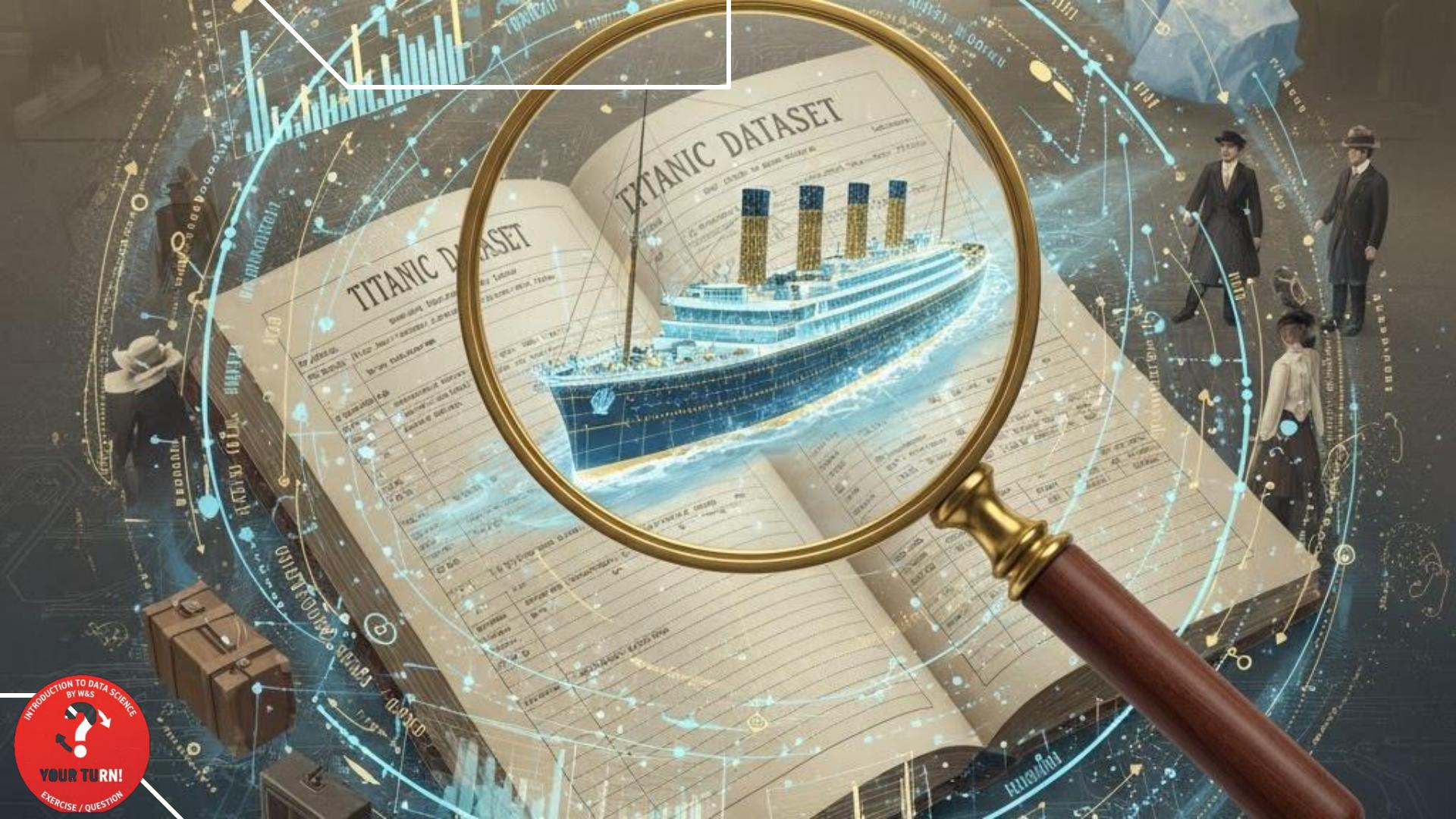
Use of Diagrams

- Usually there are data sets with many variables (characteristics) or
- we want to determine the utility value of other variables
- Application of diagrams
 - Often the same / different diagrams are used one after the other / combined
 - Decision on the type of visualization based on
 - variables,
 - dimensions in our data, and
 - the question asked



HANDS-ON TIME





Pandas Basics

datacamp

Python For Data Science

Pandas Basics Cheat Sheet

Learn Pandas Basics online at www.DataCamp.com

Pandas

The Pandas library is built on NumPy and provides easy-to-use data structures and data analysis tools for the Python programming language.

Use the following import convention:

```
>>> import pandas as pd
```

Pandas Data Structures

Series

A one-dimensional labeled array capable of holding any data type

```
Index: 0    3
      1    5
      2    7
      3    4
      4    6
```

```
>>> a = pd.Series([1, -3, 7, 4], index=[0, 1, 2, 3])
>>> a[1]
```

Dataframe

A two-dimensional labeled data structure with columns of potentially different types

Columns	Country	Capital	Population
Index	0	Brazil	1980441
	1	New Delhi	155071000
	2	China	1380000000
	3	United States	3133572905
	4	India	1283300000

```
>>> data = {'Country': 'Brazil', 'Capital': 'Brasilia', 'Population': 1980441}
>>> df = pd.DataFrame(data, index=[0, 1, 2, 3, 4])
>>> df
```

Dropping

```
>>> a.drop(1) # drop column from series (index)
>>> df.drop('Country', axis=1) # drop values from columns(axis=1)
```

Asking For Help

```
>>> help(pd.Series.sum)
```

Sort & Rank

```
>>> a.rank(method='dense').abs().sort_index(ascending=True)
>>> sort_values('Country').abs().sort_index(ascending=True)
>>> rank()
```

I/O

Read and Write to CSV

```
>>> df.to_csv('file.csv', header=True, index=False)
>>> df.to_csv('myDataFrame.csv')
```

Read and Write to Excel

```
>>> df.to_excel('file.xlsx')
>>> df.to_excel('myDataFrame.xlsx', sheet_name='Sheet1')
```

Read and Write to SQL Query or Database Table

```
>>> from sqlalchemy import create_engine
>>> engine = create_engine('sqlite:///memory:')
>>> df.to_sql('df', engine)
>>> pd.read_sql_table('df', engine)
>>> pd.read_sql_query('SELECT * FROM df', engine)
>>> read_sql('df', engine)
>>> read_sql('df', engine, index_col='Country')
>>> read_sql('df', engine, index_col='Country', schema='public')
```

Selection

Also see NumPy Array

Getting

```
>>> df['Country'] # select one column
>>> df[0] # select entire row
>>> df[['Country', 'Population']] # select multiple columns
>>> df.loc[0] # select entire row by index
>>> df.loc[0, 'Country'] # select single value by row & column labels
```

Selecting, Boolean Indexing & Setting

```
>>> df.loc[0]
>>> df.loc[[0]] # select single row by subset of rows
>>> df.loc[0, 'Country'] # select single column of subset of columns
>>> df.loc[0, :].loc['Country'] # select single value of subset of rows
```

Boolean Indexing

```
>>> df[(df['Country'] == 'Brazil') | (df['Country'] == 'New Zealand')]
>>> df[df['Country'].isin(['Brazil', 'New Zealand'])]
>>> df[(df['Country'] == 'Brazil') | (df['Country'] == 'New Zealand')]
```

Retrieving Series/DataFrame Information

Basic Information

```
>>> df.shape # returns (rows,columns)
>>> df.info() # returns descriptive stats
>>> df.dtypes # returns column data types
>>> df.describe() # returns summary statistics
>>> df.size # returns number of elements
>>> df.shape[0] # returns number of rows
>>> df.shape[1] # returns number of columns
```

Summary

```
>>> df.sum() # sum of values
>>> df.count() # dimension sum of values
>>> df.mean() # mean of values
>>> df.idxmax() # index of maximum value
>>> df.idxmin() # index of minimum value
>>> df.mode() # return mode of values
>>> df.median() # median of values
```

Applying Functions

```
>>> f = lambda x: x**2
>>> df.apply(f) # apply function to DataFrame
>>> df.applymap(f) # apply function element-wise
```

Data Alignment

Internal Data Alignment

N/A values are introduced in the indices that don't overlap:

```
>>> pd.concat([df, df, df], ignore_index=True)
>>> df
0   Brazil
1   New
2   Deli
3   India
4   United
```

Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> df.fillna(0, inplace=True)
0   1.0
1   4.0
2   5.0
3   7.0
4   8.0
>>> df
0   1.0
1   4.0
2   5.0
3   7.0
4   8.0
```

Learn Data Skills Online at
www.DataCamp.com



https://media.datacamp.com/legacy/image/upload/v1676302204/Marketing/Blog/Pandas_Cheat_Sheet.pdf

Plotly Basics

Data Visualization with Plotly Express in Python

Learn Plotly online at www.DataCamp.com

What is plotly?

Plotly Express is a high-level data visualization package that allows you to create interactive plots with very little code. It is built on top of Plotly Graph Objects, which provides a lower-level interface for developing custom visualizations.

Interactive controls in Plotly

Plotly plots have interactive controls shown in the top-right of the plot. The controls allow you to do the following:

- Download plot as a png: Save your interactive plot as a static PNG.
- Zoom: Zoom in on a region of interest in the plot.
- Pan: Pan: Move around in the plot.
- Box Select: Select a rectangular region of the plot to be highlighted.
- Lasso Select: Draw a region of the plot to be highlighted.
- Autoscale: Zoom to a "best" scale.
- Reset axes: Return the plot to its original state.
- Toggle Spikes Lines: Show or hide lines to the axes whenever you hover over data.
- Show closest data on hover: Show details for the nearest data point to the cursor.
- Compare data on hover: Show the nearest data point to the x-coordinate of the cursor.

Plotly Express code pattern

The code pattern for creating plots is to call the plotting function, passing a data frame as the first argument. The x argument is a string naming the column to be used on the x-axis. The y argument can either be a string or a list of strings naming column(s) to be used on the y-axis.

```
px.plotting(px.DataFrame, # A DataFrame holding visual data
           x="=x", # Accepts a string or a list of strings
           y="=y", # Accepts a string or a list of strings
           title="=title", # Accepts a string
           xaxis_title="=xaxis_title", # Accepts a string
           yaxis_title="=yaxis_title", # Accepts a string
           width==width, # Accepts an integer
           height==height) # Accepts an integer
```

Common plot types

Import plotly

```
# Import plotly express as px
import plotly.express as px
```

Scatter plots

Create a scatterplot on a DataFrame named classical_data px.scatter(classical_data, x="Experiment_1", y="Experiment_2")



Line plots

Create a lineplot on a DataFrame named stock_data px.line(stock_data, x="Date", y="PBOE", "MSDN")



Bar plots

Create a barplot on a DataFrame named commodity_data px.bar(commodity_data, x="Country", y="Gold", "Silver", "Bronze", color_discrete_map={"Gold": "#yellow", "Silver": "#white", "Bronze": "#brown"})



Histograms

Create a histogram on a DataFrame named billt_data px.histogram(billt_data, x="Total_Bill")



Heatmaps

Create a heatmap on a DataFrame named tips_data px.imshow(tips_data, x="Time", y="Day", color_continuous_scale="magma")



Customizing plots in plotly

The code pattern for customizing a plot is to save the figure object returned from the plotting function, call its update_traces() method, then call its show() method to display it.

```
# Create a plot with plotly (can be of any type)
fig = px.line(px.DataFrame())
# A context manager to make it work with update_traces() and .show()
with fig.update_traces():
    fig.update_traces()
fig.show()
```

Customizing markers in Plotly

When working with visualizations like scatter plots, treelists, and more, you can customize markers according to certain properties. These include:

- size: set the marker size
- color: set the marker color
- opacity: set the marker transparency
- line: set the width and color of a border
- symbol: set the marker shape

In this example, we're updating a scatter plot named fig_scatter

```
fig_scatter.update_traces(marker={"size": 24, "color": "#green", "opacity": 0.5, "line": {"width": 2, "color": "cyan"}, "symbol": "square"})

fig_scatter.show()
```



Customizing lines in Plotly

When working with visualizations that contain lines, you can customize them according to certain properties. These include:

- color: set the line color
- dash: set the dash style ("solid", "dash", "longdash", "dotdash", "longdashdot")
- width: set the line width
- shape: set how values are connected ("linear", "spline", "vh", "vhv")
- smooth: set the line smoothness

In this example, we're updating a scatter plot named fig_line

```
fig_line.update_traces(line={"color": "#red", "dash": "solid", "width": 3})

fig_line.show()
```



Customizing bars in Plotly

When working with barplots and histograms, you can update the bars themselves according to the following properties:

- color: set the marker color
- opacity: set the marker transparency
- line: set the width and color of a border
- symbol: set the shape of the marker

In this example, we're updating a scatter plot named fig_bar

```
fig_bar.update_traces(marker={"color": "#green", "opacity": 0.5, "line": {"width": 2, "color": "cyan"}})

fig_bar.show()
```



In this example, we're updating a histogram named fig_hist

```
fig_hist.update_traces(marker={"color": "#blue", "opacity": 0.5, "line": {"width": 2, "color": "cyan"}})

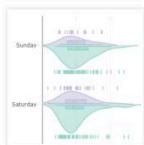
fig_hist.show()
```



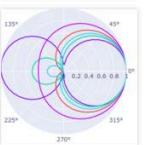
Learn Data Skills Online at
www.DataCamp.com


Plotly Documentation

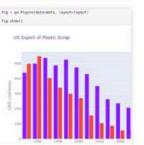
Fundamentals



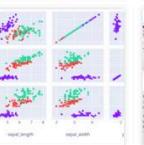
The Figure Data Structure



Creating and Updating Figures



Displaying Figures



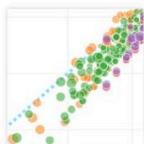
Plotly Express



Analytical Apps with Dash

[More Fundamentals »](#)

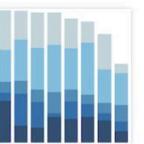
Basic Charts



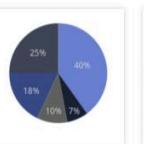
Scatter Plots



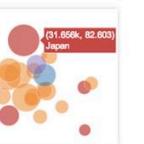
Line Charts



Bar Charts



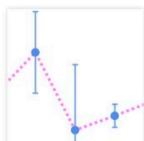
Pie Charts



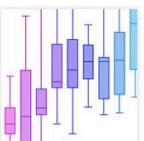
Bubble Charts

[More Basic Charts »](#)

Statistical Charts



Error Bars



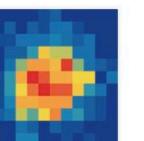
Box Plots



Histograms



Distplots



2D Histograms

[More Statistical Charts »](#)

<https://plotly.com/python/>

Challenges



Whats hard about data science

- Getting the data (usually)
- Overcoming assumptions
- Communication
 - With domain experts
 - Expectation management for client
- Making ad-hoc explanations of data patterns
- Overgeneralizing
- Not checking enough (validate models, data pipeline integrity, etc.)
- Using statistical tests correctly
- Prototype to Production transition
- Data pipeline complexity (team boundary)

Data Challenges (1)

	Variety	Volume	Velocity
Challenge	Handling multiplicity of types, sources and formats	Dealing with large volumes of data	Streams, sensors, near real-time data
Impacted Tasks	Data integration	Storage, processing & analytics	Processing & analytics
Solution	Semantic technologies are a good fit	Distributed storage & parallel processing	Real-time technologies

Data Challenges (2)

- **Data veracity:** coping with uncertainty, imprecision, missing values, misstatements or untruths.
- **Data quality:** determining the quality of datasets and relevance to particular issues. Depends on the use case:
 - How broad/complete is the data?
 - How fine is the sample resolution? How timely are the readings?
 - Does the data contain any “noise” (errors)? Is it representative?

Data Challenges (3)

- **Data discovery:** finding relevant data from enormous amount of data available on the Web.
- **Data dogmatism:** analysis of Big Data can offer remarkable insights. However, data analysis should not entirely replace domain expert knowledge, but act as a tool to support /confirm facts.
 - E.g., Google Flu Trends

Data Challenges (3)

Example: Data dogmatism

google.org Flu Trends

[Google.org home](#)

[Dengue Trends](#)

[Flu Trends](#)

[Home](#)

[Germany](#) 

[National](#) 

[Download data](#)

[How does this work?](#)

[FAQ](#)

Explore flu trends - Germany

We've found that certain search terms are good indicators of flu activity. Google Flu Trends uses aggregated Google search data to estimate flu activity. [Learn more »](#)

National



Source: <https://www.google.org/flutrends/de/#DE>

“By counting how often we see these [flu-related topics] **search queries**, we can estimate **how much flu is circulating** in different countries and regions around the world.”

Data Challenges (3)

Example: Data dogmatism

Comparison of Google Flu Trends model against medical reports.



The Parable of Google Flu: Traps in Big Data Analysis

David Lazer^{1,2,*}, Ryan Kennedy^{1,3,4}, Gary King³, Alessandro Vespignani^{5,6,3}

 Author Affiliations

 Corresponding author. E-mail: d.lazer@neu.edu.

In February 2013, Google Flu Trends (GFT) made headlines but not for a reason that Google executives or the creators of the flu tracking system would have hoped. *Nature* reported that GFT was predicting more than double the proportion of doctor visits for influenza-like illness (ILI) than the Centers for Disease Control and Prevention (CDC), which bases its estimates on surveillance reports from laboratories across the United States (1, 2). This happened despite the fact that GFT was built to predict CDC reports. Given that GFT is often held up as an exemplary use of big data (3, 4), what lessons can we draw from this error?

 [Read the Full Text](#)

Source: <http://www.sciencemag.org/content/343/6176/1203>

SCIENCE BIG DATA

Google's Flu Project Shows the Failings of Big Data

Bryan Walsh @bryanrwalsh | March 13, 2014



A new study shows that using big data to predict the future isn't as easy as it looks—and that raises questions about how Internet companies gather and use information



Source: <http://time.com/23782/google-flu-trends-big-data-problems/>

Process Challenges (1)

Big-Data Scientists “Janitor Work”

Data acquisition:

- Data availability
- Data permissions

Aligning/integrating data from different sources:

Syntactic challenge: Data in different formats

Semantic challenges: Resolving when two objects are the same, describing relationships between data points, resolving inconsistencies

Transforming, cleaning and organizing the data into a form suitable for analysis

50% - of data scientists' time

80% - spent in “Data wrangling”

Process Challenges (2)

Modeling data:

Mathematically: mathematical models to describe the data Statistical methods or Machine Learning

Simulations

Knowledge representation: ontologies and rules

Understanding the output:

Interpreting the results

Visualizing

Sharing the results

Data privacy, security, and Governance:

- Ensuring that data is used correctly (abiding by its intended uses and relevant laws).
- Tracking how the data is used, transformed, derived, ...
- Managing data lifecycle.

“Many data warehouses contain sensitive data such as personal data. There are legal and ethical concerns with accessing such data. So the data must be secured and access controlled as well as logged for audits”.

Michael Blaha, Modelsoft Consulting Corporation, 2012

Source: <http://www.odbms.org/blog/2012/03/data-modeling-for-analytical-data-warehouses-interview-with-michael-blah/>

Data science – two sides of the same coin

Opportunities	Risks
<ul style="list-style-type: none">• Discover potential• Develop new services• Informed decisions based on forecasts• Service and product improvement	<ul style="list-style-type: none">• Surveillance• Manipulation• Uselessness• Data protection

Bis Donnerstag

samuel.schlenker@lehre.dhbw-stuttgart.de

