

Exploratory Data Analysis and Visualization

1. Introduction

LI Xinke

Department of Data Science
City University of Hong Kong

<http://xinke.li>

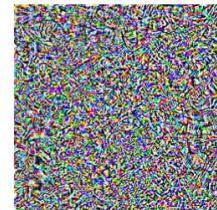
LI Xinke

- Global Research Assistant Professor
 - Department of Data Science, City University of Hong Kong
- Ph.D., Industry System Engineering, National University of Singapore
- Research: 3D Computer Vision, Trustworthy Geometry Learning, Medical AI applications.



“pig”

+ 0.005 x



=
“airliner”



香港城市大學
City University of Hong Kong



This course

Major objectives:

- Introduce exploratory data analysis (EDA).
- Introduce popular visualization techniques.

Tools

- The main coding environment is **Python** (consistent with other courses – popular).
- **Tableau** will also be introduced.
- There are several lab sessions right after the lecture (for Python and Tableau only).

This course

Tentative Schedule

- Mid-term exam in Week 10. No final.
- There is a group project
 - Find your friends before Week 9. (6-8 members)
 - Project presentation in Week 11-13.

This course

- Grading:
 - Group project: 40%
 - Individual coursework: 25%
 - Quizzes: on time submission - 2pt,
late submission - 1pt
 - HWs: based on performance,
full grade - 10pt
 - Midterm test: 35%

This course

- Canvas: 202509SDSC5002C62
(not 202509SDSC5002 !!!)

Same content as 202509SDSC5002C61 Prof. Lijia Wang

- You can find:
 - Announcement
 - Slides
 - Tutorials
 - Recordings
 - Quiz
 - Assignment
 - Project details
- ...

If you need more information to be posted on Canvas,
contact me!

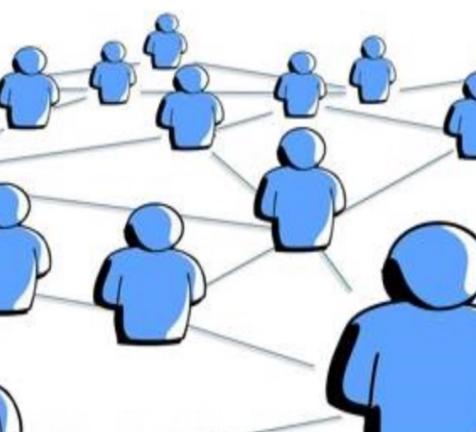
This course

Tentative Schedule:

- Weeks 1-3: key concepts of exploratory data analysis and data visualization
- Week 4: data analytics and data visualization for machine learning
- Week 5: National day
- Weeks 6-7: data analytics and data visualization for machine learning
- Weeks 8: High-dimensional data visualization
- Weeks 9: Chung Yeung Festival (Project Proposal Submission)
- Week 10: mid-term
- Weeks 11: Advanced topics in EDA, start of group project presentation (tentative)
- Week 12: group project presentation
- Week 13: group project presentation
- Instructor: LI, Xinkel
 - xinkeli@cityu.edu.hk
- TAs:
 - Mr. PAN Jiming(Organize Group Project)
 - jmpan3-c@my.cityu.edu.hk
 - Mr. ZHANG Wenlin (Python)
 - wl.z@cityu.edu.hk
 - Ms. LI Yanru(Tableau)
 - yanru.li@my.cityu.edu.hk

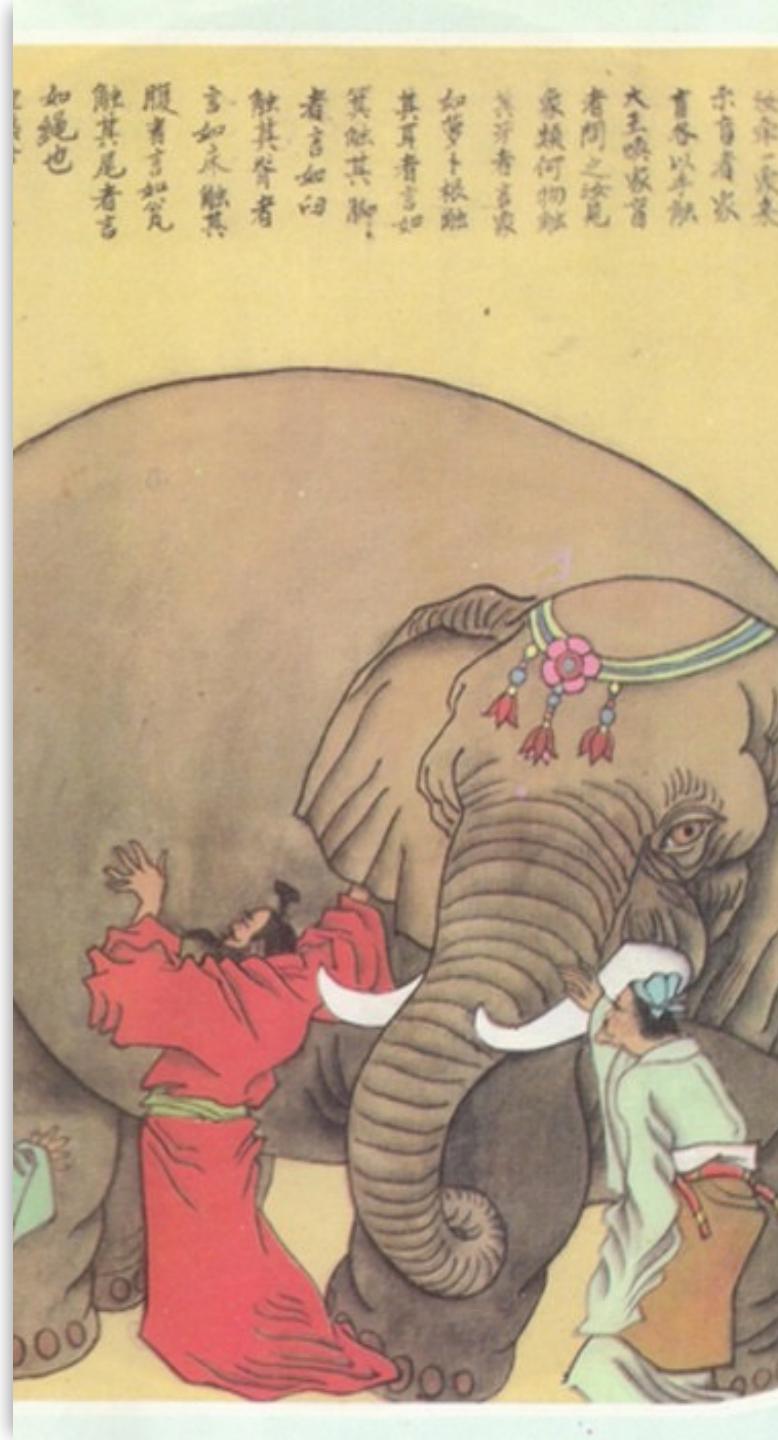


Data Analytics

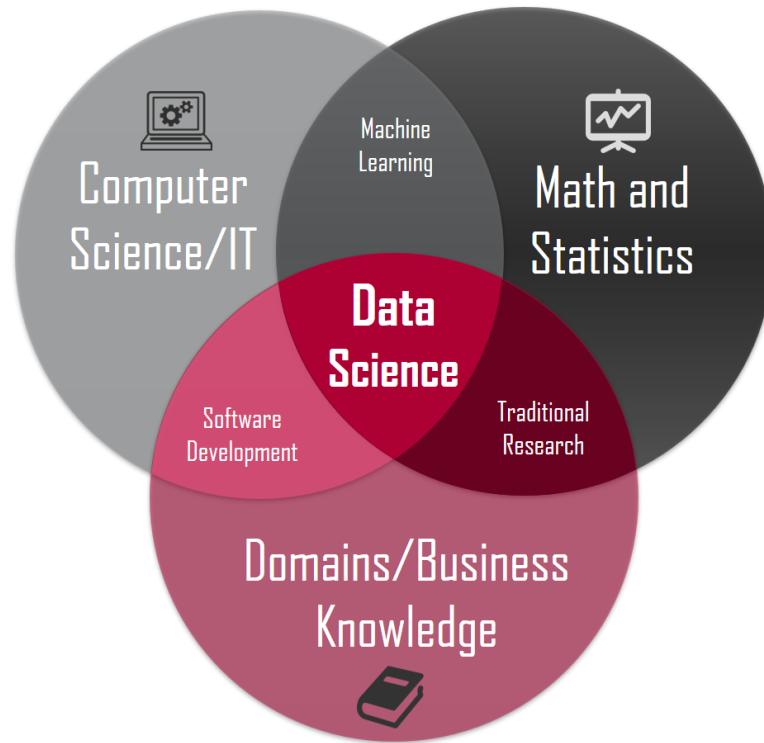
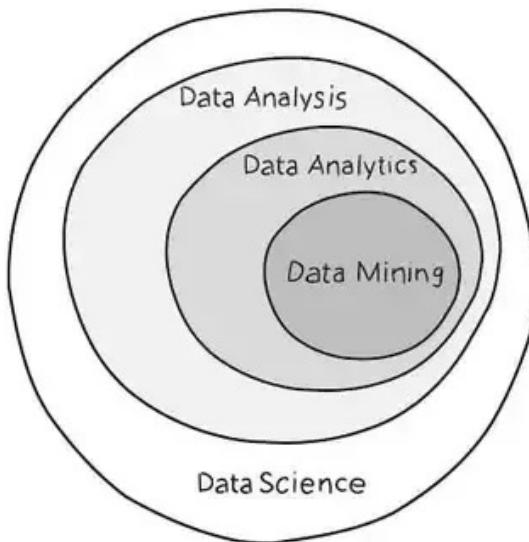
- 
- Systems (social, economic, ecological, biological, political, educational, transportation and other types of systems) generate huge amount of data.
 - In today's world, we have the necessary hardware and software to utilize this data effectively.
 - The time has come to **capitalize** on the potential of data analytics.
- 

Data Analytics

- The world is a mystery for everyone.
 - e.g., animals/humans remain insufficiently studied in the realms of biology, psychology, and social sciences.
- Animals/Humans' behaviors are **stochastic** and **SUPER complex**. It is difficult to build a **PHYSICAL** model!
- Fortunately, we can extract data, analyze the data, and generate insights and knowledge about the system.
- “盲人摸象” (Blind men and an elephant)



Data Analytics



- <https://www.mo-data.com/what-is-the-difference-between-data-analytics-data-analysis-data-mining-data-science-machine-learning-big-data-and-predictive-analytics/>

Question

- What is your current primary smart phone?
 - Android
 - iPhone
 - Others
- What was your last primary smart phone?
 - Android
 - iPhone
 - Others

Very simple model

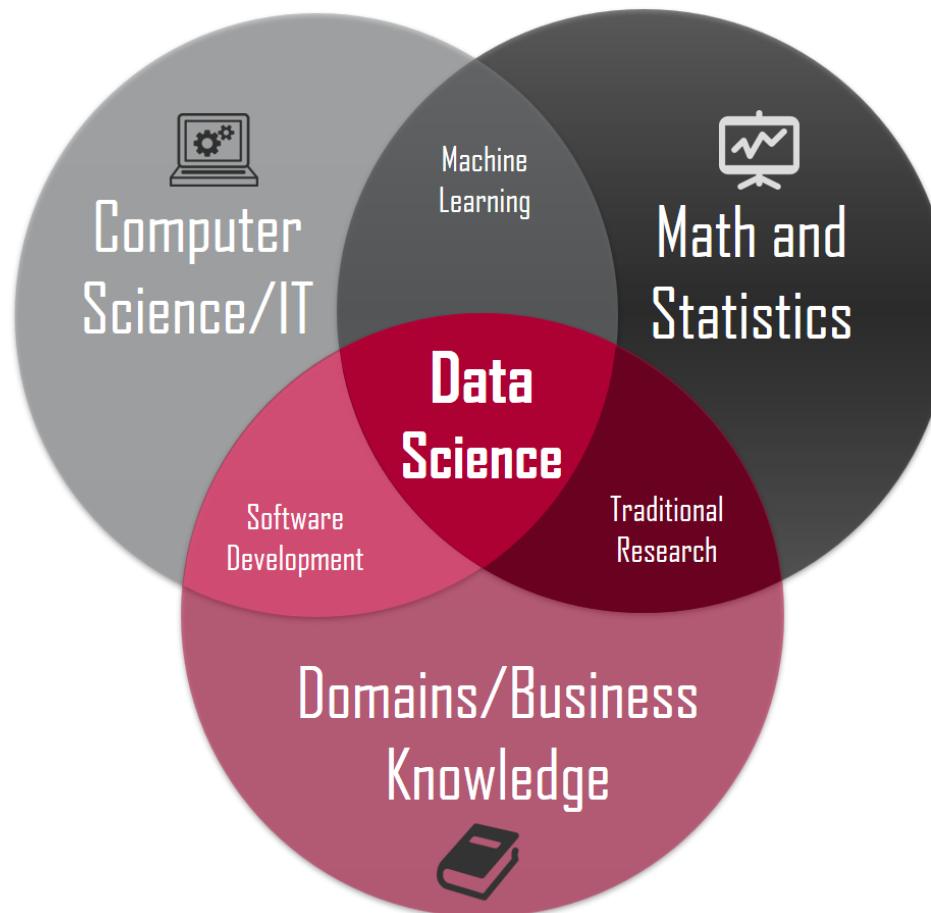
- # of people using iPhone, x_1 ; # of people using Android, x_2
- # of people used iPhone, y_1 ; # of people used Android, y_2
- # of people who switched to iPhone from Android, $x_{1,2}$
 - %: $p_{1,2} = x_{1,2} / y_2$
- # of people who switched from Android to iPhone, $x_{2,1}$
 - %: $p_{2,1} = x_{2,1} / y_1$
- % of people who remained (stickiness of the product)
 - iPhone: $p_{1,1} = (x_1 - x_{1,2}) / y_1$
 - Android: $p_{2,2} = (x_2 - x_{2,1}) / y_2$

Prediction

- Next time when you upgrade your phone, how many of you will be using
 - iPhone: $x'_1 = x_1 p_{1,1} + x_2 p_{2,1}$
 - Android: $x'_2 = x_2 p_{2,2} + x_1 p_{1,2}$
 - $x'_1 + x'_2 = x_1 + x_2$

Domain Knowledge

- Domain knowledge is important!



If we have domain knowledge

- We would ask for more questions:
 - Gender, g
 - Income, m
 - Interest in digital products, d
- Then, we build a regression model:
 - $y_{t+1} = \beta_0 + \beta_1 y_t + \beta_2 g + \beta_3 m + \beta_4 d + \epsilon$

Additional domain knowledge

- Social influence
 - % of friends who are using Android, a
- $y_{t+1} = \beta_0 + \beta_1 y_t + \beta_2 g + \beta_3 m + \beta_4 d + \beta_5 a + \epsilon$
- The more the better “多多益善”?
 - How to evaluate the performance of the prediction?
 - EDA and data visualization
- If too much – we need to “engineer” features
 - Data mining

AI and Statistics

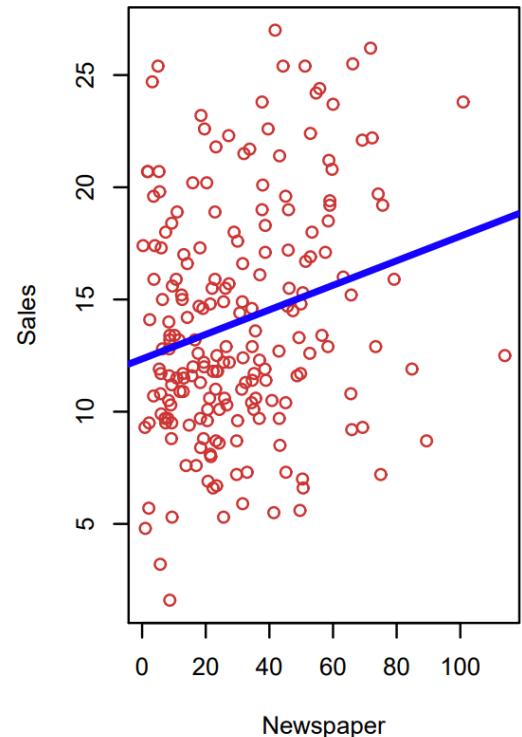
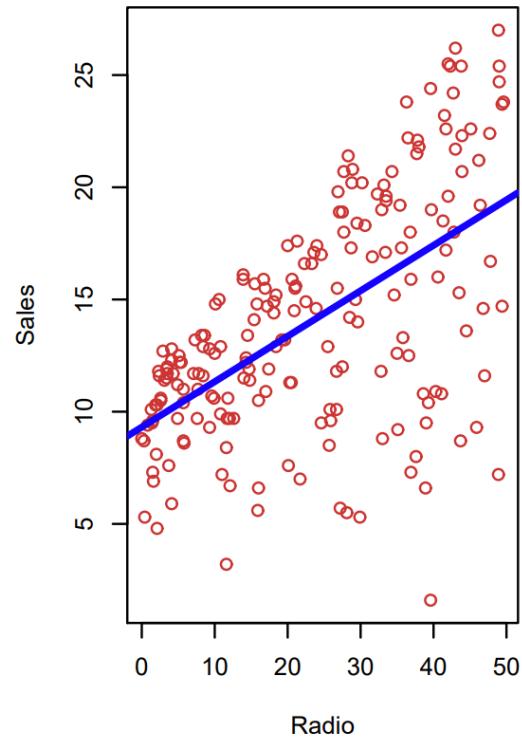
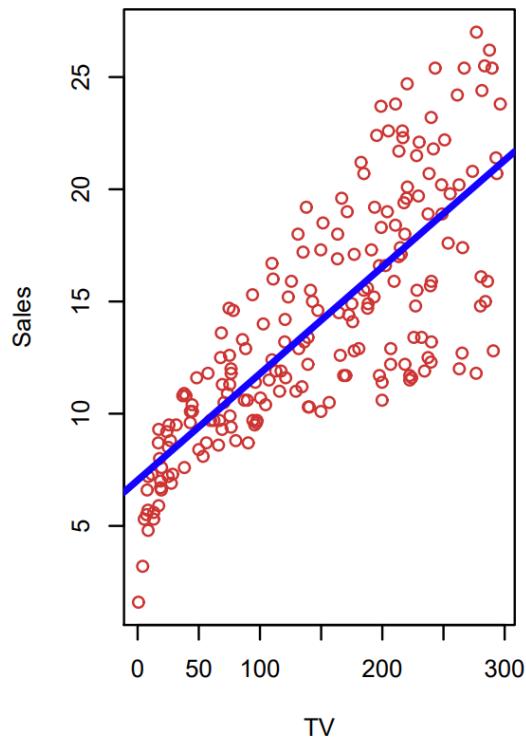
- Statistics is the key tool for machine learning. (Key philosophy)
- Machine learning is packaged as “AI” for now.

Let's build a skyscraper step by step!



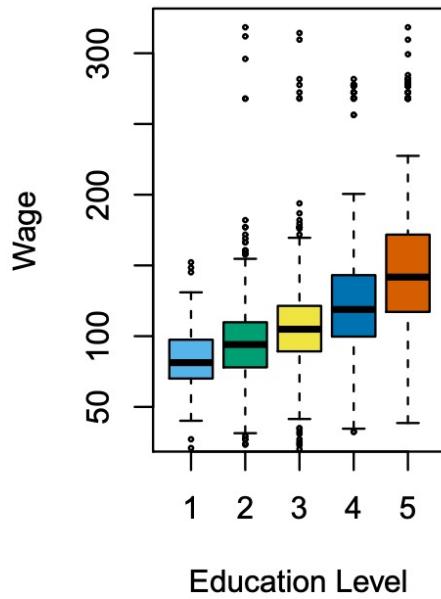
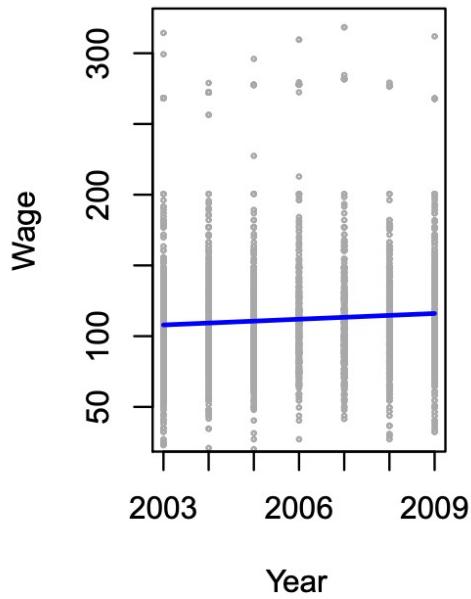
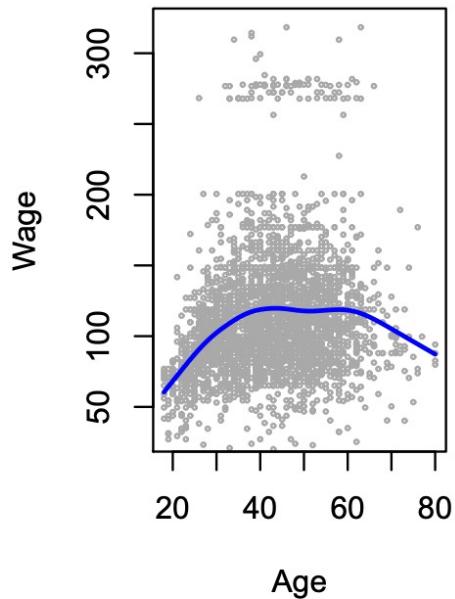
Statistical Learning

- Example (Advertising data)



Statistical Learning

- Example

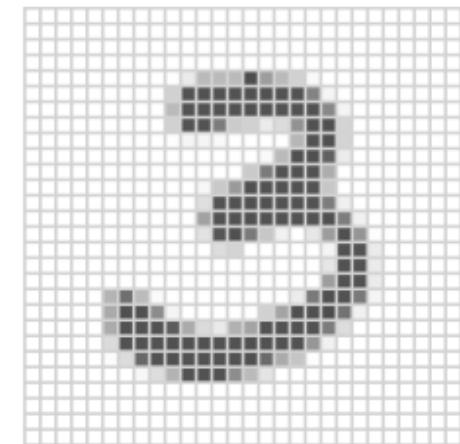


Income survey data for males from the central Atlantic region of the USA in 2009.

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

Statistical Learning

- Example



[https://en.wikipedia.org/
wiki/MNIST_database](https://en.wikipedia.org/wiki/MNIST_database)

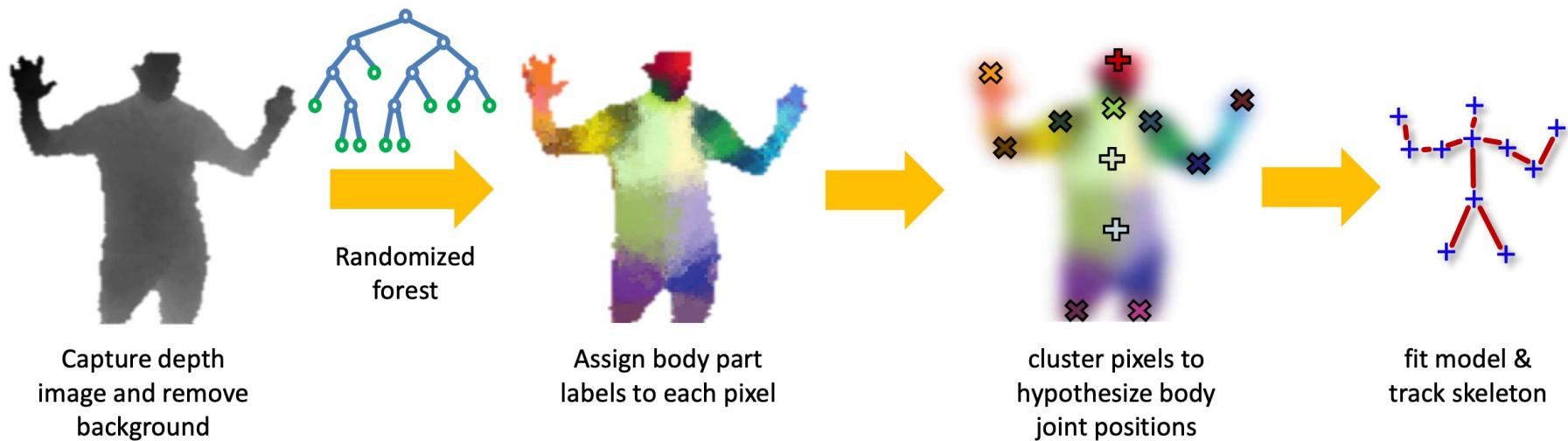
Statistical Learning

- Example



Statistical Learning

- Kinect
- https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/ks_book_2012.pdf



Supervised Learning

- Outcome measurement Y
 - dependent variable, response, target
- Vector of p predictor measurements X
 - inputs, regressors, covariates, features, independent variables
- In **regression problem**, Y is quantitative
 - price, blood pressure
- In **classification problem**, Y takes discrete values
 - survived/died, digit 0-9, cancer class of tissue sample
- We have training data $(x_1, y_1), \dots, (x_N, y_N)$. These are observations (examples, instances) of these measurements.

Supervised Learning

Objective

- On the basis of the training data we would like to:
 - Accurately **predict** unseen test cases.
 - **Understand** which inputs affect the outcome, and how.
 - **Assess** the quality of our predictions and inferences.

Supervised Learning

A few points

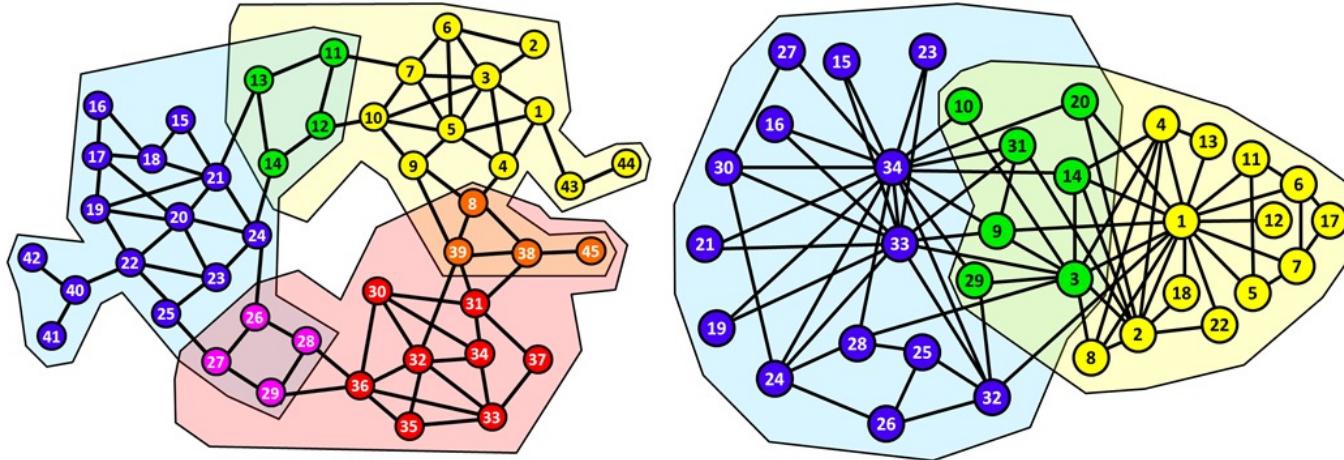
- It is important to **understand the ideas behind the various questions**, in order to know how and when to use them.
 - Domain knowledge – IE/Biz students are good at industrial/biz problems
- One has to understand **the simpler methods first**, in order to master the more complicated ones.
- It is important to accurately **assess the performance** of a method, to know how well or how badly it is working
 - Simpler methods often perform as well as fancier ones!
- This is an exciting research area, having important applications in science, industry and finance.

Unsupervised Learning

- **No outcome variable**, just a set of predictors (features) measured on a set of samples.
- Objective is more general –
 - find groups of samples that behave similarly,
 - find features that behave similarly,
 - find linear combinations of features with the most variation.
- Difficult to know how well you are doing.
- Different from supervised learning but can be useful as a pre-processing step for supervised learning.

Unsupervised Learning

- Community detection in social networks



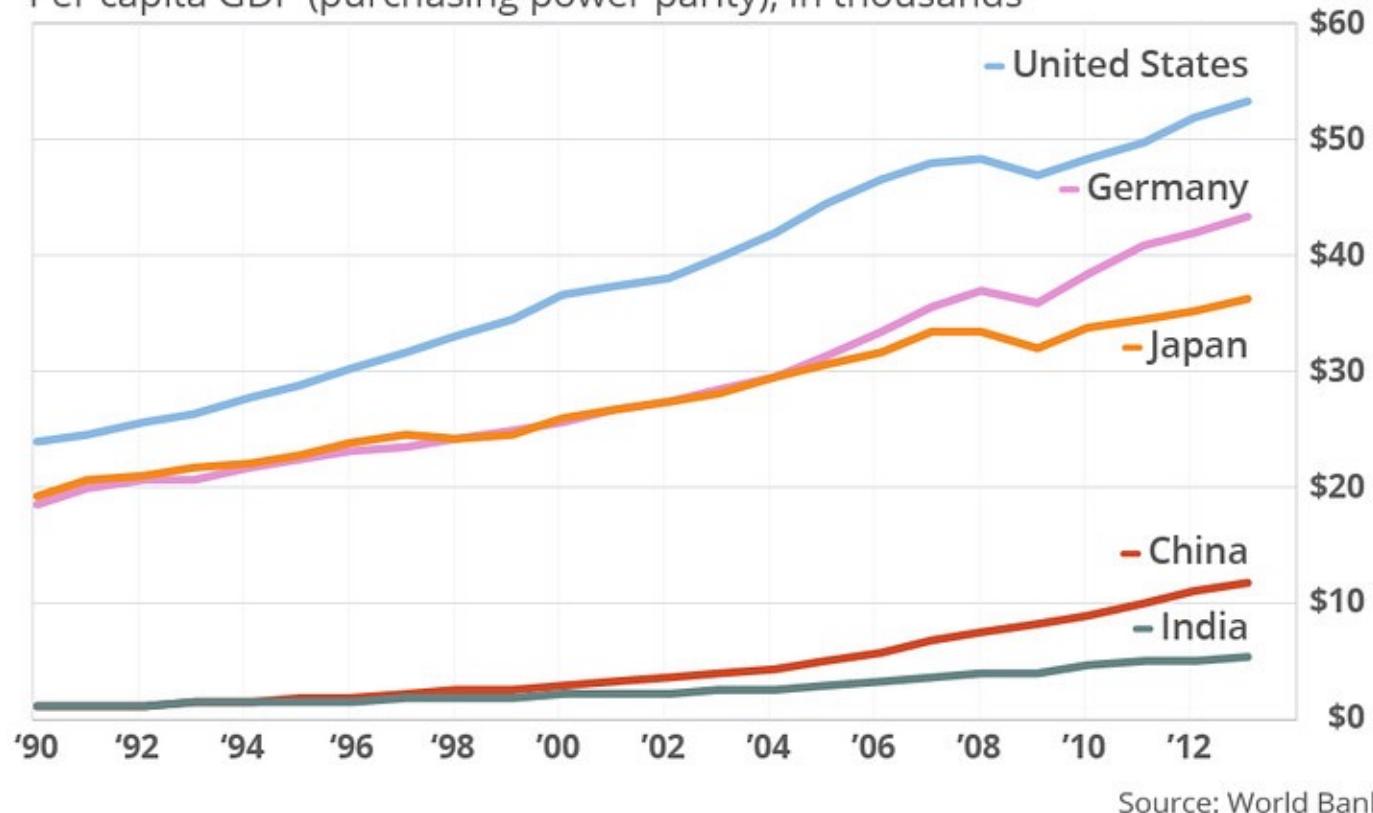
Focus of this course

- Exploratory data analysis and visualization
 - Exploratory data analysis
 - Data visualization

Example: GDP

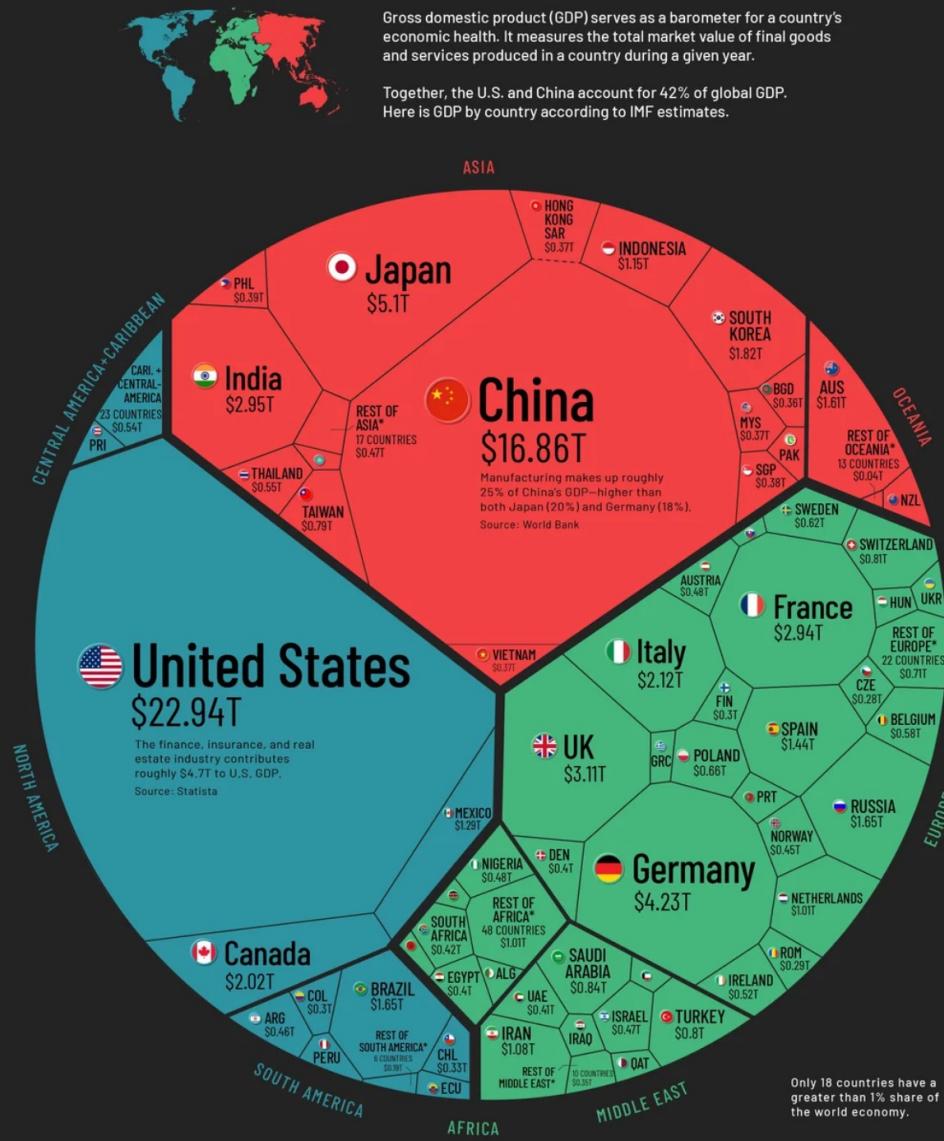
Americans are far richer than Chinese

Per capita GDP (purchasing power parity), in thousands



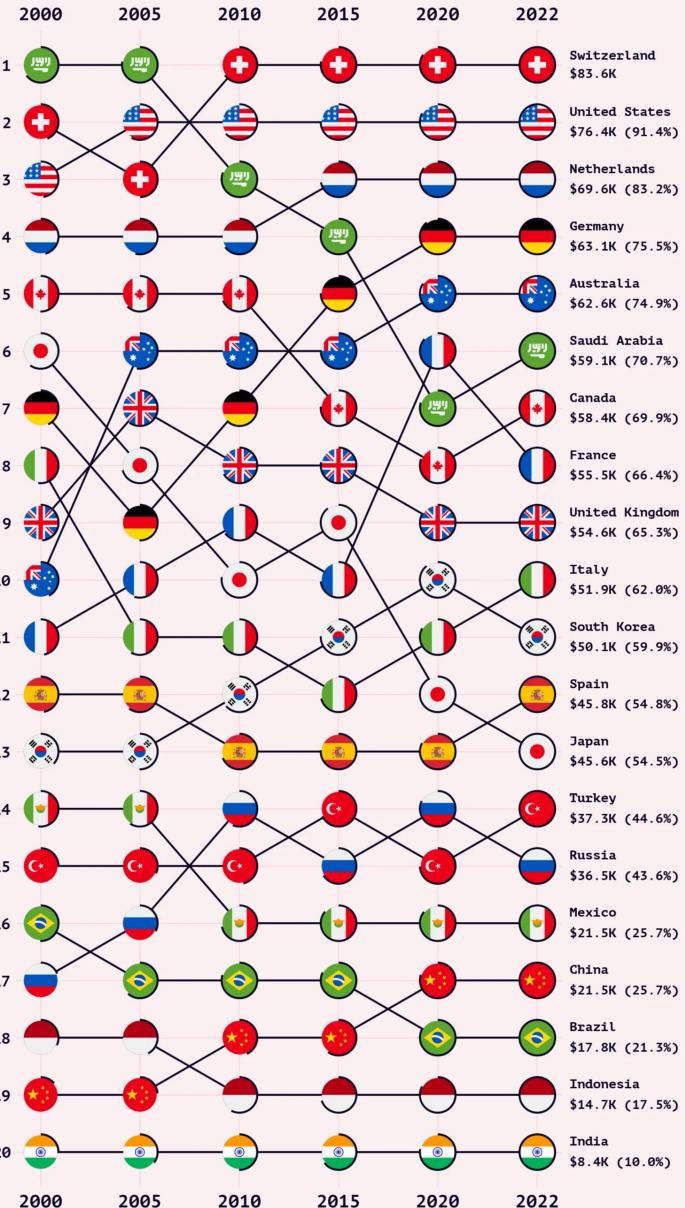
Example: GDP

GLOBAL GDP 2021

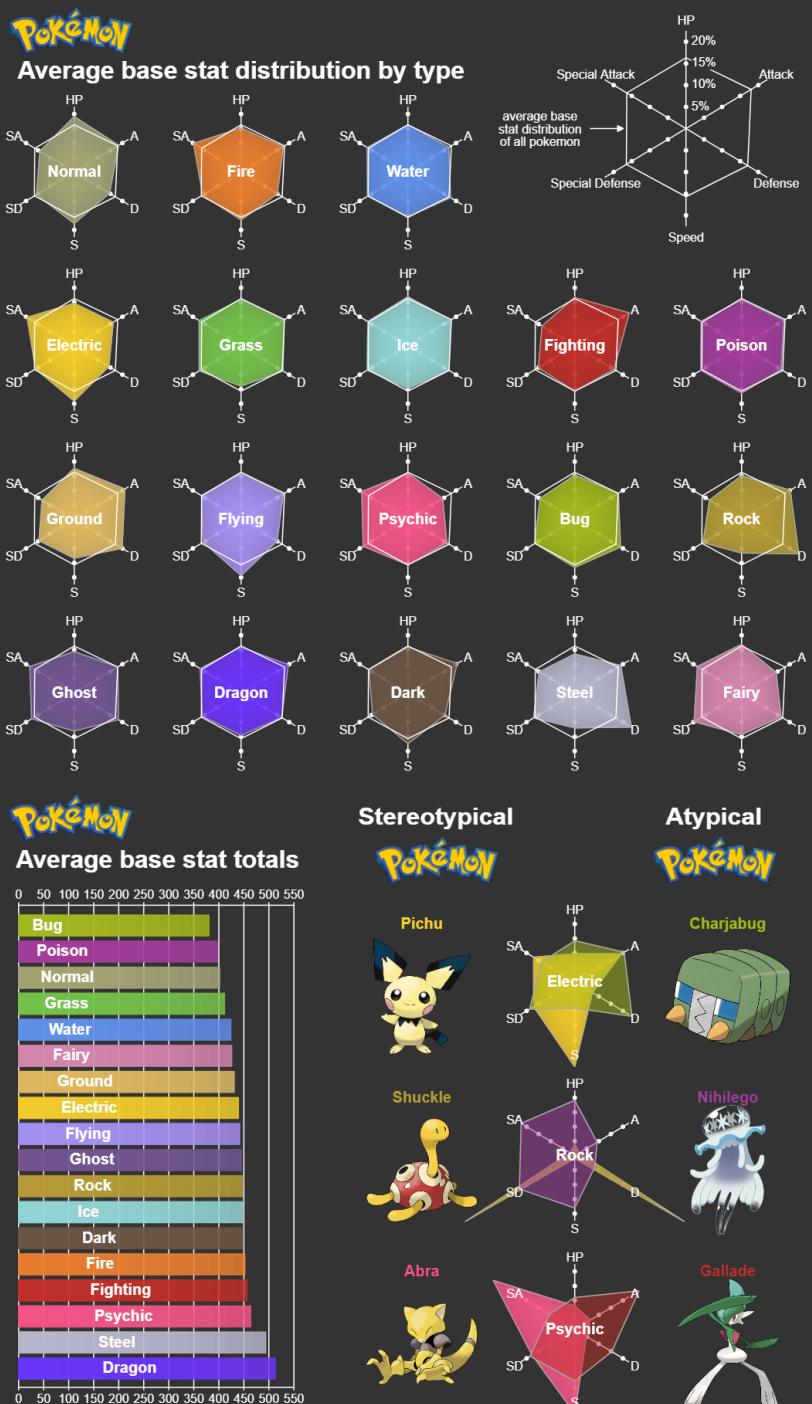
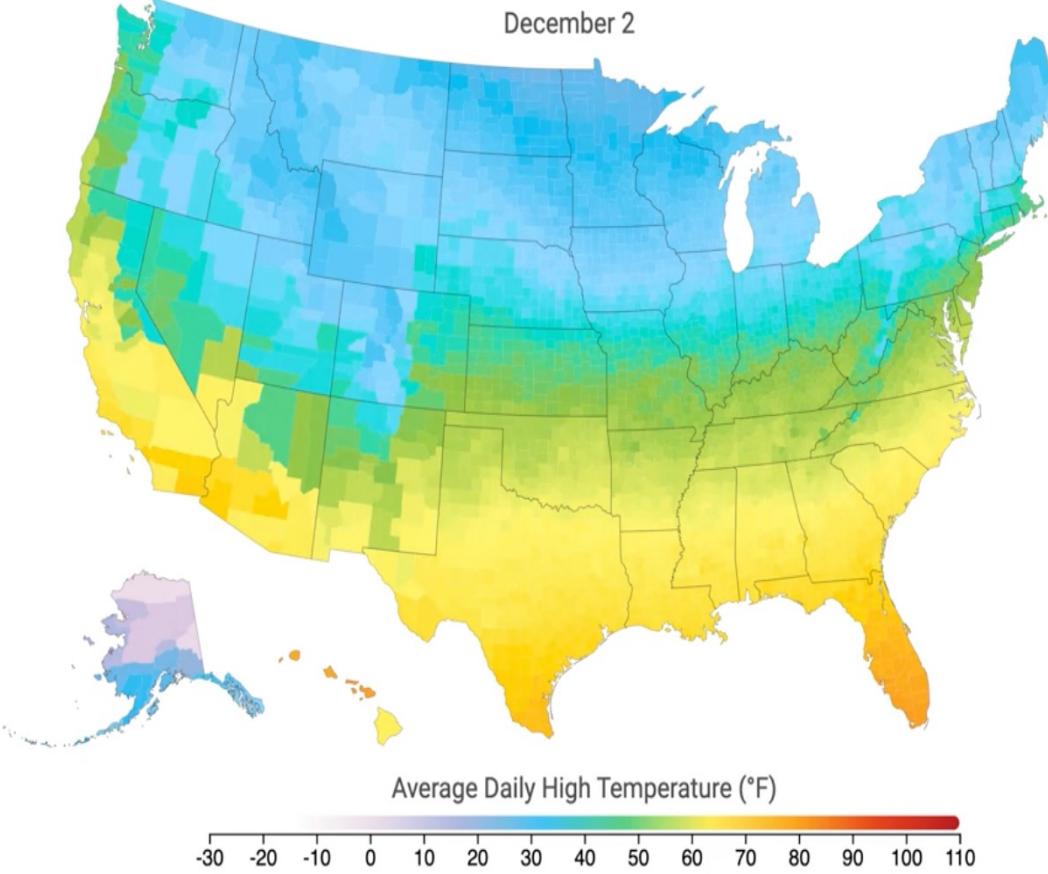


Example: GDP

Comparing Today's Largest Economies
GDP per capita, PPP (current \$)



Data source: <https://data.worldbank.org/>



Core

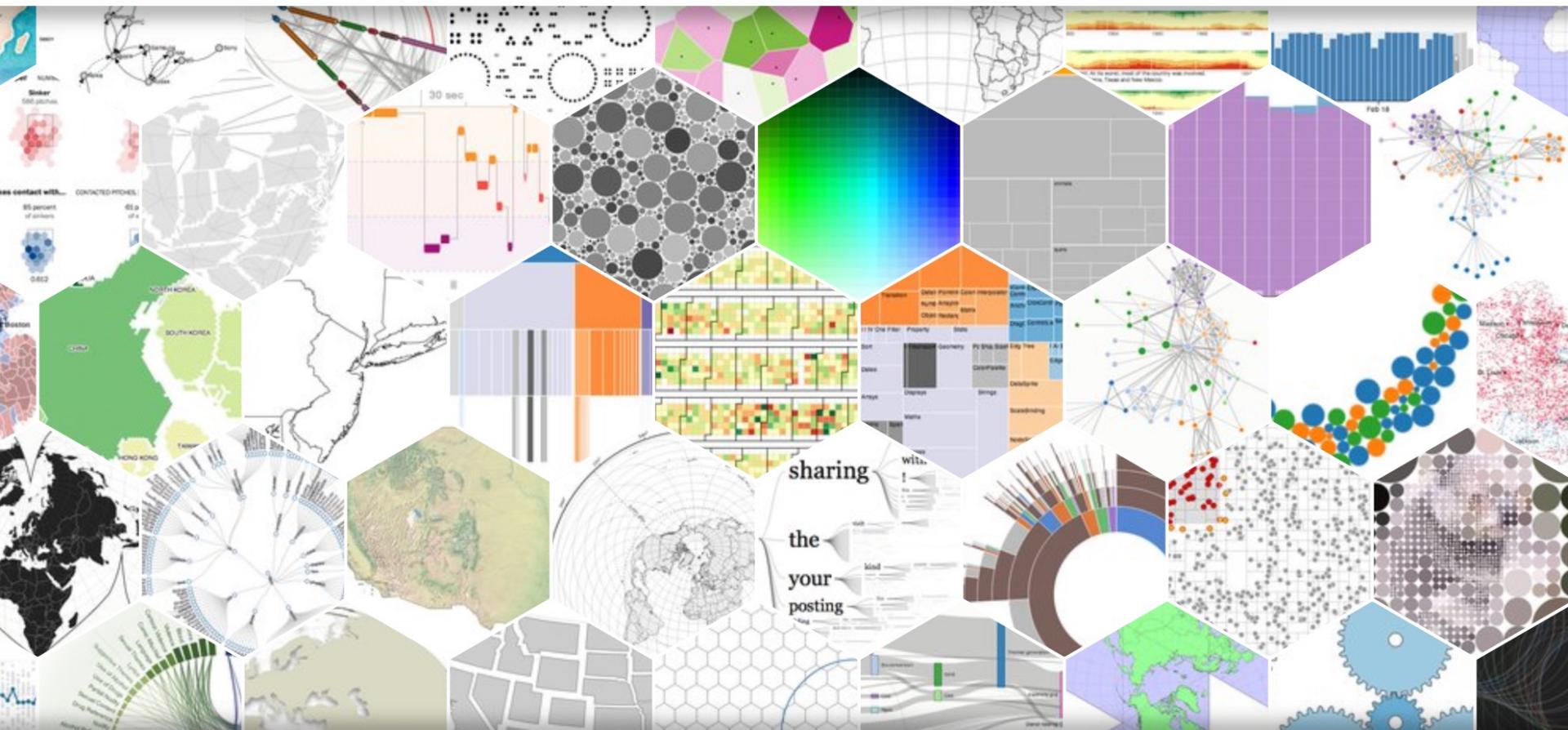
- **Data Foundations** - information in numbers
- **Context and Understanding** - contributes to problem solving
- A visualization has a mission:
- Put it simple,
 - we can characterize the data from a straightforward perspective.
 - We would like to visualize the analytics results to help others understand.

Tools

- **Excel** - the business standard
- **SPSS** and **SAS** – commercial
 - Heavily used in social sciences
- **R** – statistician (data scientist) standard
- **Python** – data scientist standard
- **D3** – computer scientist standard
 - Based on JavaScript - interactive
- **Tableau** – a beautiful enterprise standard



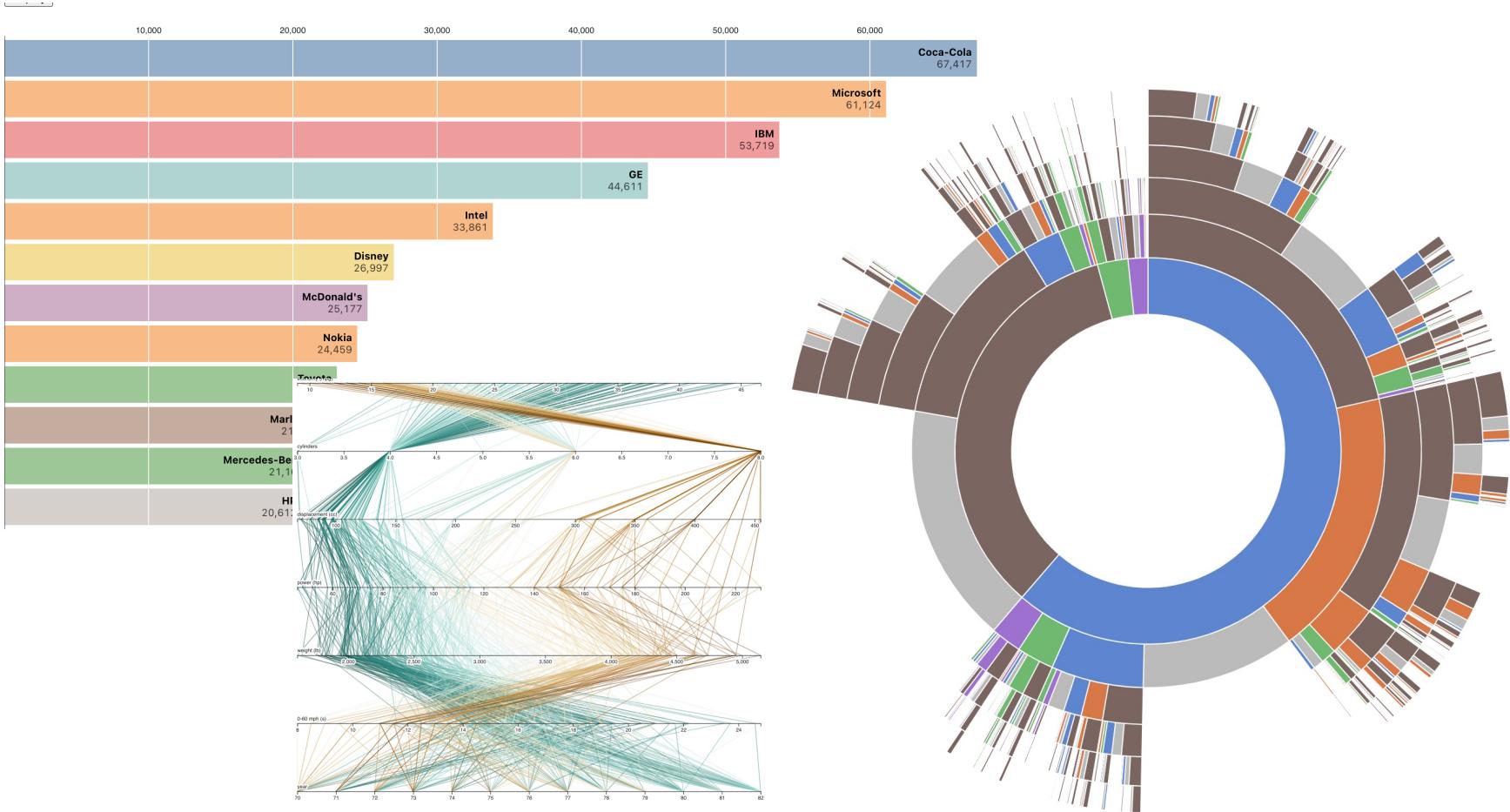
Data-Driven Documents



<https://d3js.org/>

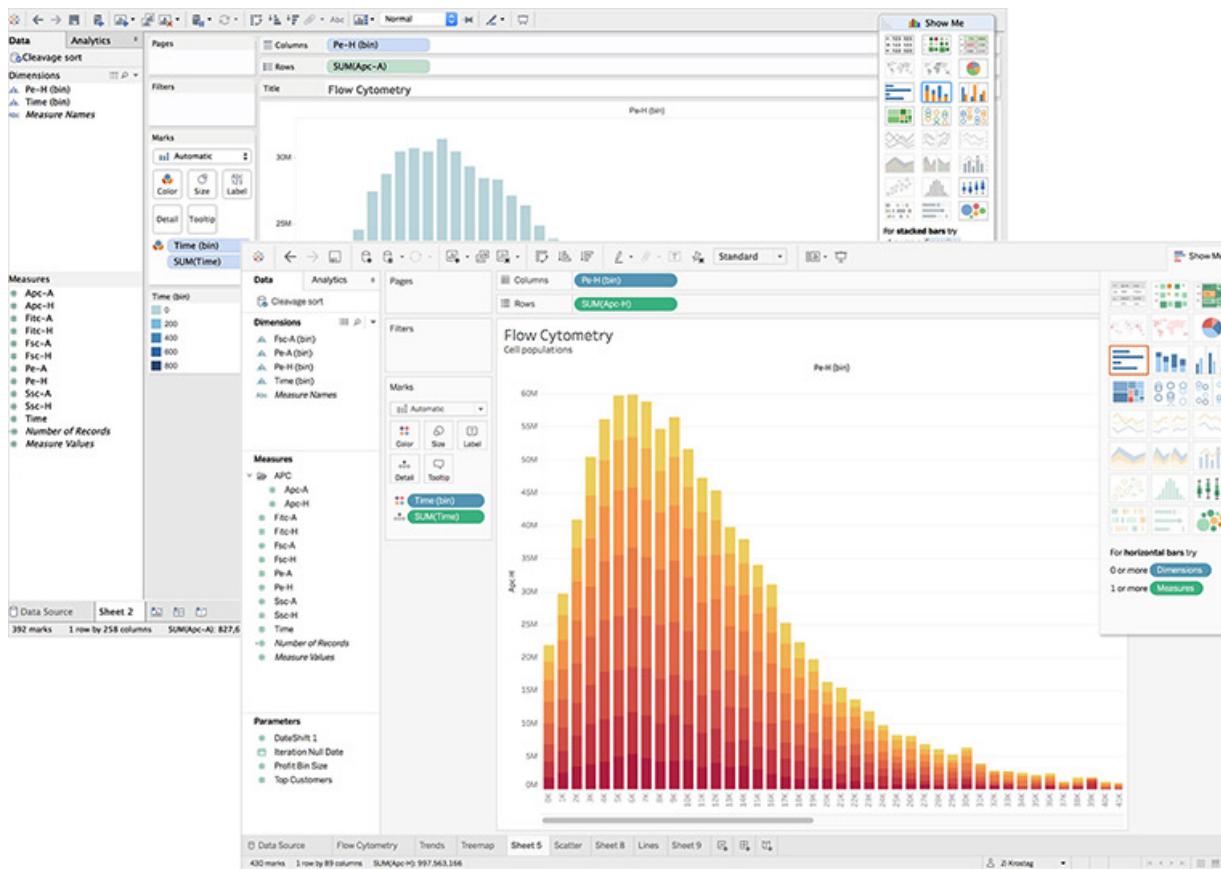
Example:

- https://observablehq.com/@d3/gallery?utm_source=d3js-org&utm_medium=hero&utm_campaign=try-observable



Commercial software

- Tableau (Free for student)



Pie Chart

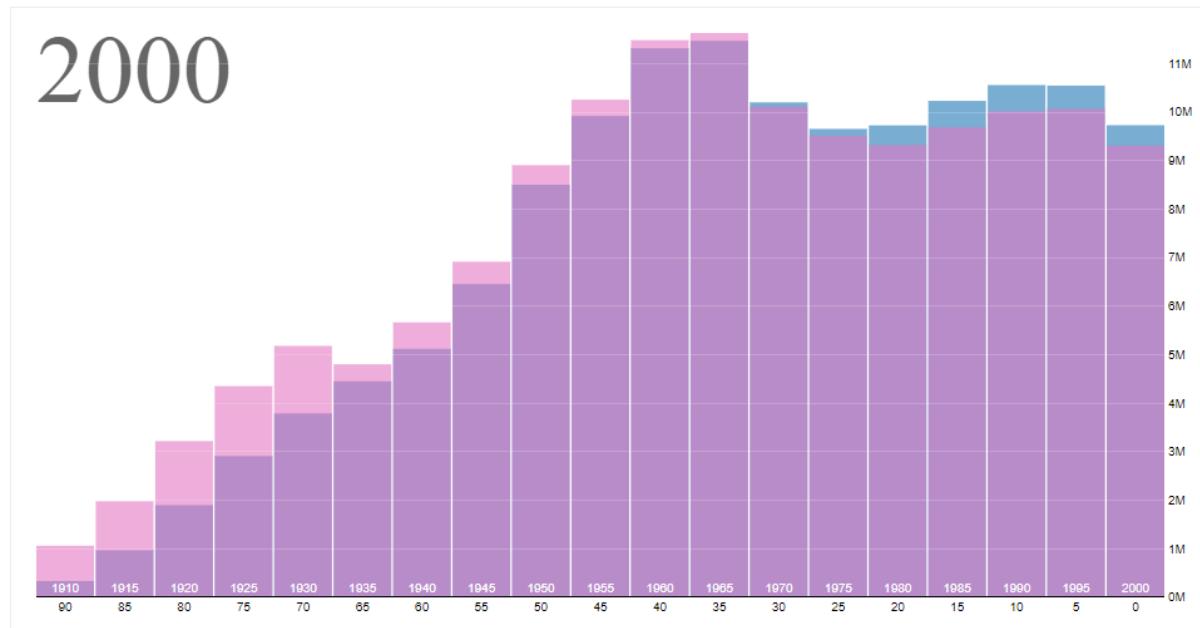
Coca Cola, PepsiCo and Keurig Dr Pepper together hold almost 94% of the carbonated soft drink (CSD) market, owning 50+ brands between them



Bar Chart

- This diagram shows the distribution of age groups in the United States.

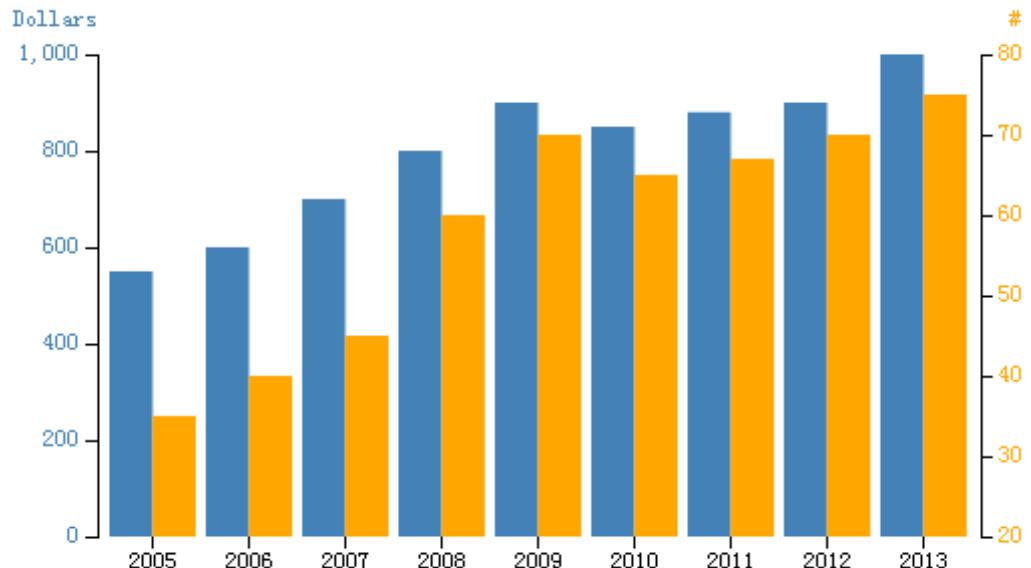
Population Pyramid



- <https://bl.ocks.org/mbostock/4062085>

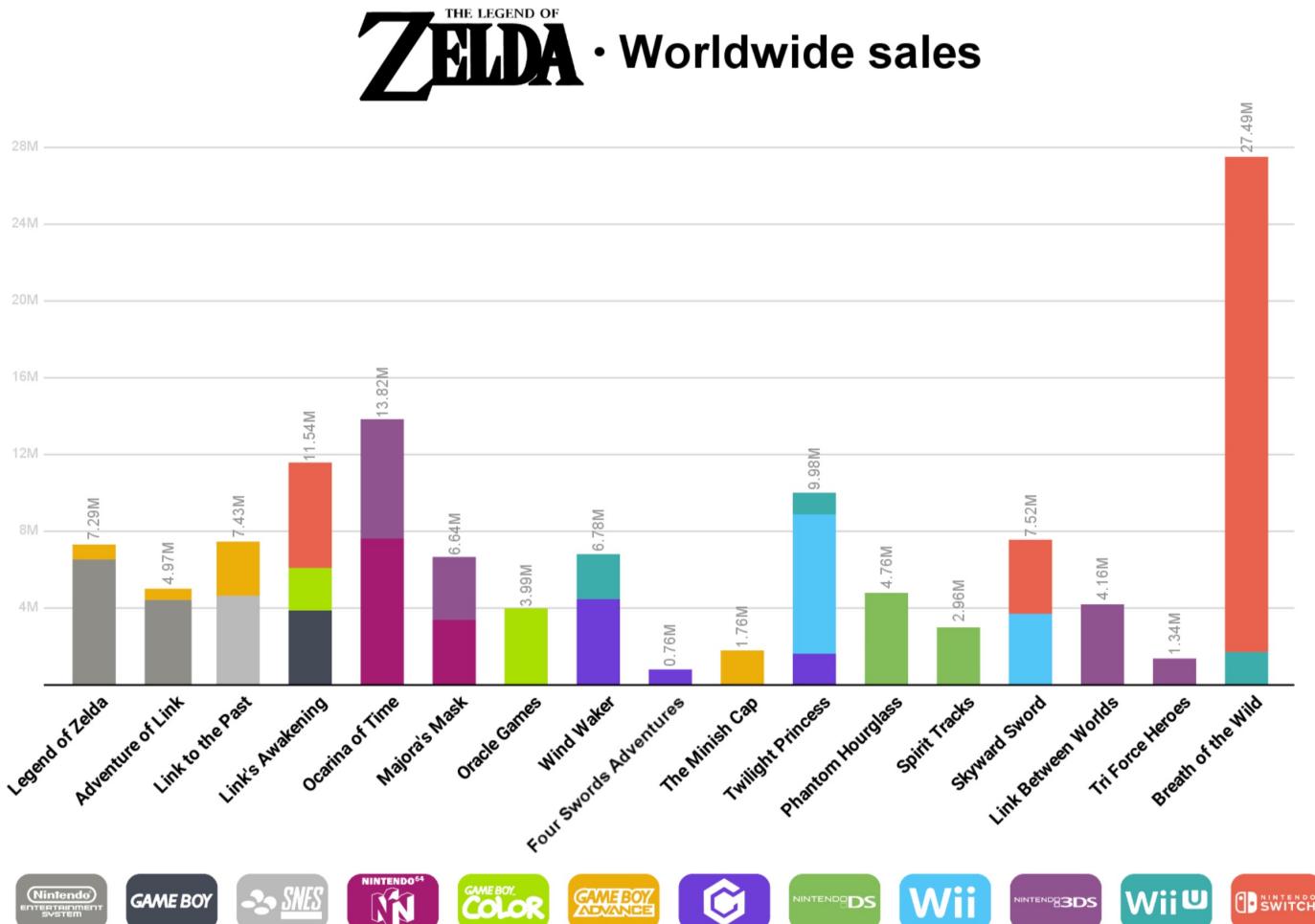
Bar Chart

- Double Bar Char



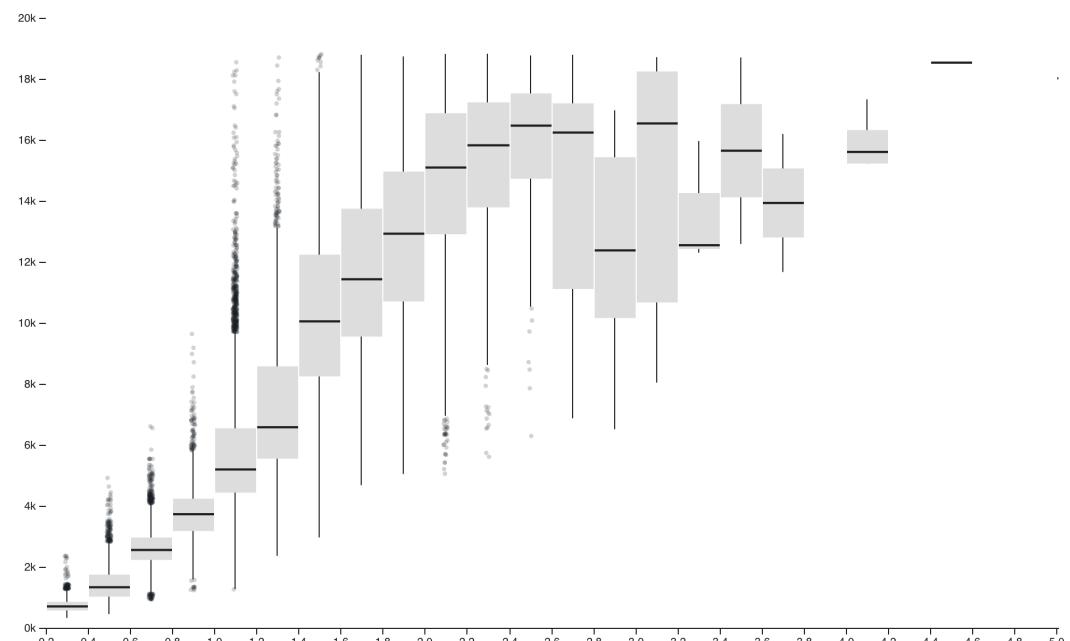
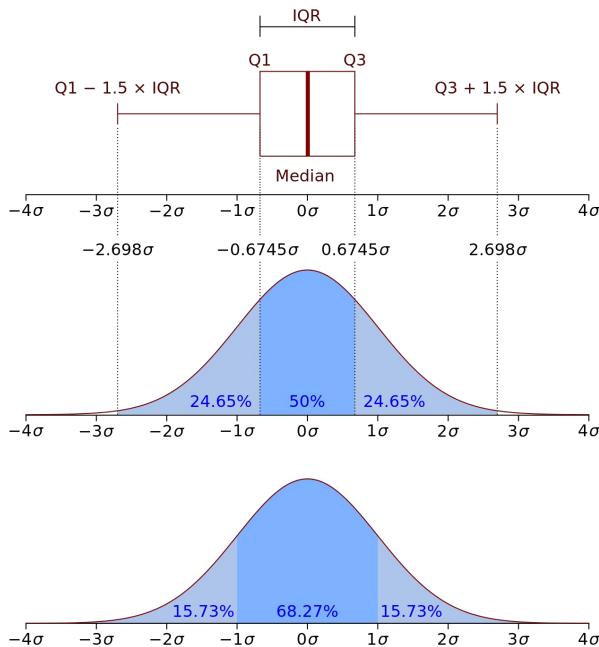
- <https://github.com/liufly/Dual-scale-D3-Bar-Chart>

Bar Chart



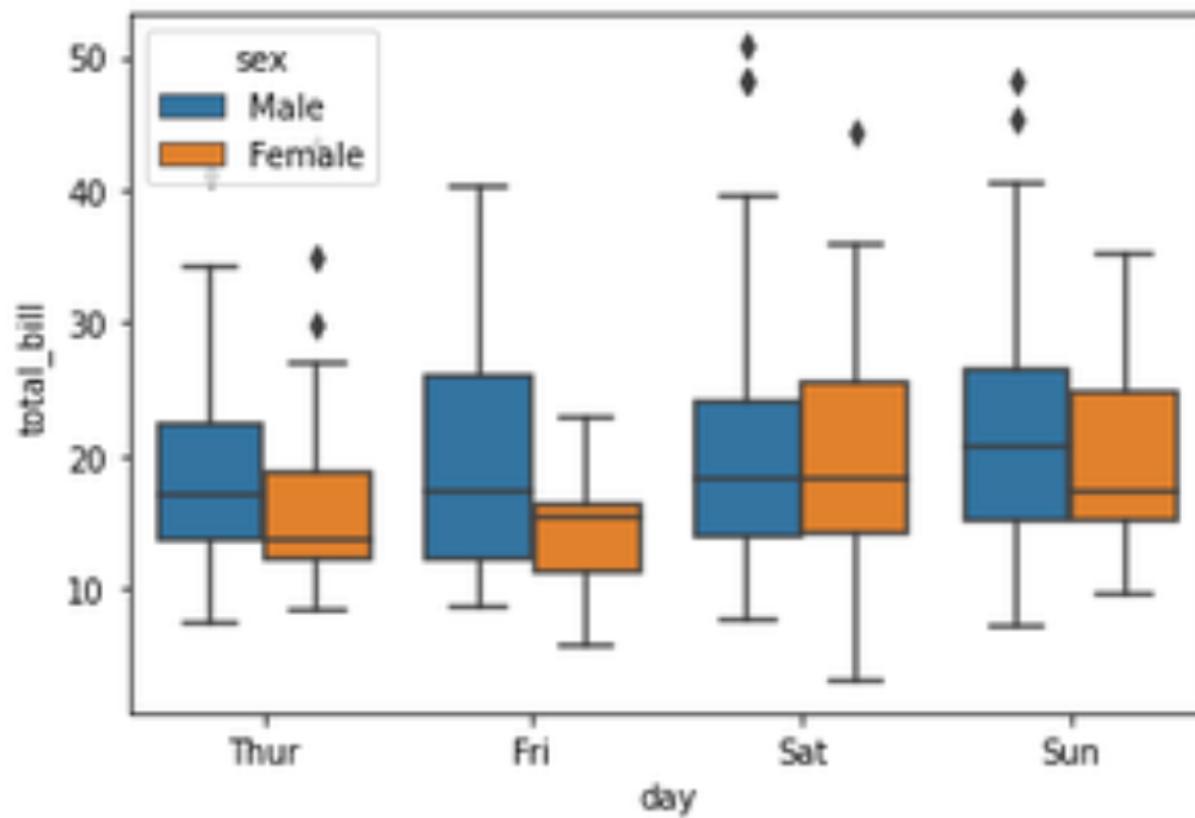
Box Plot

- In descriptive statistics, a box plot or boxplot is a method for graphically depicting groups of numerical data through their quartiles.



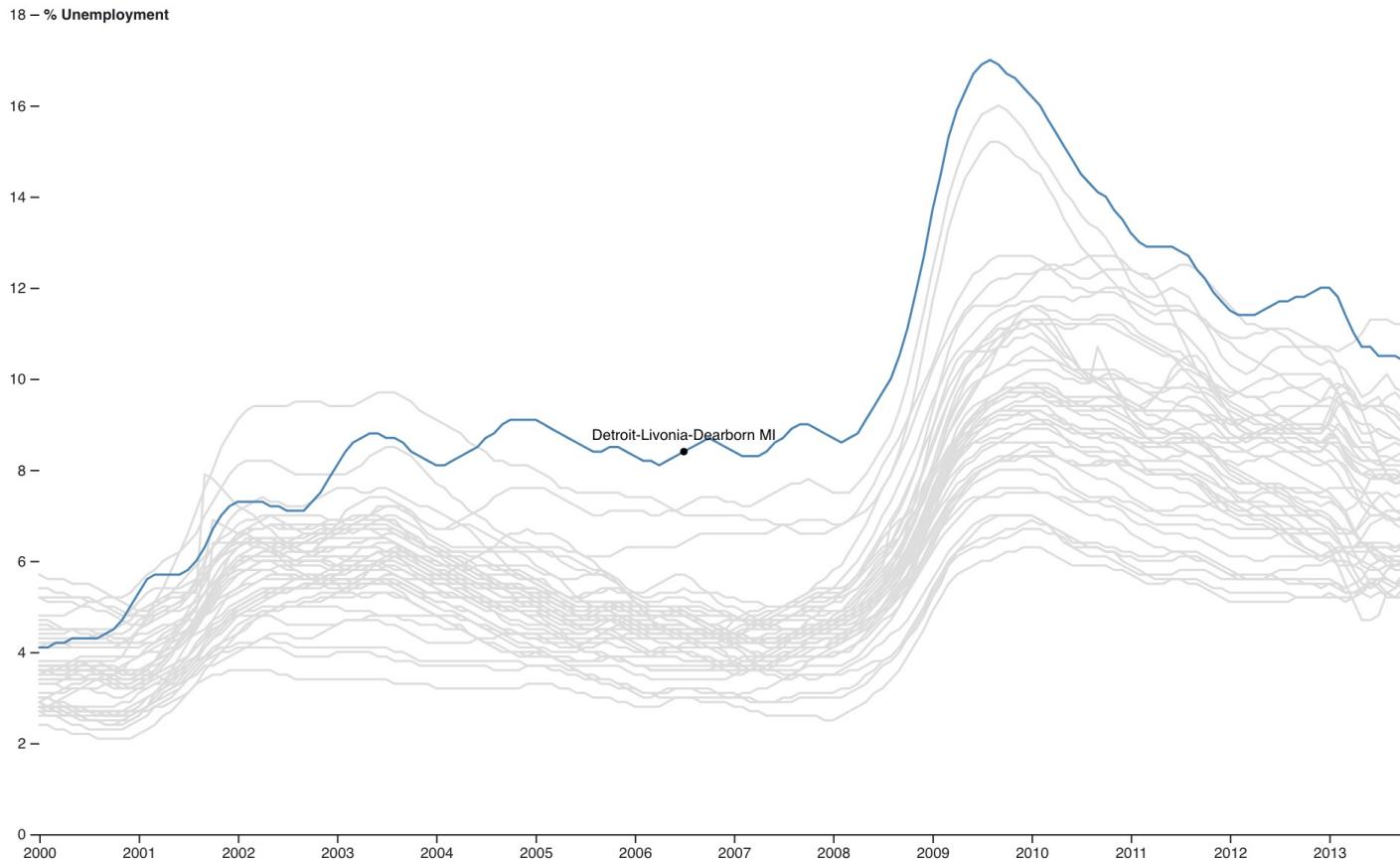
- <https://beta.observablehq.com/@mbostock/d3-box-plot>

Box Plot



<https://rdrr.io/cran/reshape2/man/tips.html>

Time Series

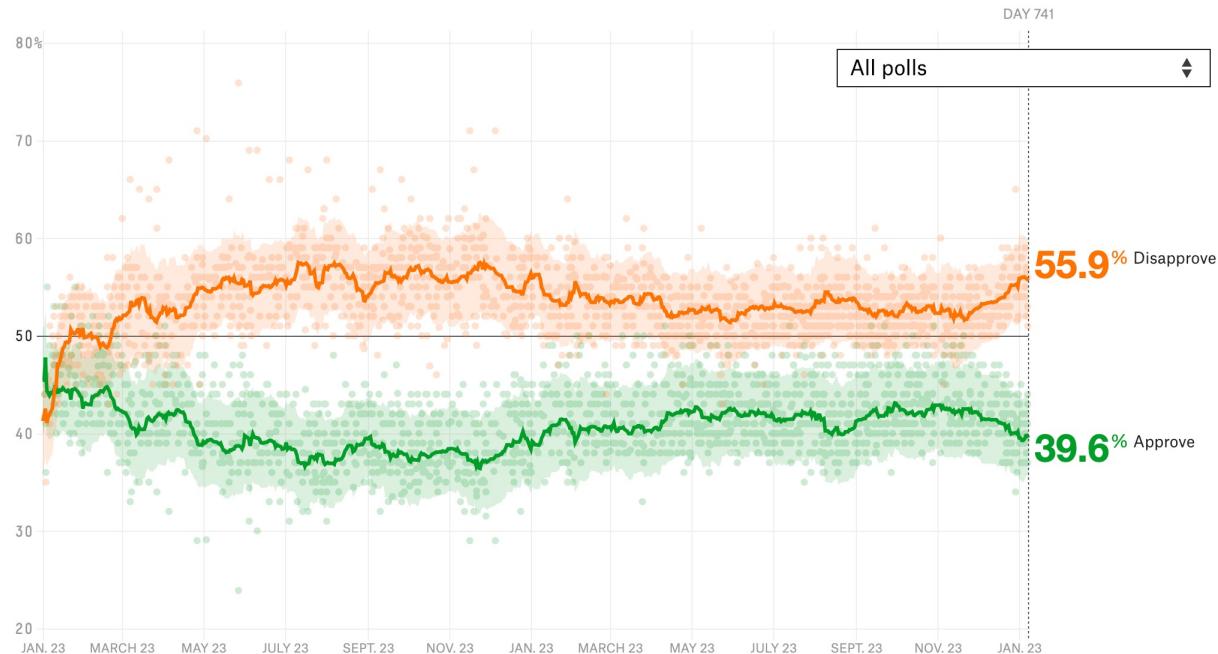


- <https://beta.observablehq.com/@mbostock/d3-multi-line-chart>

Time Series

How unpopular is Donald Trump?

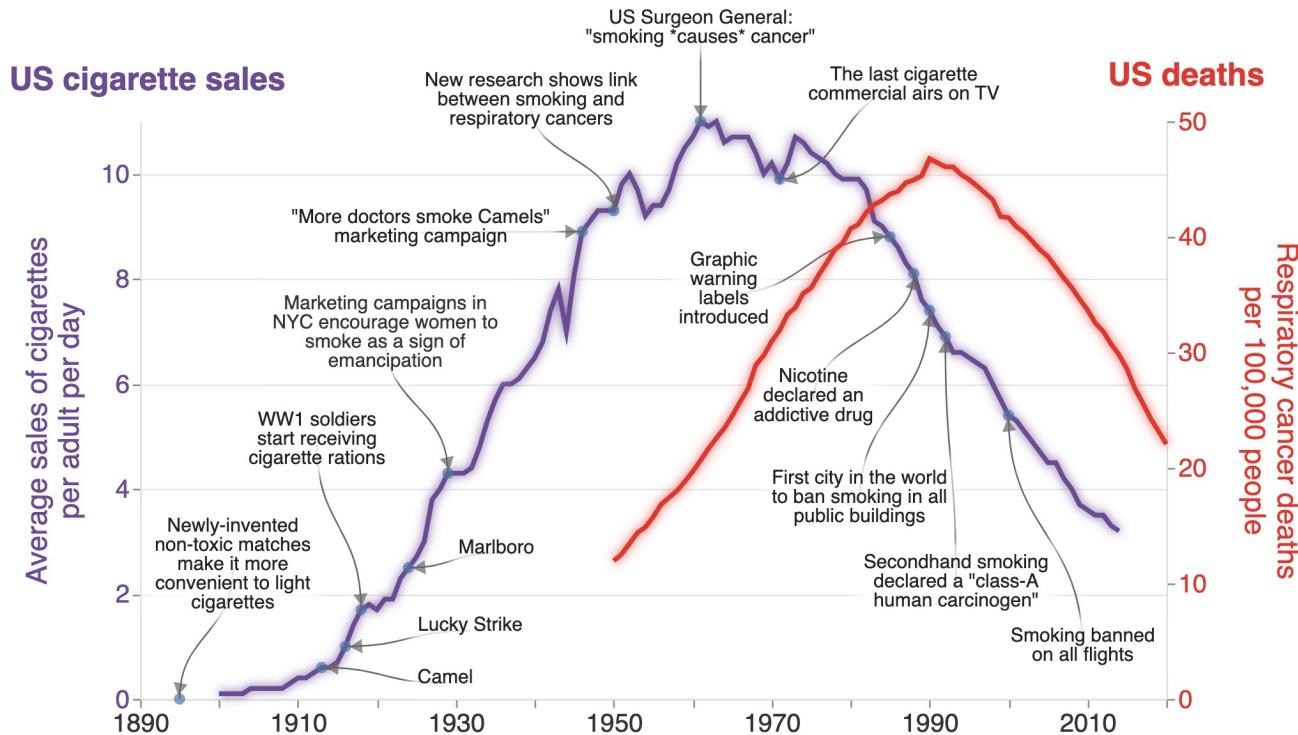
An updating calculation of the president's approval rating, accounting for each poll's quality, recency, sample size and partisan lean. [How this works »](#)



- <https://projects.fivethirtyeight.com/trump-approval-ratings>

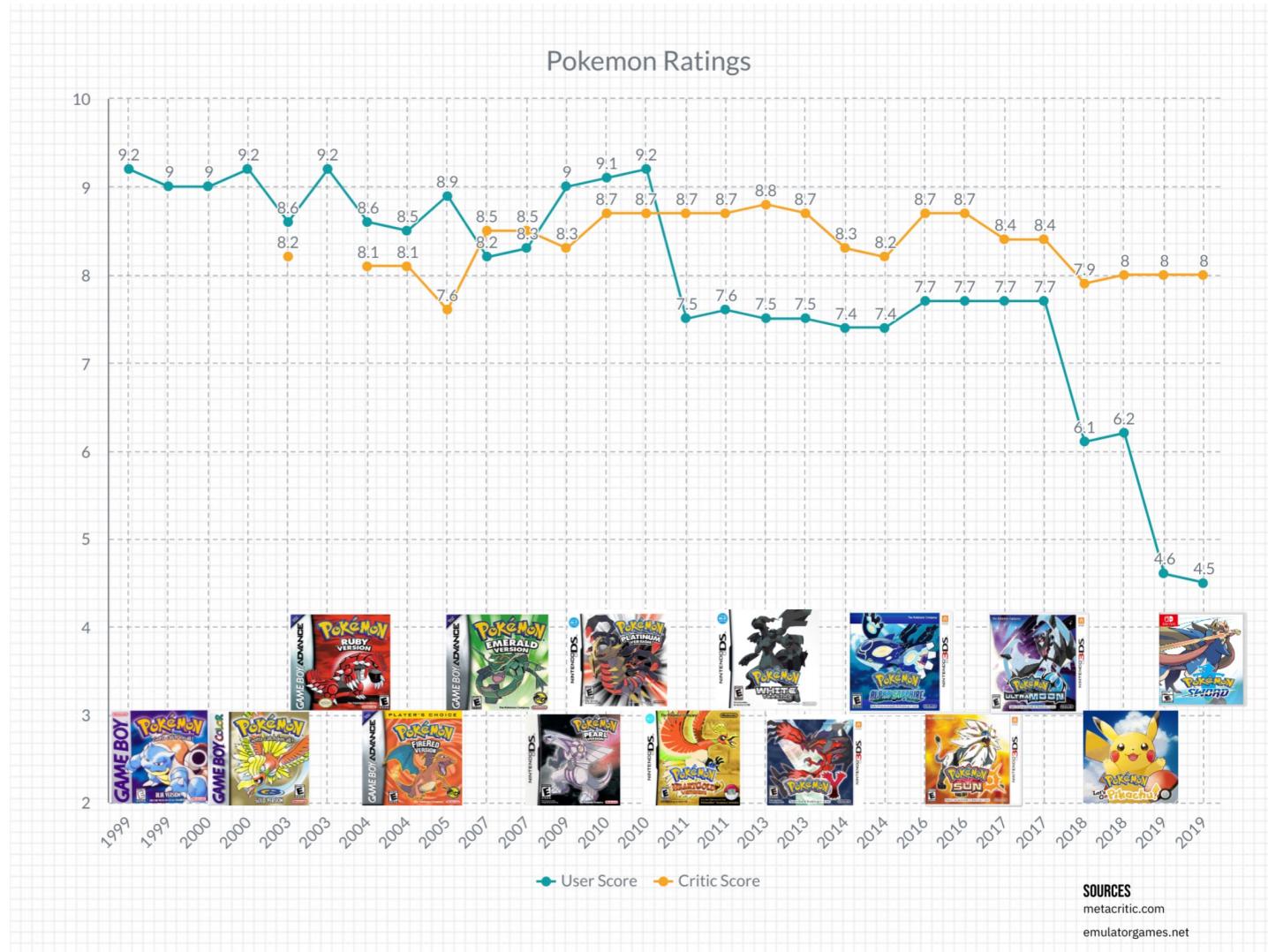
Time Series

A History of Cigarette Sales and Lung Cancer Deaths in the US

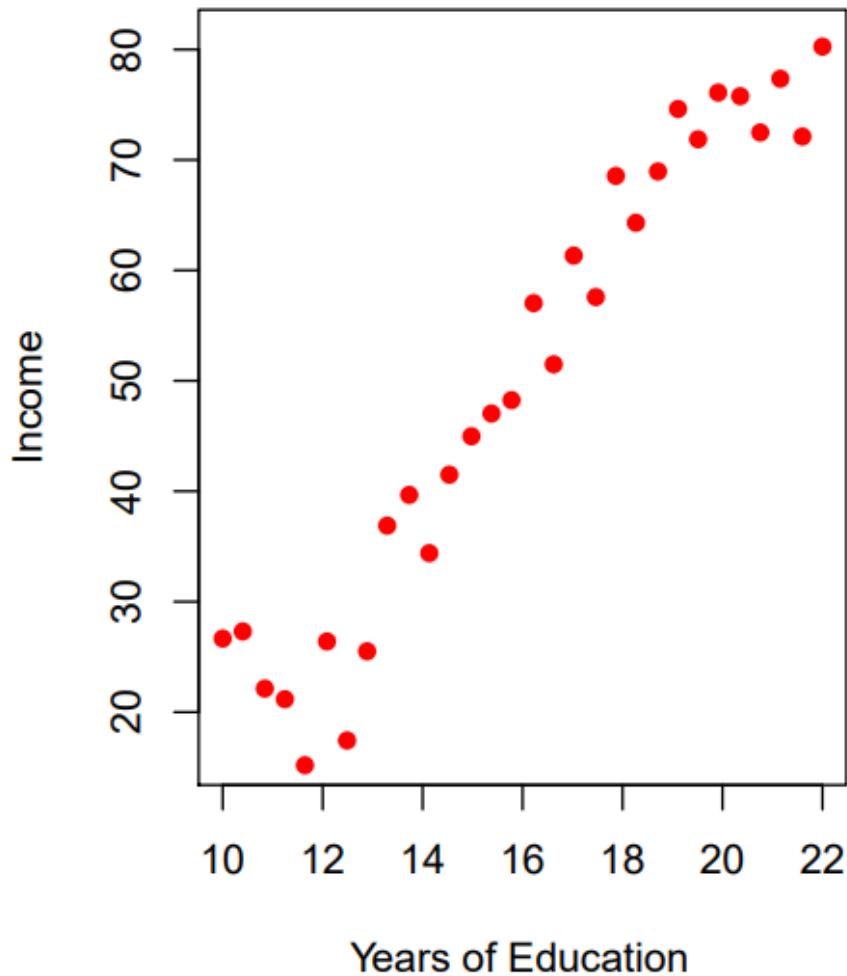


<https://www.contextualize.ai/mpereira/a-history-of-cigarette-sales-and-lung-cancer-deaths-in-the-us-ce8dea7a>

Time Series

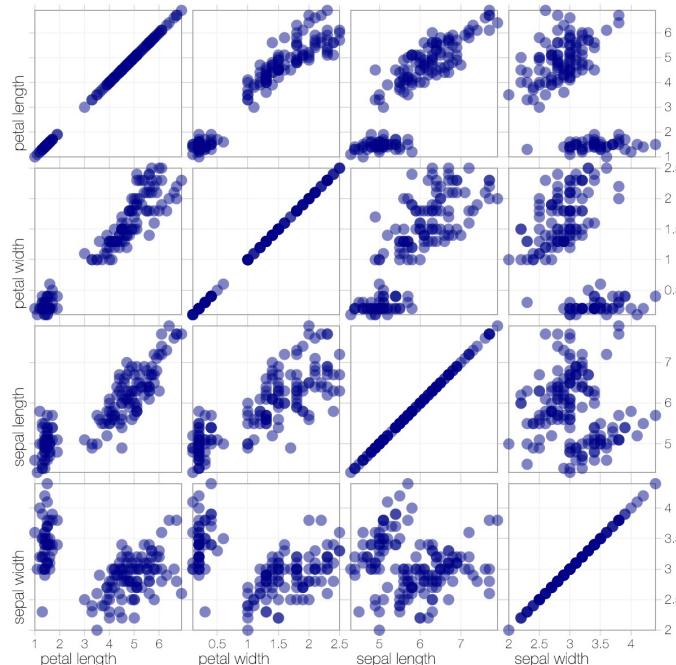


Scatterplot



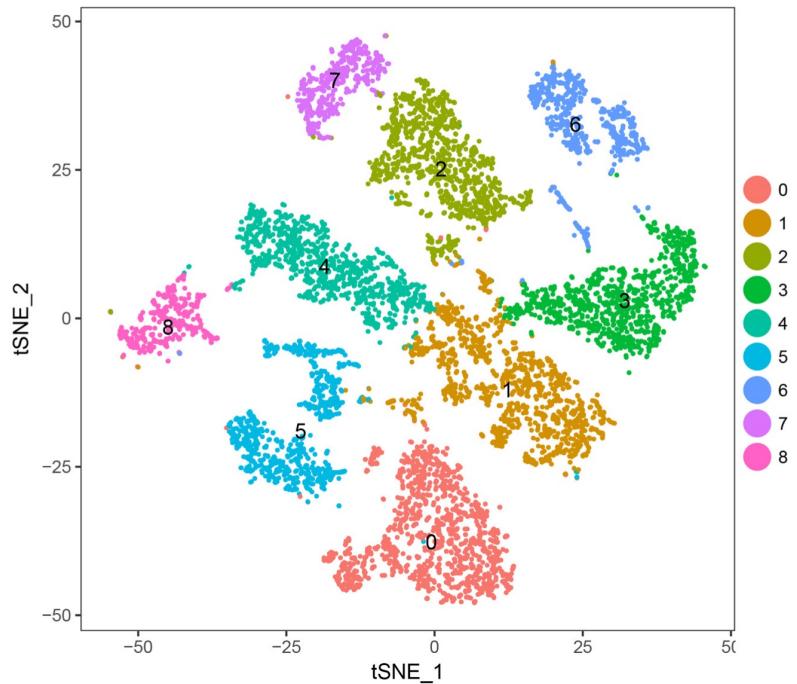
Scatterplot Matrix

- Pairwise correlations for multi-dimensional data; (Thus, the diagonal cells are diagonal lines.)
- This matrix shows Anderson's data on iris flowers on the Gaspé Peninsula.

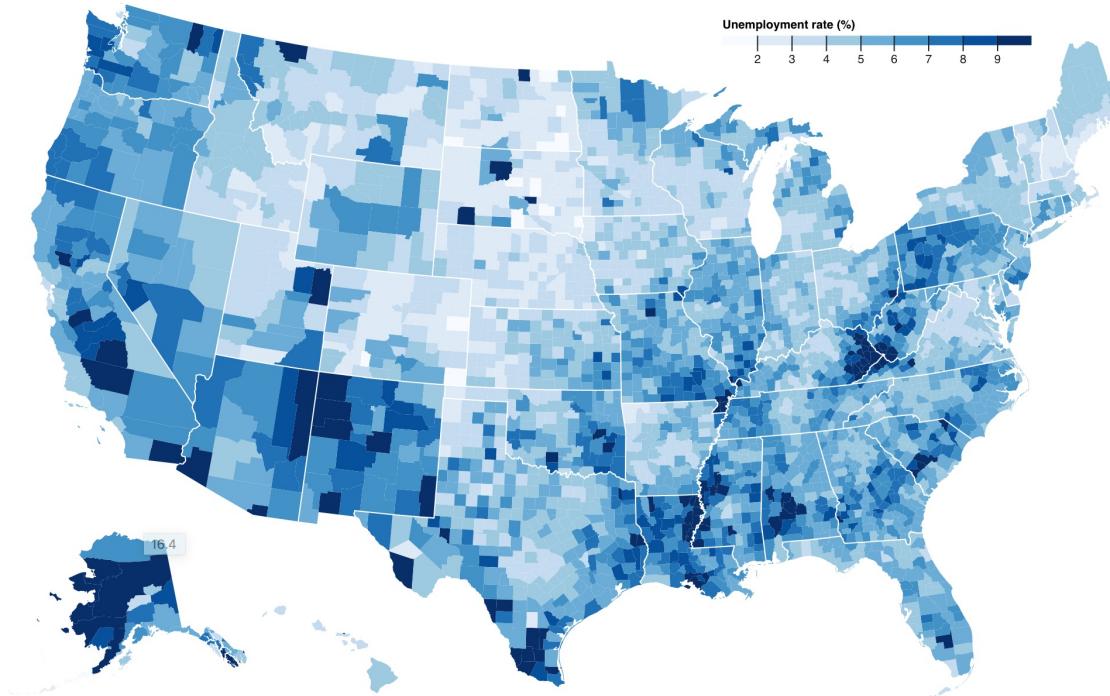


- <http://benjiec.github.io/scatter-matrix/demo/demo.html>

Scatter Plot for clustering



Map: Choropleth map



- <https://beta.observablehq.com/@mbostock/d3-choropleth>

Choropleth map: U.S. Disability Population Statistics

Data Dashboard of Disability Population Distributions

Map Value Format

Percent

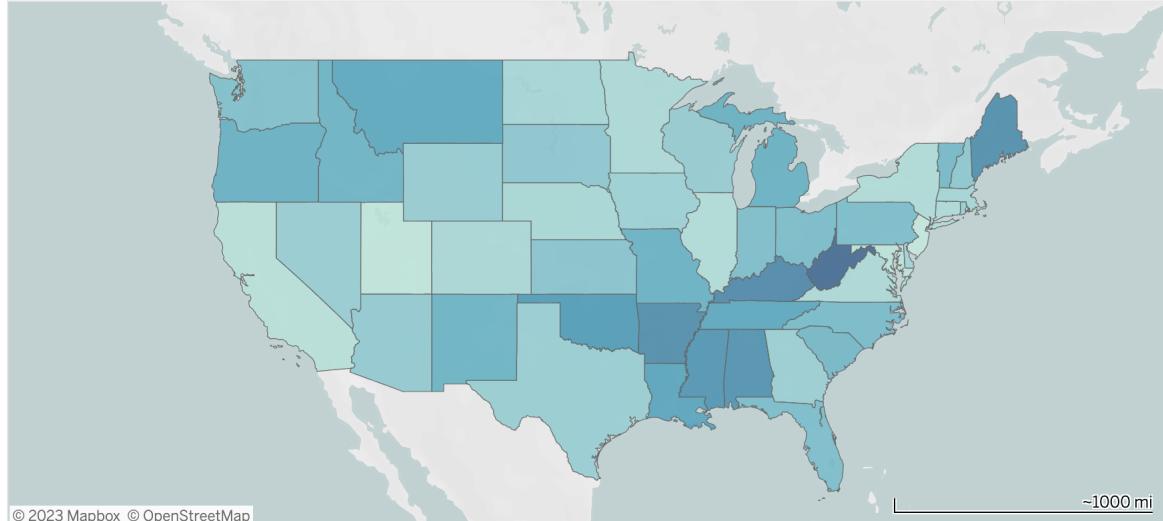
Number

Percent

4.44% 9.30%

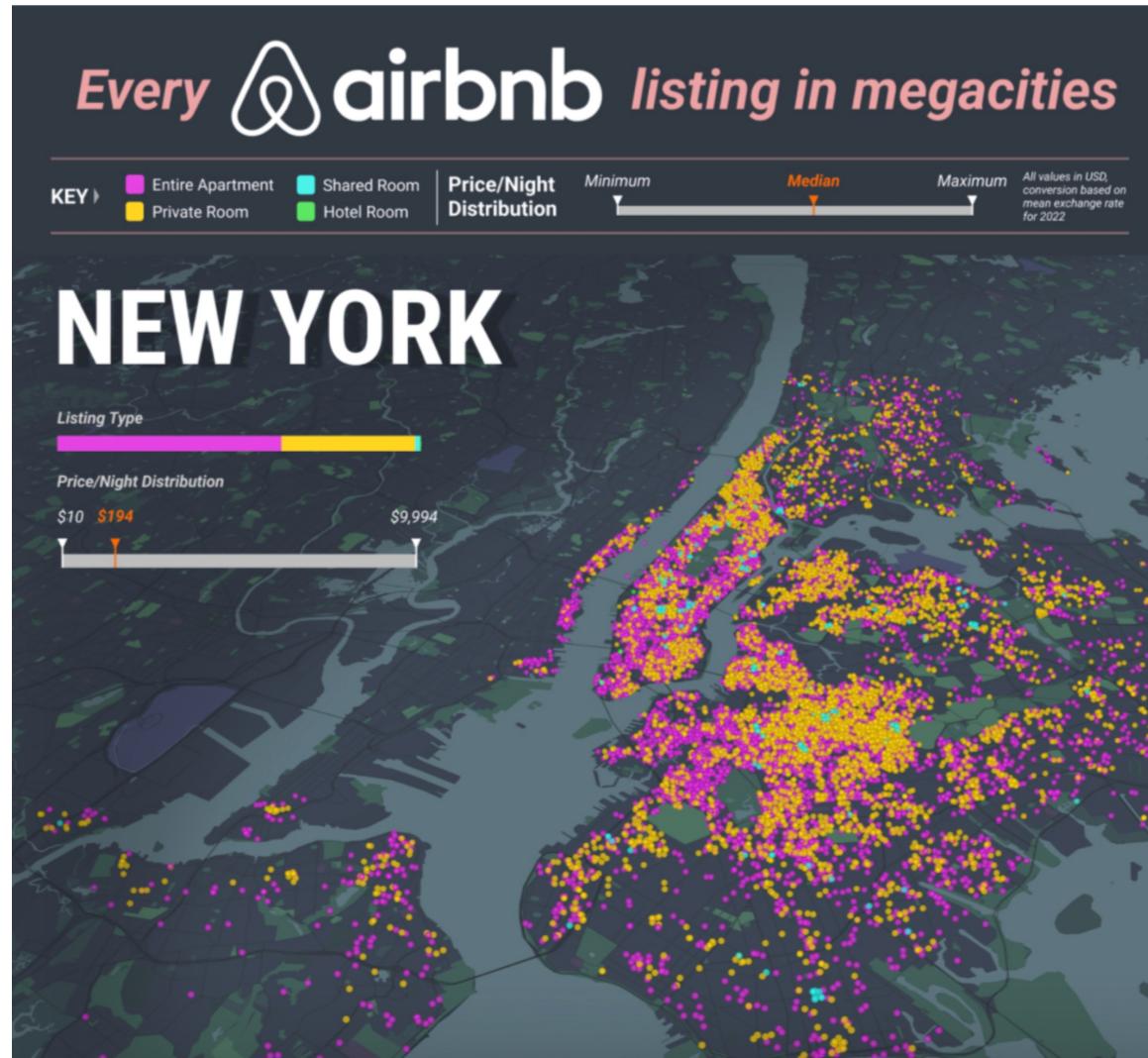
Number

Null



<https://public.tableau.com/app/profile/jiawei.wang8332/viz/DisabilityPopulationStatistics-Year11-14/MainDashboard>

Map plot



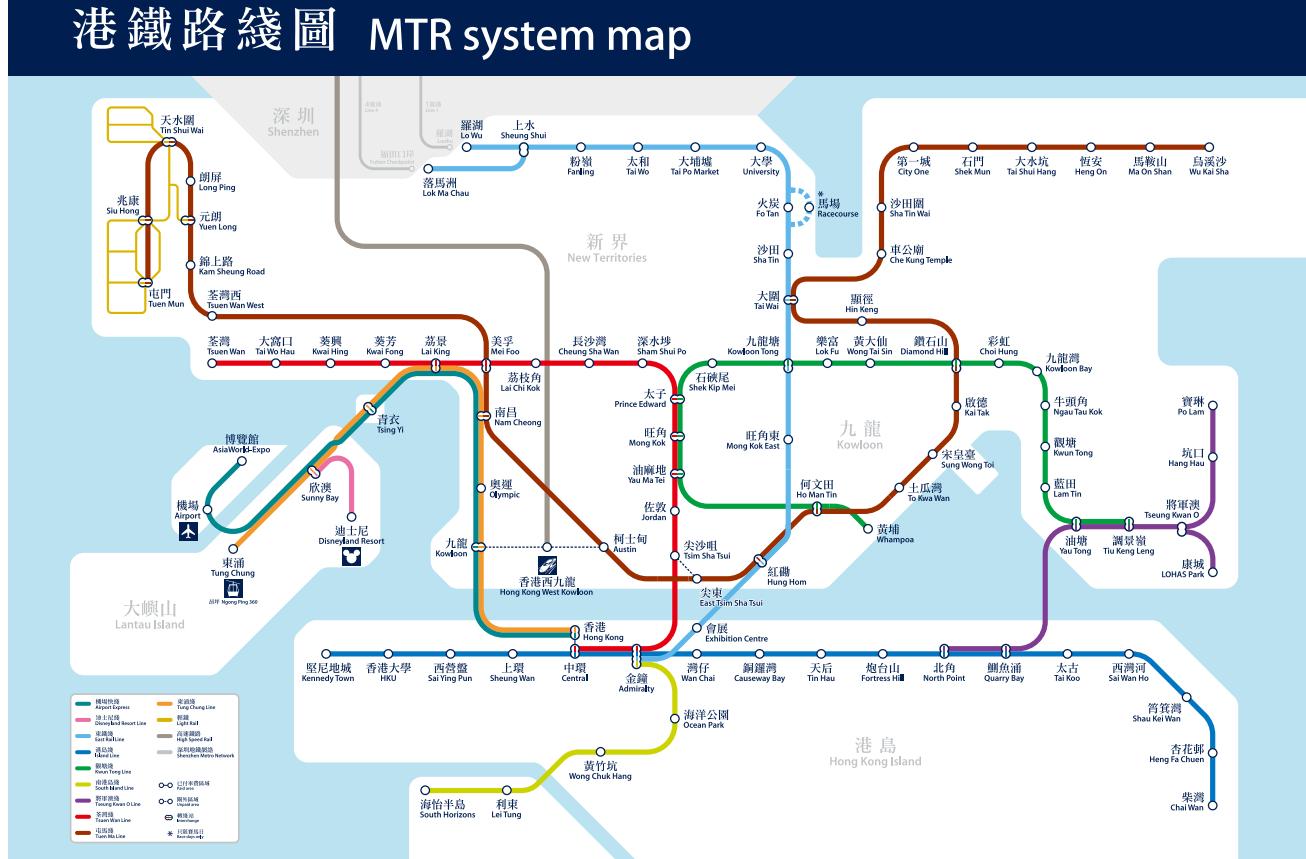
Transportation Map



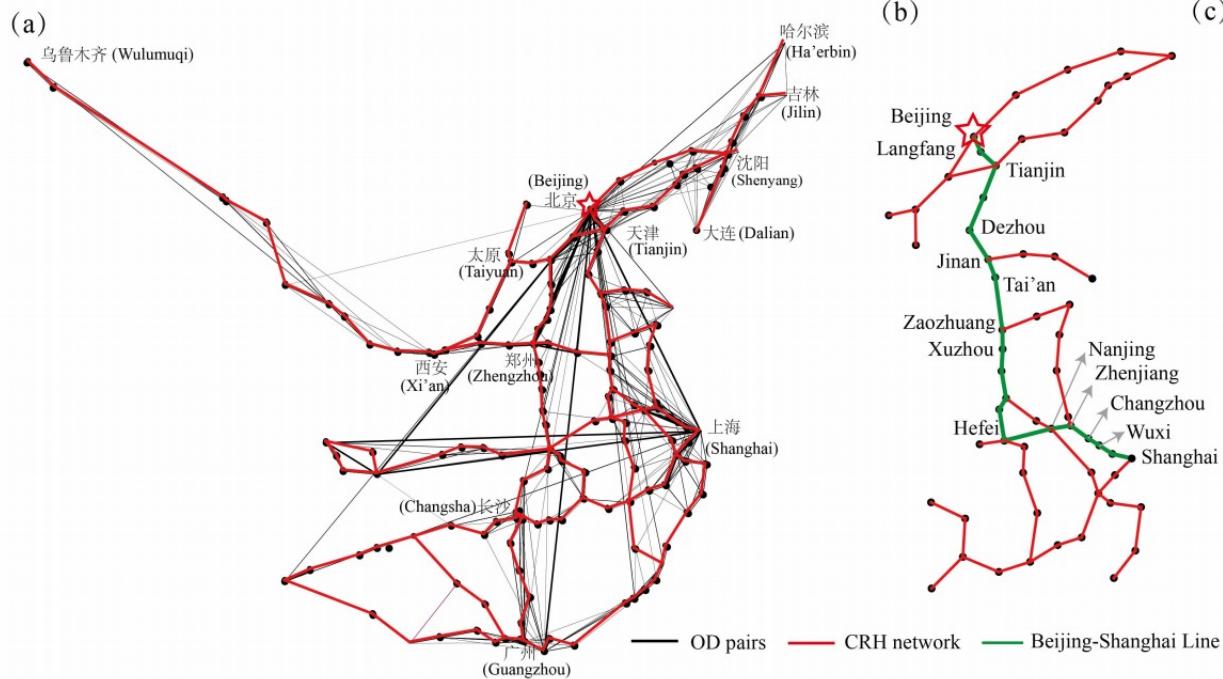
https://www.cathaypacific.com/cx/en_HK.html

Transportation Map

港鐵路線圖 MTR system map



Transportation Map

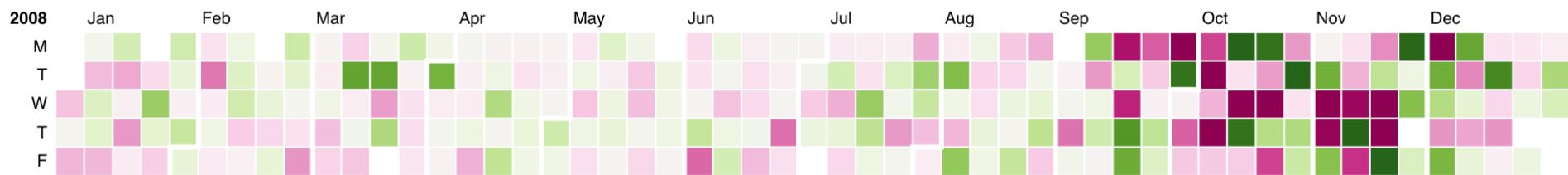
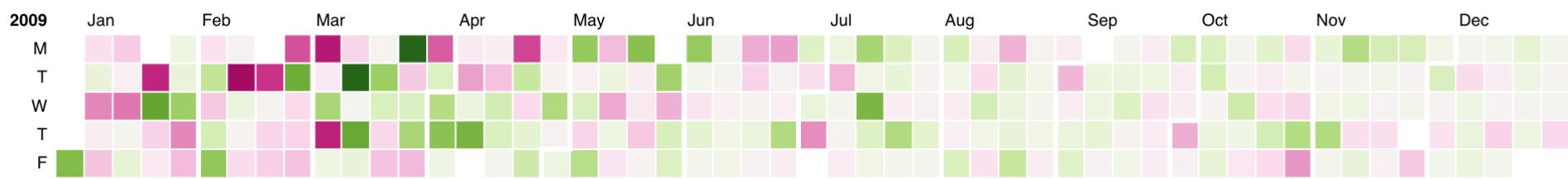


From CITY	To Beijing	From Beijing To CITY	
CITY	# of trips	CITY	# of trips
Langfang	7220	Langfang	8577
Baoding	6144	Tianjin	7675
Jinan	8532	Guangzhou	7619
Dezhou		Baoding	7219
Tai'an		Shenzhen	6585
Zaozhuang		Chengdu	6551
Xuzhou		Shanghai	9626
Nanjing		Chengdu	7278
Zhenjiang		Wuhan	6609
Changzhou		Shijiazhuang	6549
Wuxi		Zhengzhou	7829
Shanghai		Shijiazhuang	6856
		Shenzhen	7370
		Zhengzhou	5688

- <https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=8710615>

Calendar View

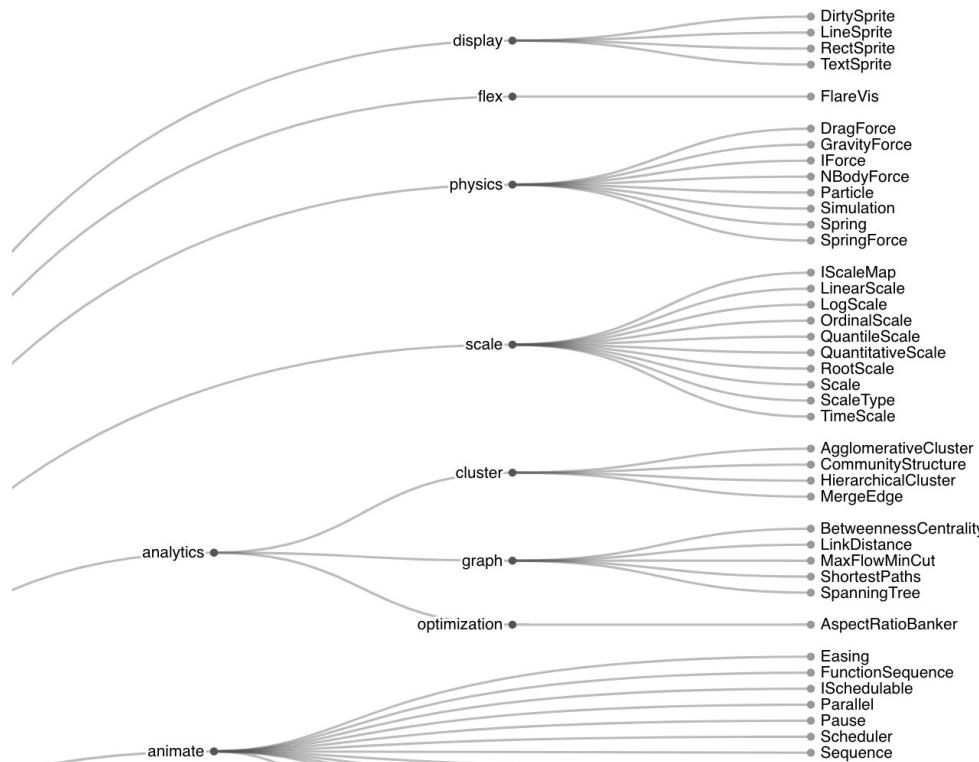
- This chart shows daily changes of the Dow Jones Industrial Average from 1990 to 2010. Days the index went up are green; days the index went down are pink.



- <https://beta.observablehq.com/@mbostock/d3-calendar-view>

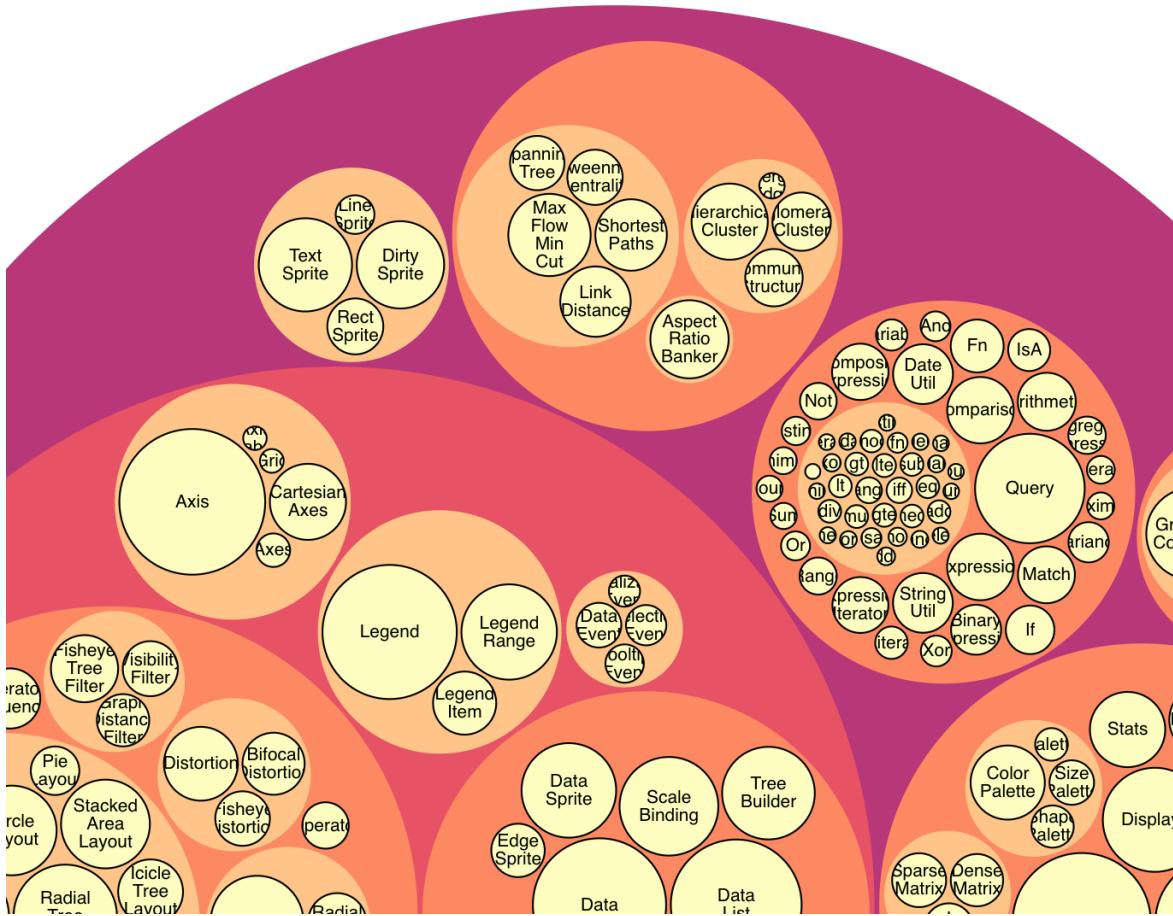
Hierarchical Diagram

- Tree – shows the organizational structure



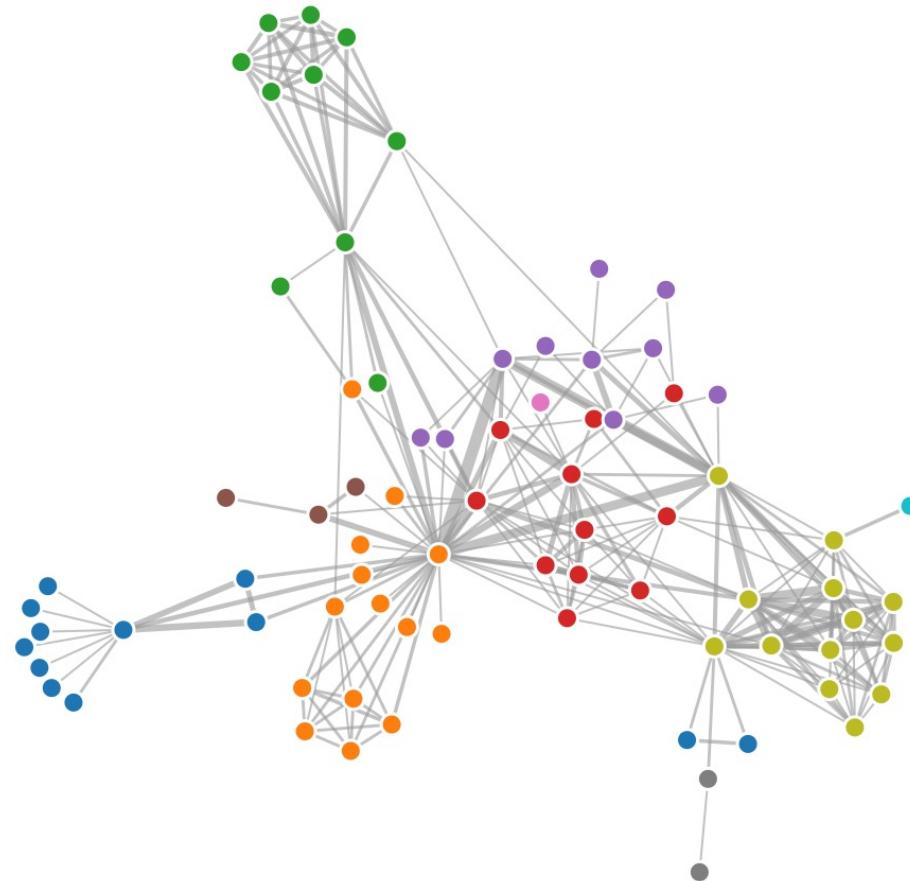
- <https://beta.observablehq.com/@mbostock/d3-cluster-dendrogram>

Circle Packing



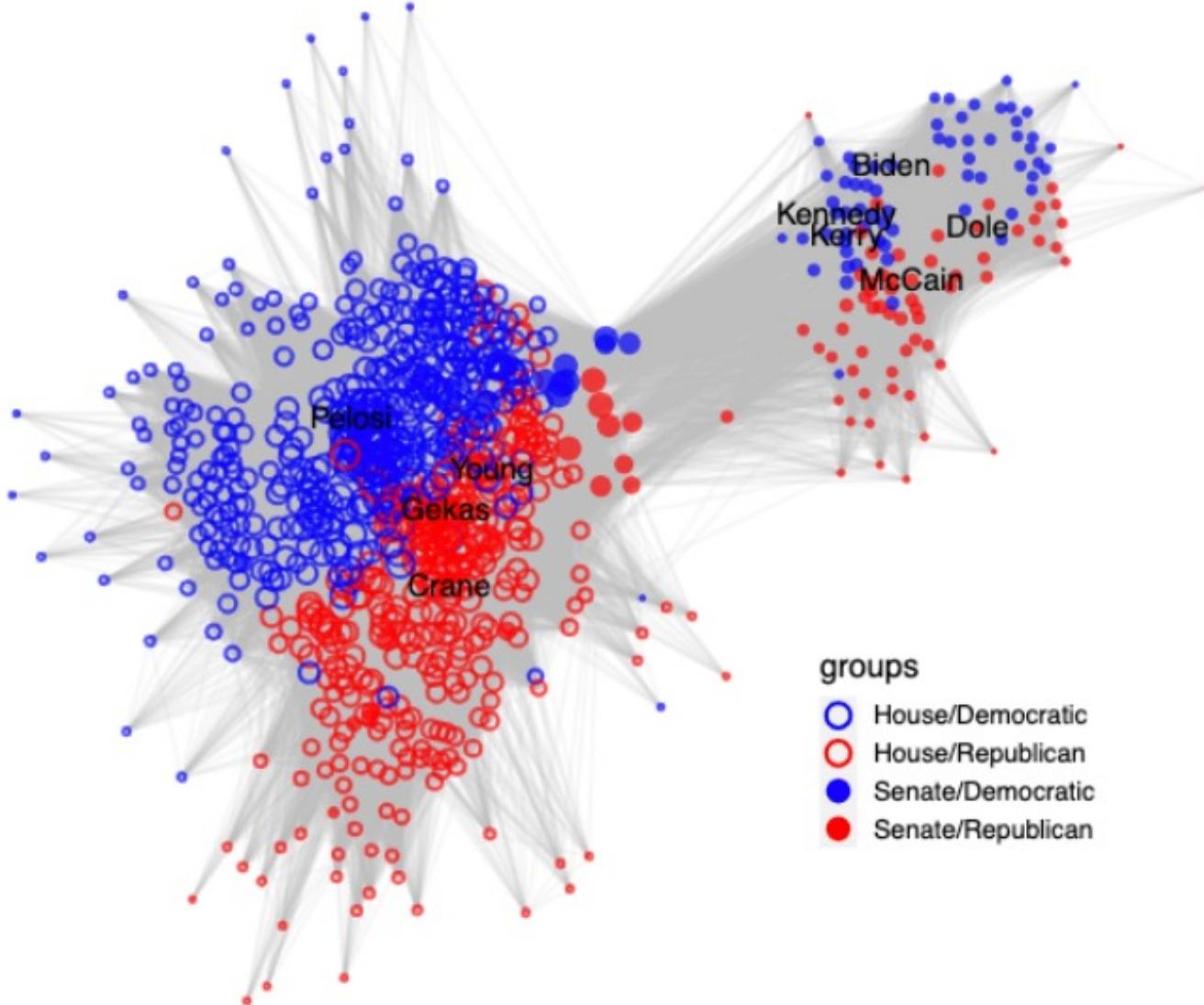
- <https://beta.observablehq.com/@mbostock/d3-circle-packing>

Social Network

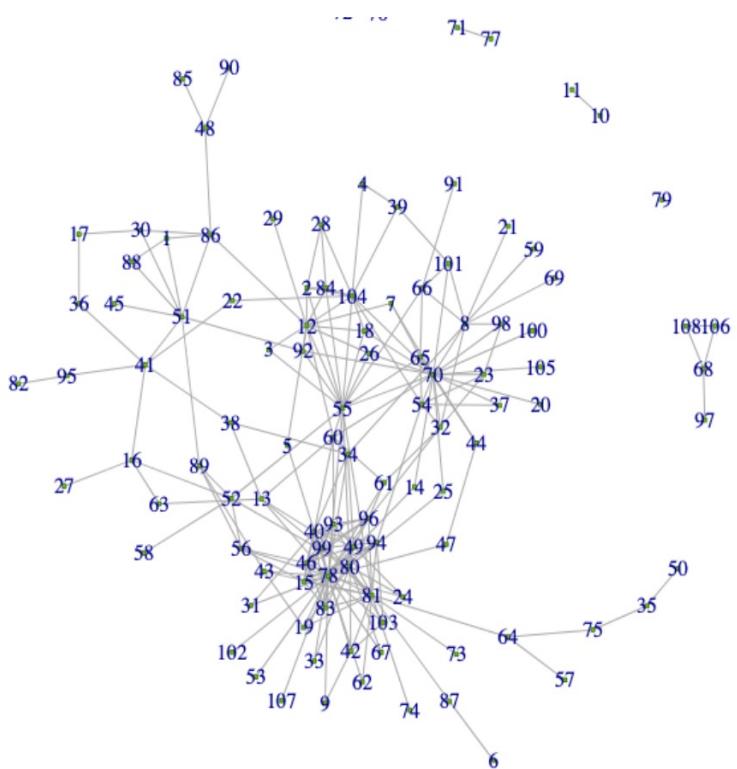


- <https://beta.observablehq.com/@mbostock/d3-force-directed-graph>
- <http://vax.herokuapp.com/game>

Social Network

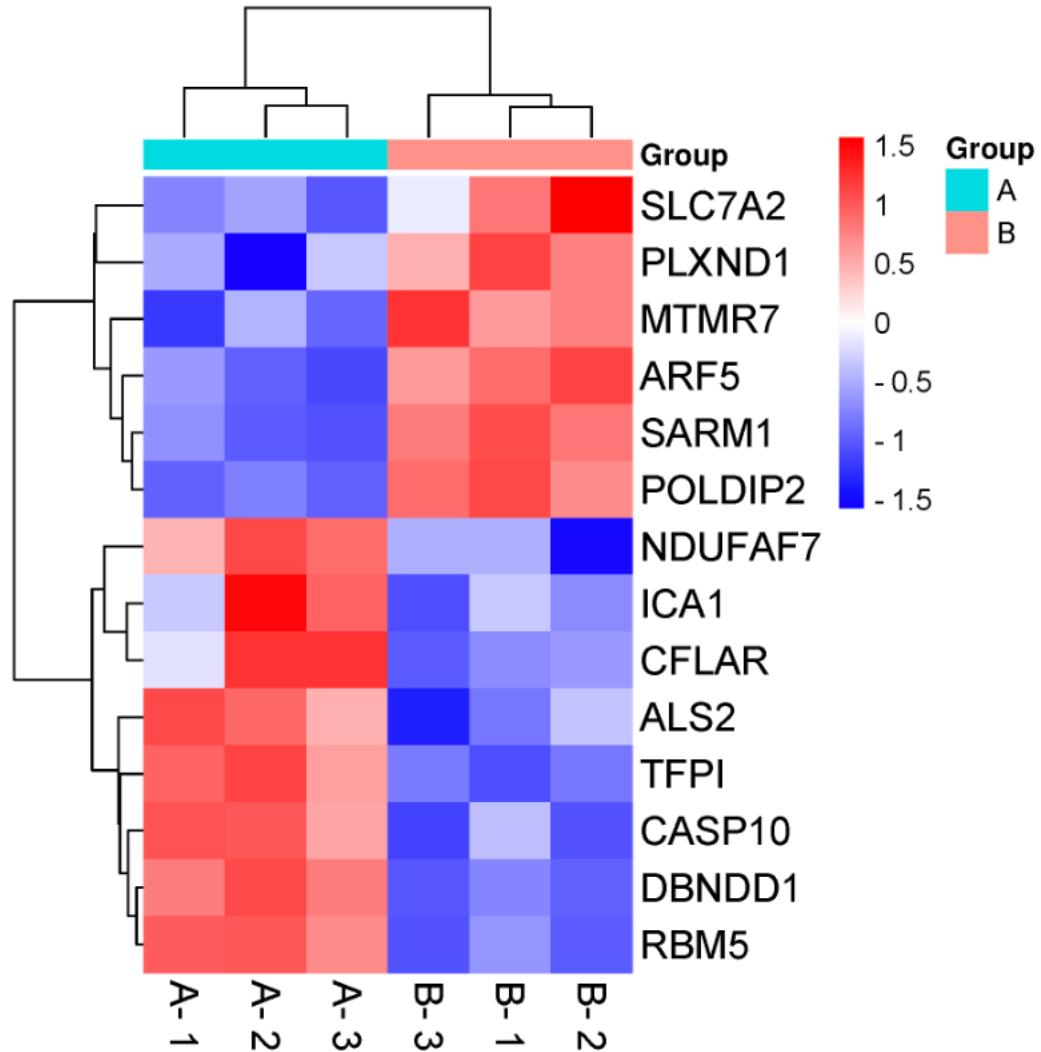


Citation Network



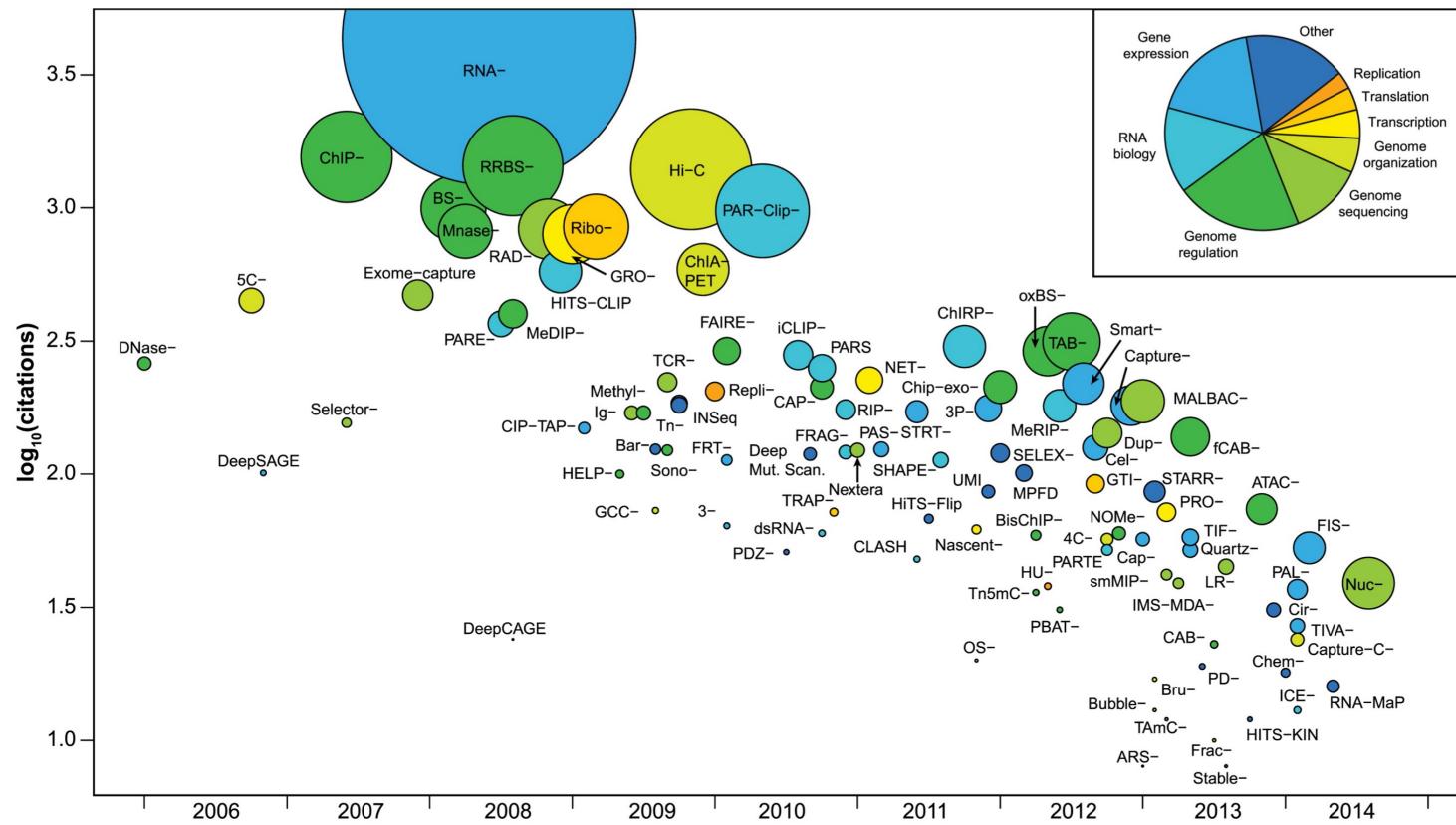
<https://www.sciencedirect.com/science/article/pii/S2666389922001295>

Heat map



https://www.bioinformatics.com.cn/plot_basic_cluster_heatmap_plot_024_en

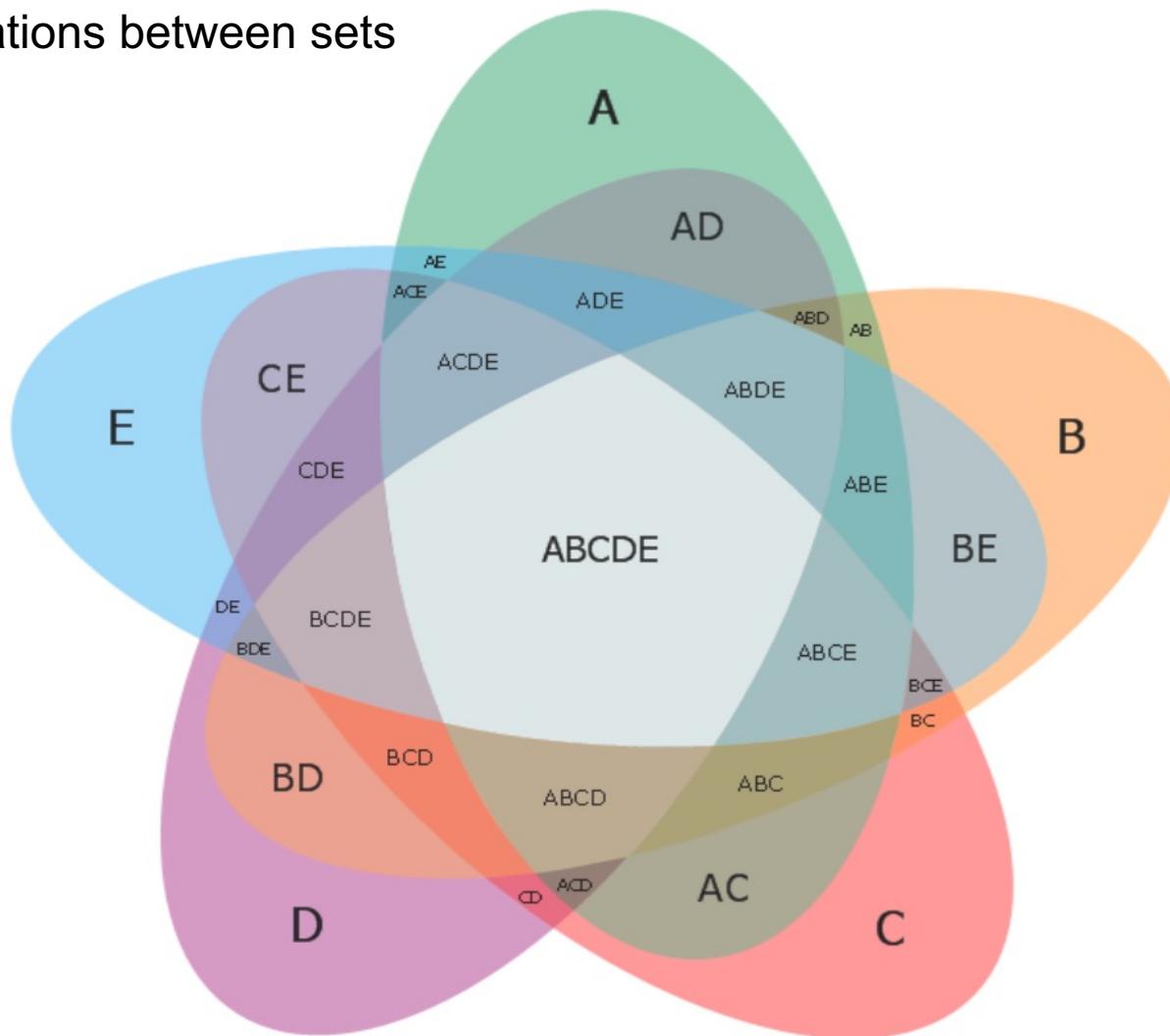
Bubble plots



[https://www.cell.com/molecular-cell/pdf/S1097-2765\(15\)00340-8.pdf](https://www.cell.com/molecular-cell/pdf/S1097-2765(15)00340-8.pdf)

Venn diagram

- Presents relations between sets



Question

- Have you learnt Machine Learning Courses Before?
- Are you familiar with Python?

Getting prepared

- Install Tableau academic version
 - <https://www.tableau.com/academic>
 - Optional: using the practice datasets to generate some plots, e.g., bar, line and pie plots.
- Setup Python programming environment.
 - Preparing for the tutorial in week 2
- Finish **QIZE 1** before Wednesday 11pm.