

GE 461

Introduction to Data Science

Spring 2024



Course Website

All course related material will be provided in the course website

<http://www.cs.bilkent.edu.tr/~ge461/2024Spring>

Check regularly for announcements!

Weekly topics, instructors are stated.

Slides will be provided here.

Assignments released on Moodle.

Various external links to other similar courses and online textbooks.

Instructors

Cross-department Course with Multiple Instructors.

CS Department

S. Aksoy,
C. Alkan,
S. Arashloo,
F. Can,
A.E. Cicek,
H. Dibeklioglu,
A. Dundar,
I. Korpeoglu,
E. Tuzun

EE Department

- T. Cukur,
- C. Tekin

IE Department

- S. Dayanik

TAs will be announced on the Course Website. They will be from all 3 departments.

Location & Time

When: Mon 13:30 – 15:20 and Thursday 8:30 – 10:20.

Where: EE-317.

What: A lot! Introduction to data science fundamentals, techniques and applications; data collection, preparation, storage and querying; parametric models for data; models and methods for fitting, analysis, evaluation, and validation; dimensionality reduction, visualization; various learning methods, classifiers, clustering, data and text mining; applications in diverse domains such as business, medicine, social networks, computer vision; breadth knowledge on topics and hands-on experience through projects and computer assignments.

[See weekly coverage.](#)

Grading Policy

Final: 40%

Project: 60%

Multiple computer/programming/exercise assignments of various sizes.

A project can be assigned earlier than the indicated date on the weekly plan.

Projects can be individual or group based (Python, Java, R or Matlab).

Projects will be uploaded to Moodle.

Grades will be announced on SAPS.

Attendance:

A student who misses more than **9 hours** will fail the course.

What is Data Science?

The field of study that uses various **methods** to extract useful insights and knowledge from the **data** to make data-driven decisions.

Methods can include/require, domain expertise, programming skills (i.e., scripting to process data), statistical modeling (i.e., machine learning algorithms), visualization techniques.

Usually performed on **big** data.



DATA

Data Scientist: The Sexiest Job of the 21st Century

by [Thomas H. Davenport](#) and [D.J. Patil](#)

From the October 2012 Issue

Recommended readings:

[http://cdn.oreilly.com/radar/2010/06/What is Data Science.pdf](http://cdn.oreilly.com/radar/2010/06/What_is_Data_Science.pdf)

<https://hbr.org/2012/10/data-scientist-the-sexiest-job-of-the-21st-century>

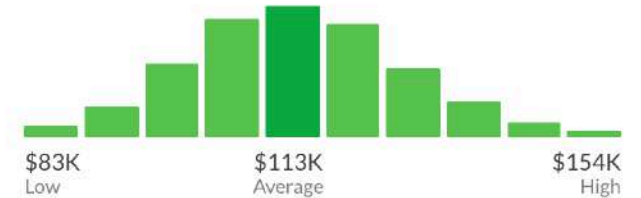
Data Scientist Salaries

6,606 Salaries Updated Jan 22, 2020

Industries Company Sizes Years of Experience

Average Base Pay

\$113,309 /yr



Additional Cash Compensation [?]

Average \$11,258

Range \$3,850 - \$26,084

How much does a Data Scientist make?
The national average salary for a Data Scientist is \$113,309 in United States. Filter by location to see... [More](#)

VS

Computer Engineer Salaries

256,924 Salaries Updated Jan 22, 2020

Industries Company Sizes Years of Experience

Average Base Pay

\$92,046 /yr



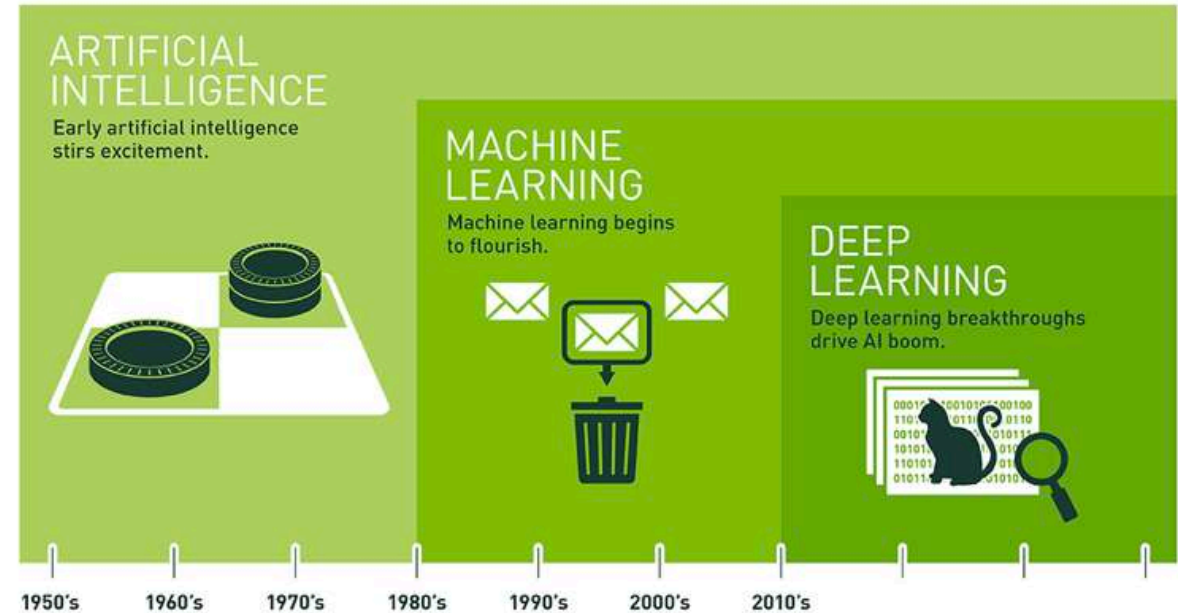
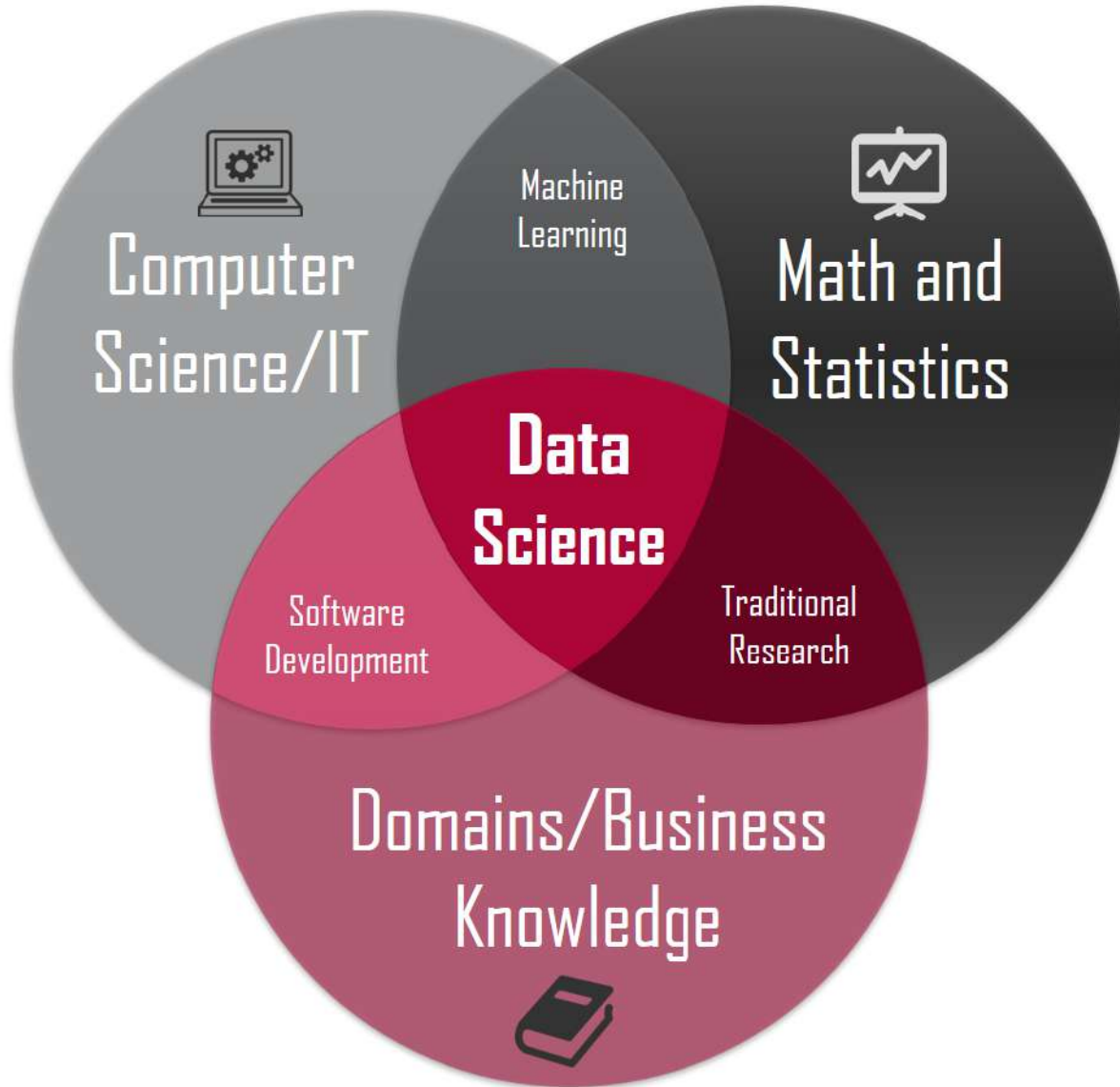
Additional Cash Compensation [?]

Average \$7,871

Range \$1,810 - \$20,486

How much does a Computer Engineer make?
The national average salary for a Computer Engineer is \$92,046 in United States. Filter by location to see... [More](#)

What is NOT Data Science?



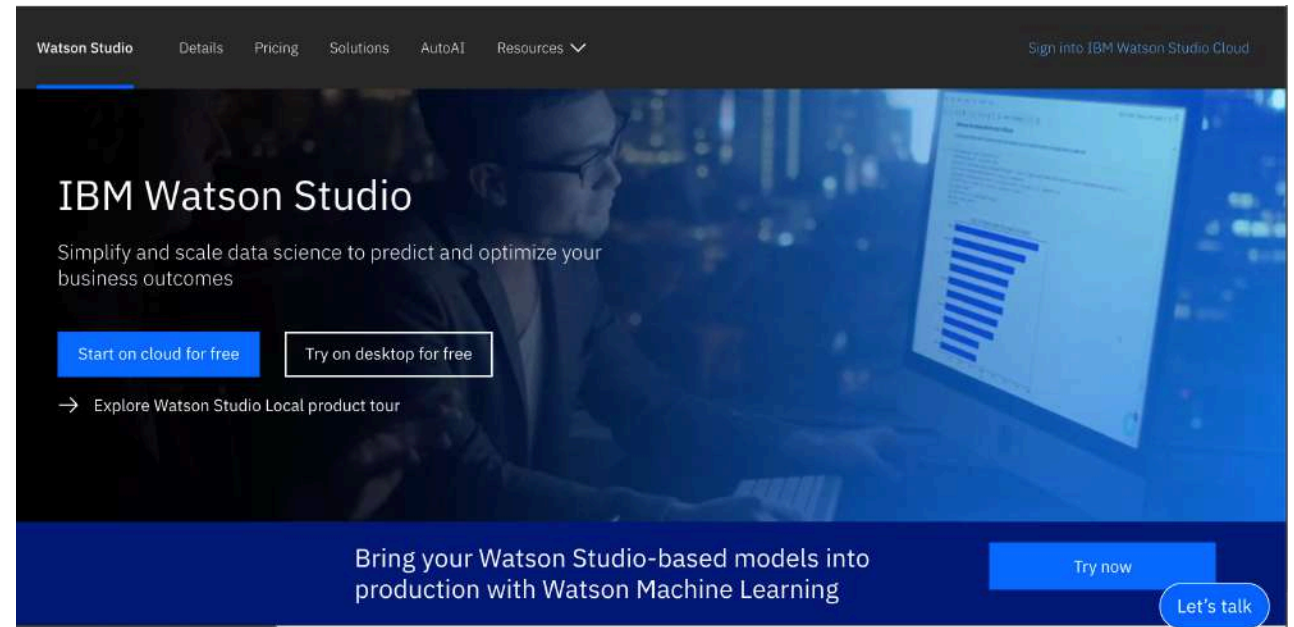
Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

Data Science makes use of AI, ML, DL

<https://blogs.nvidia.com/blog/2016/07/29/whatsdifference-artificial-intelligence-machine-learningdeep-learning-ai/>

What is NOT Data Science? Example

An AI breakthrough in 2011, now empowers Data Science.



Data Science vs Other Related Terms

Many terms are used interchangeably; vague definitions.

Data Science aims at finding the right questions, more predictive analysis. Somewhat involves creativity. On the other hand, **Business Intelligence** aims helping in the decision making of a business based on past data.

Data mining is a technique that searches for patterns in the data and can be considered as a tool of Data Science.

For example: Baby diapers and beer are frequently bought together.

Data analytics aims at analyzing data to find answers to concrete questions.

For instance, optimizing the teller processes at the bank to serve more customers.

It is a tool for **Business Intelligence**.

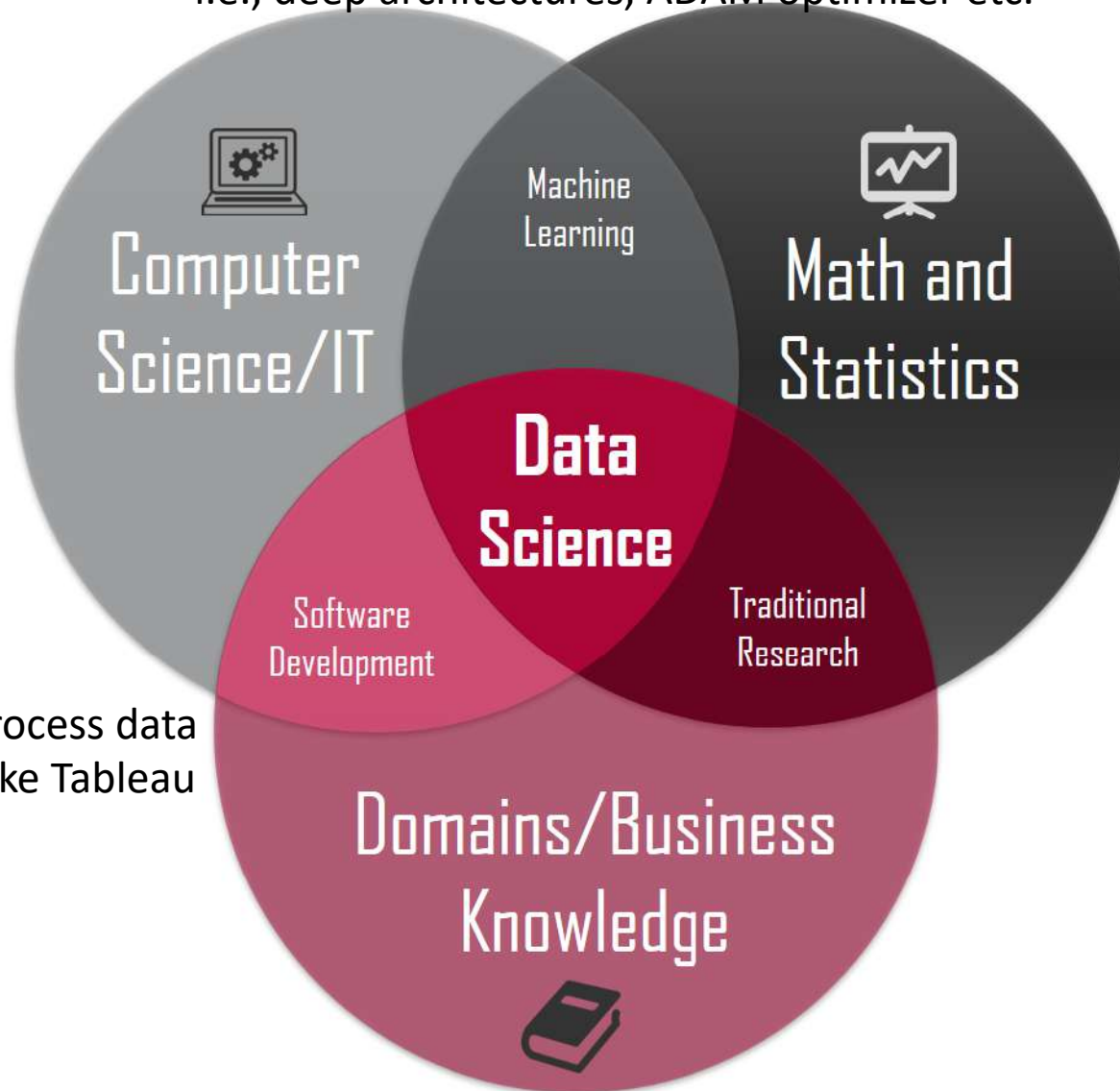
Why Now? Some advances

Better machine learning algorithms
i.e., deep architectures, ADAM optimizer etc.

Faster Computers

GPU power to crunch large datasets

Better ways (NoSQL) to manage
Data (Hadoop, Hive, HBase)



+ big data

Data is ubiquitous
Cheap to produce and
store

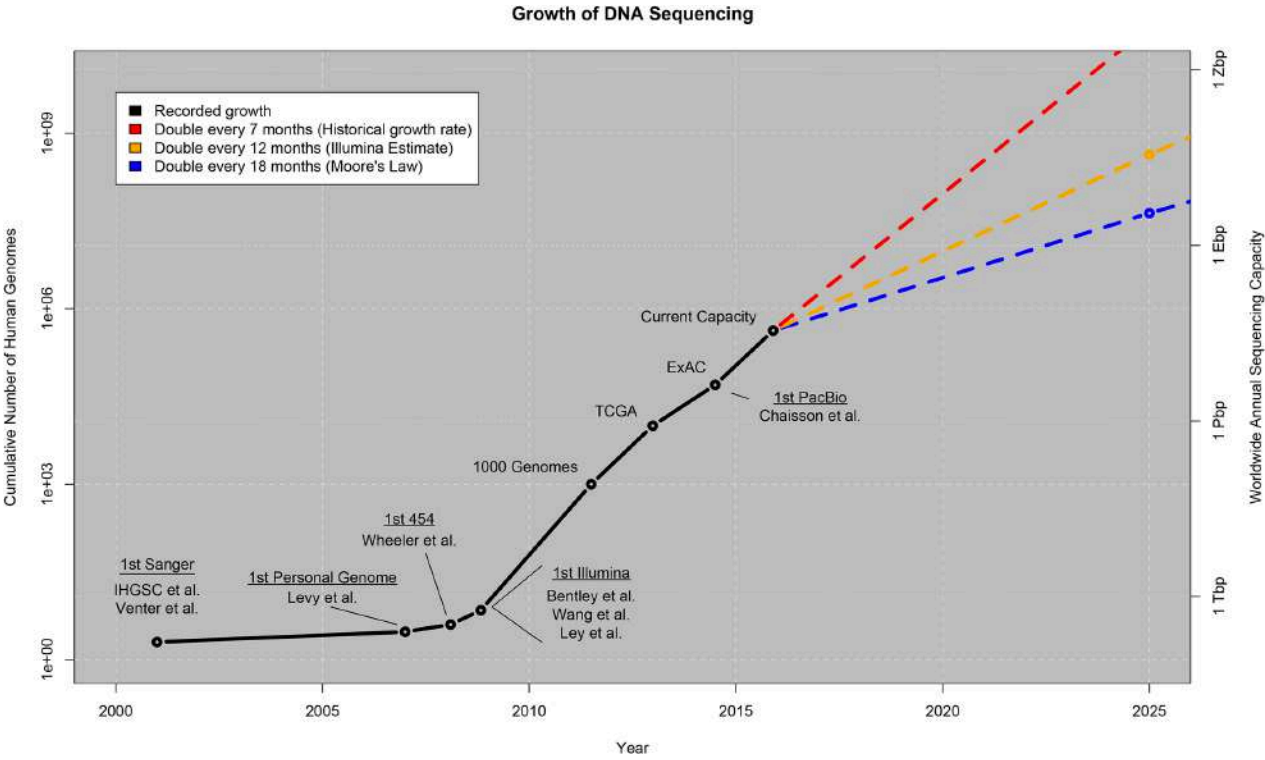
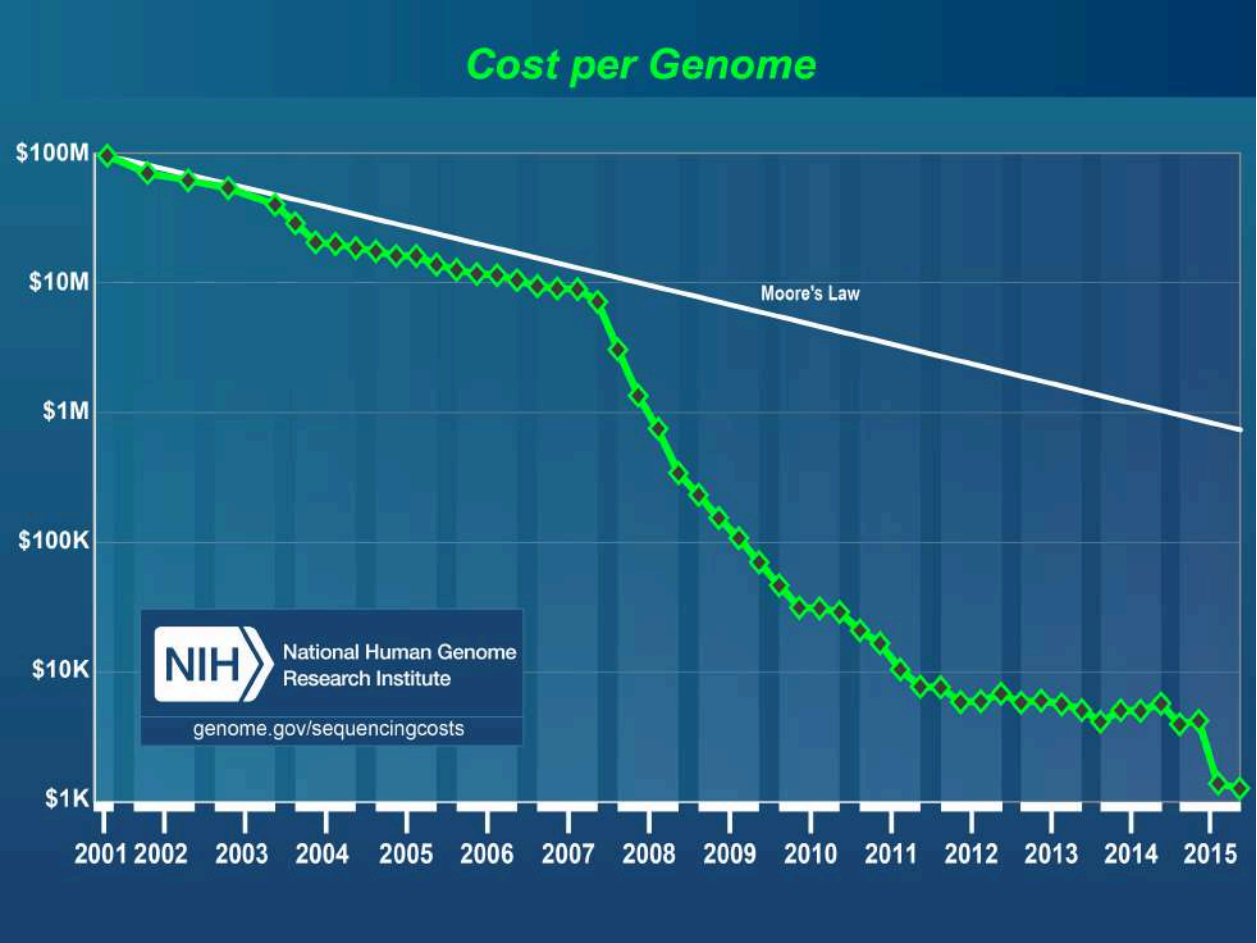


amazon
S3

Python and R vs SAS and SPSS to process data
Advanced data visualization tools like Tableau

Big Data

Data is easy to produce, cheap to store. One example from genomics.

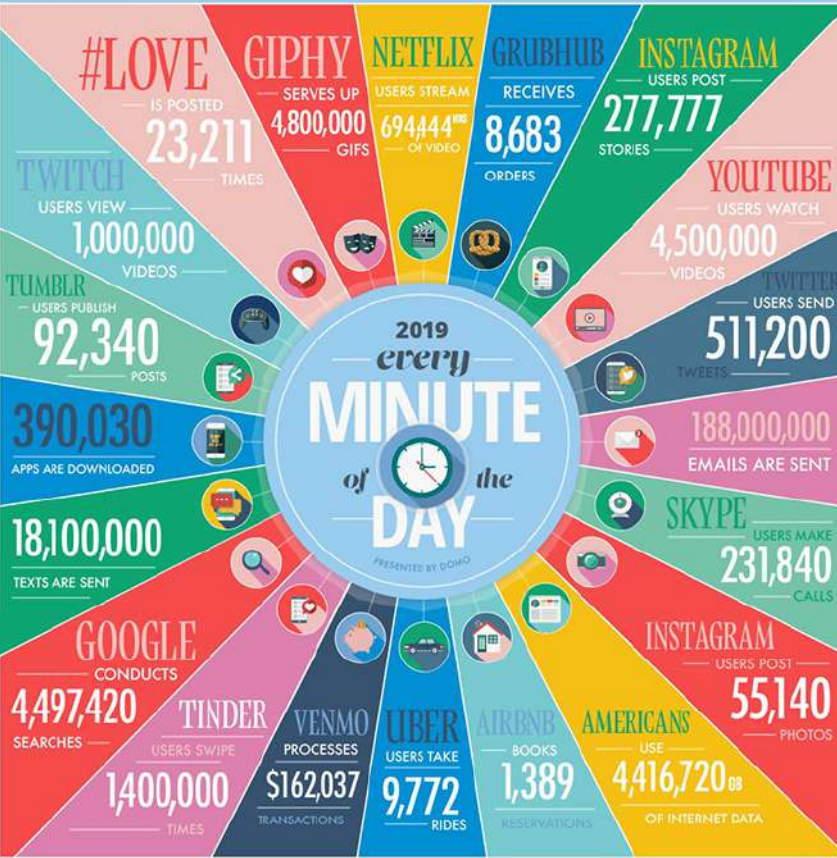




DATA NEVER SLEEPS 7.0

How much data is generated every minute?

There's no way around it: big data just keeps getting bigger. The numbers are staggering, and they're not slowing down. By 2020, there will be 40x more bytes of data than there are stars in the observable universe. In our 7th edition of Data Never Sleeps, we bring you the latest stats on how much data is being created in every digital minute — and the numbers are staggering.



The world's internet population is growing significantly year-over-year. As of January 2019, the internet reaches 36% of the world's population and now represents 4.3 billion people — a 19% increase from January 2018.



GLOBAL INTERNET POPULATION GROWTH 2012-2018 (IN BILLIONS)

SOURCES: STATISTA, INTERNET LIVE STATS, EXPANDED BANKINGS, NATIONAL ASSOCIATION OF CITY TRANSPORTATION OFFICIALS, WIREX

The ability to make data-driven decisions is crucial to any business. With each click, swipe, share, and like, a world of valuable information is created. Domo puts the power to make those decisions right into the palm of your hand by connecting your data and your people at any moment on any device, so they can make the kind of decisions that make an impact.

Learn more at domo.com



Data Never Sleeps 9.0

How much data is generated every minute?

The 2020 pandemic upended everything, from how we engage with each other to how we engage with brands and the digital world. At the same time, it transformed how we eat, how we work and how we entertain ourselves. Data never sleeps and it shows no signs of slowing down. In our 9th edition of the "Data Never Sleeps" infographic, we bring you a glimpse of how much data is created every digital minute in our increasingly data-driven world.



As of July 2021, the internet reaches 65% of the world's population and now represents 5.17 billion people—a 10% increase from January 2021. Of this total, 92.6 percent accessed the internet via mobile devices. According to Statista, the total amount of data consumed globally in 2021 was 79 zettabytes, an annual number projected to grow to over 180 zettabytes by 2025.

Global Internet Population Growth (IN BILLIONS)



As the world changes, businesses need to change too—and that requires data. Domo gives you the power to make data-driven decisions at any moment, on any device, so that you can make smart choices in a rapidly changing world. Every click, swipe, share, or like tells you something about your customers and what they want, and Domo is here to help you and your business make sense of all of it.

Learn more at domo.com

SOURCES: LOCAL IQ, BUSINESS OF APPS, DEEPEN SOURCE, THOUGHTS EXPANDED BANKING, WE FRANK WORLD, STATISTA, CNBC, DATAWATTS, VEE, THE SHARE FILE, PACTIVE, RIVALTA, THE BRIDGE, MANAGEMENT CONSULTING GROUP, A CASE ADVANTAGE, FINANCIAL, INTERNET LIVE STATS, SOCIAL, STATISTA



Data Never Sleeps 11.0

Domo has been keeping tabs on the world's data usage—in a minute—for over a decade now. What the numbers consistently show is that how we use data is always evolving—and that data isn't slowing down. We're also seeing some big changes. The rise of Artificial Intelligence (AI) is reshaping the way we communicate, work, and create. Digital payments continue to replace traditional transactions. Taylor Swift streams in countless headphones. And a rash of cybercrime grows alongside these digital experiences.

In Domo's 11th edition of Data Never Sleeps, we take the pulse of our digital age, where every click, swipe, and stream fuels an ever-expanding digital universe. These are not just numbers; they are the heartbeat of a world where data reigns supreme.



The world's internet population continues to grow significantly year-over-year. As of November 2023, the internet represents 5.2 billion people—approximately 64.6% of the global population. According to Statista, the total amount of data predicted to be created, captured, copied, and consumed globally in 2023 is 120 zettabytes, a number projected to grow to 181 zettabytes by 2025.

Global Internet Population Growth (IN BILLIONS)



As data grows and evolves, businesses need to grow and evolve, too. Domo helps you harness the power of data so you can change as quickly as the world changes and make data-driven decisions that set you apart from the crowd. Let Domo help you make sense of all the clicks, swipes, and shares so you can see the big picture that a lot of small decisions make.

Learn more at domo.com

SOURCES: EARTRIVE, DUSTIN STOUT, DEMANDSAGE, HOOTSUITE, BUSINESSOFAPPS, DOORDASH, SOCIALPLOT, X, TWITTER.COM, GITNIX, INVIGATE, THINKIMPACT, SIMA.ORG, STATISTA, PR NEWSWIRE, NETSCOUT



Database (old) vs Data Science (new)

	Databases	Data Science
Data Value	“Precious”	“Cheap”
Data Volume	Modest	Massive
Examples	Bank records, Personnel records, Census, Medical records	Online clicks, GPS logs, Tweets, Building sensor readings
Priorities	Consistency, Error recovery, Auditability	Speed, Availability, Query richness
Structured	Strongly (Schema)	Weakly or none (Text)
Properties	Transactions, ACID*	CAP* theorem (2/3), eventual consistency
Realizations	SQL	NoSQL: Riak, Memcached, Apache River, MongoDB, CouchDB, Hbase, Cassandra,...

Modelling vs Data-Driven Solutions

Scientific modelling

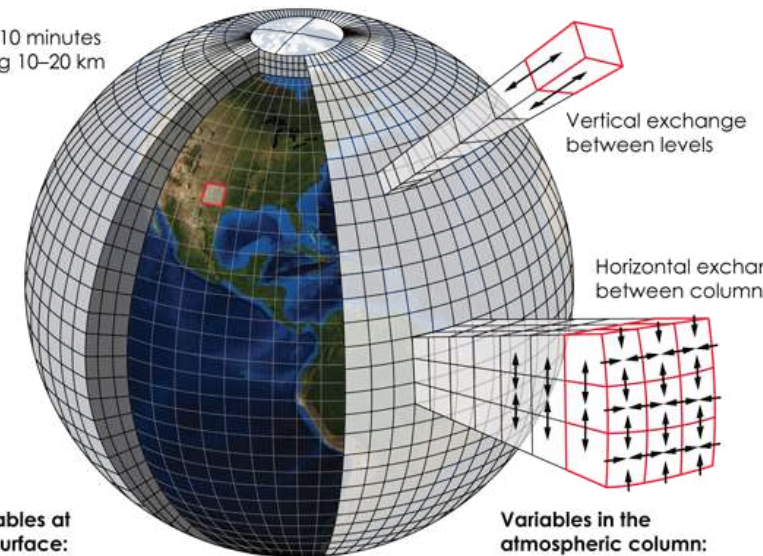
Background knowledge, set of rules, principles, representations etc.
Example: Weather forecasting.

Data-Driven Solutions

No or little apriori model, which is replaced by an inference algorithm (e.g., Neural Network, SVM etc.).
Example: Image classification.

Weather forecast modeling

Timestep 5-10 minutes
Grid spacing 10-20 km

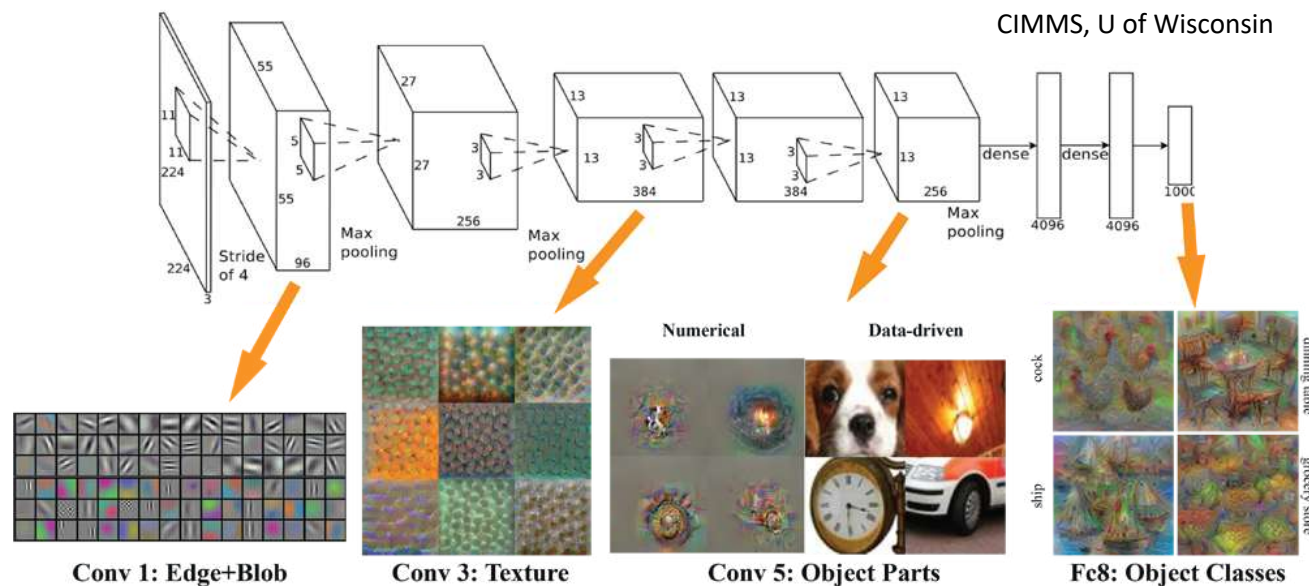


Variables at the surface:

- Temperature
- Humidity
- Pressure
- Moisture fluxes
- Heat fluxes
- Radiation fluxes

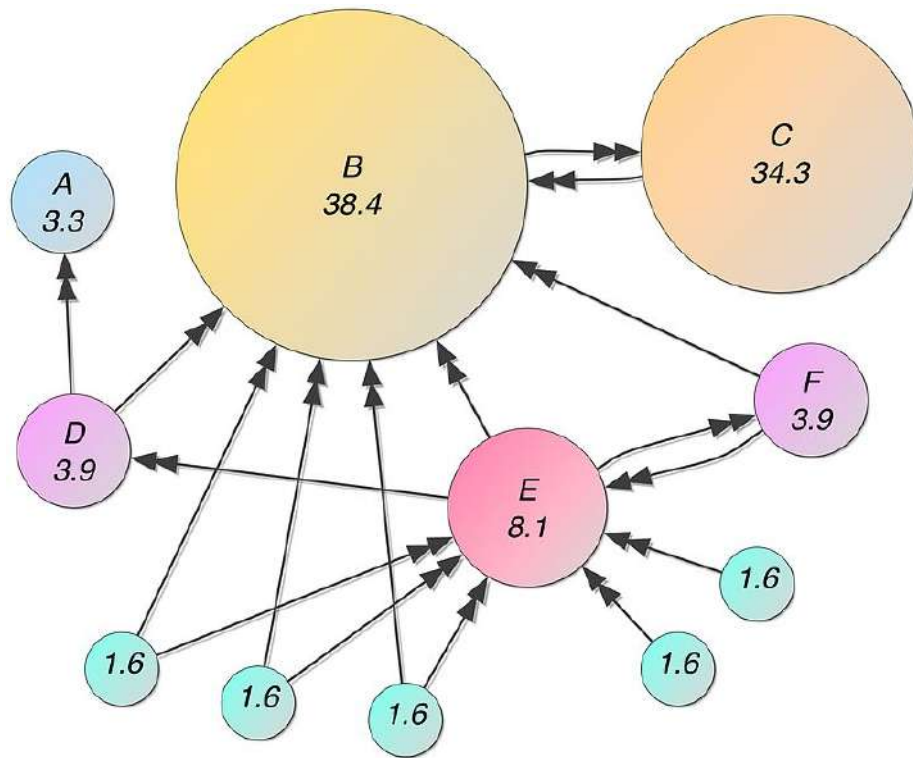
Variables in the atmospheric column:

- Wind vectors
- Humidity
- Clouds
- Temperature
- Height
- Precipitation



Some examples - Search

Google PageRank Algorithm



[PDF] [The PageRank Citation Ranking - Stanford InfoLab ...](#)

[ilpubs.stanford.edu](#) > ... ▼

by L Page - 1999 - Cited by 12987 - Related articles

Original Paper



Cornell University

[arXiv.org](#) > [cs](#) > [arXiv:1503.01331](#)

[Computer Science](#) > [Social and Information Networks](#)

PageRank Approach to Ranking National Football Teams

[Verica Lazova](#), [Lasko Basnarkov](#)

(Submitted on 4 Mar 2015 (v1), last revised 21 Apr 2015 (this version, v2))

*Used in many applications to
have data driven answers to various problems*

Some examples – Recommendation Systems

amazon.com

Recommended for You

Amazon.com has new recommendations for you based on items you purchased or told us you own.



[The Little Big Things: 163 Ways to Pursue Excellence](#)



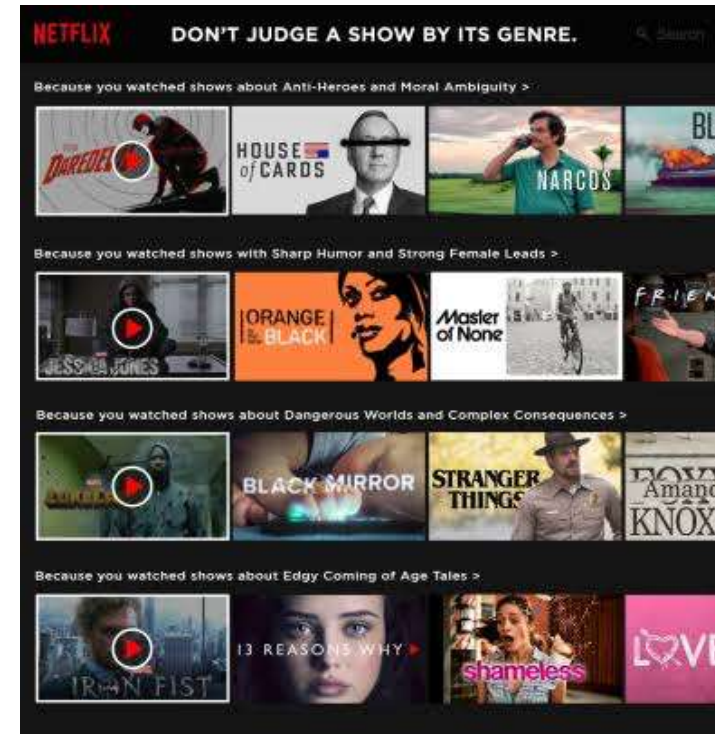
[Fascinate: Your 7 Triggers to Persuasion and Captivation](#)



[Sherlock Holmes \[Blu-ray\]](#)



[Alice in Wonderland \[Blu-ray\]](#)



	← users →					
↑ movies ↓	1		?	3	5	?
	?	1				2
		4		4	5	?

Some examples – Flu Trends

Google Flu Trends

nature

Letter | Published: 19 February 2009

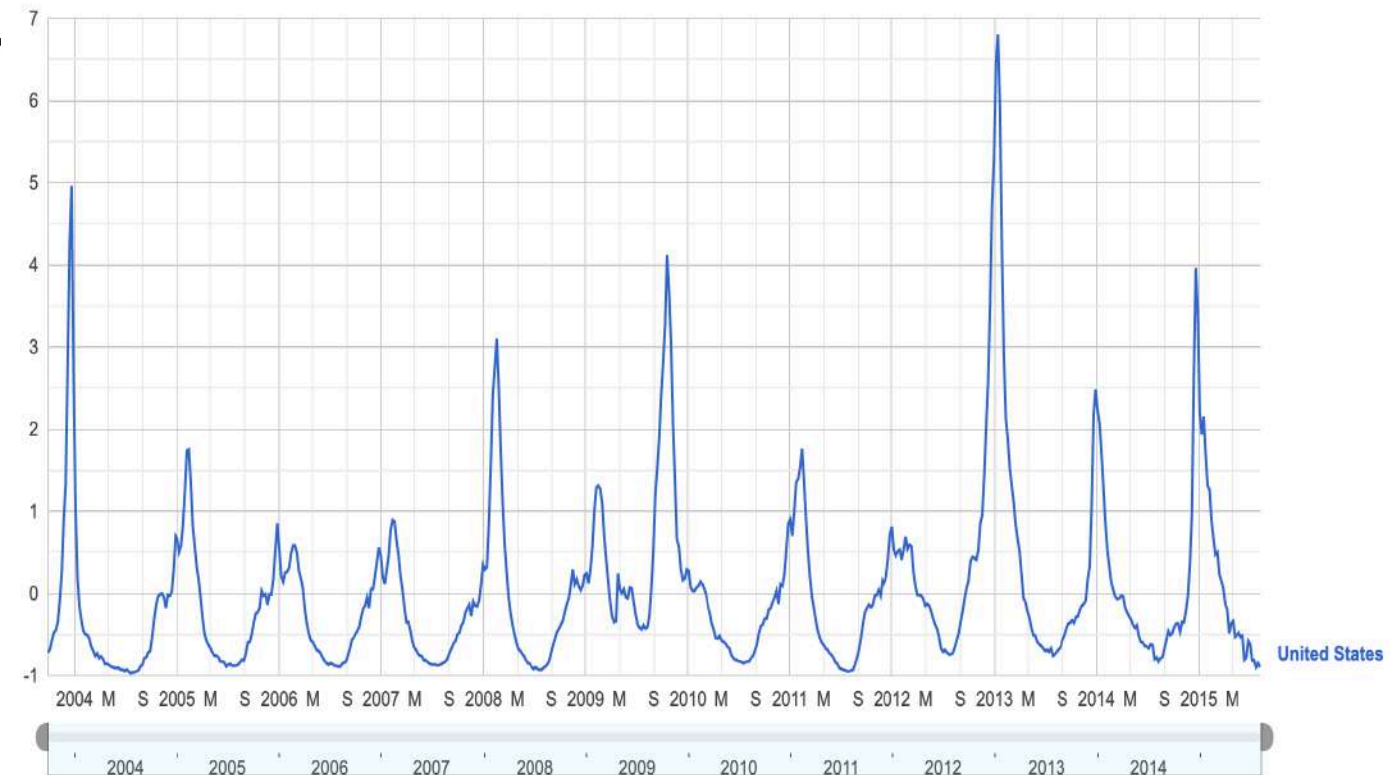
Detecting influenza epidemics using search engine query data

Jeremy Ginsberg, Matthew H. Mohebbi , Rajan S. Patel, Lynnette Brammer, Mark S. Smolinski & Larry Brilliant

Nature **457**, 1012–1014(2009) | [Cite this article](#)

5195 Accesses | **1876** Citations | **474** Altmetric | [Metrics](#)

Flu search activity (standard deviation from baseline) ?



Data from Google Inc. Last updated: Aug 19, 2015

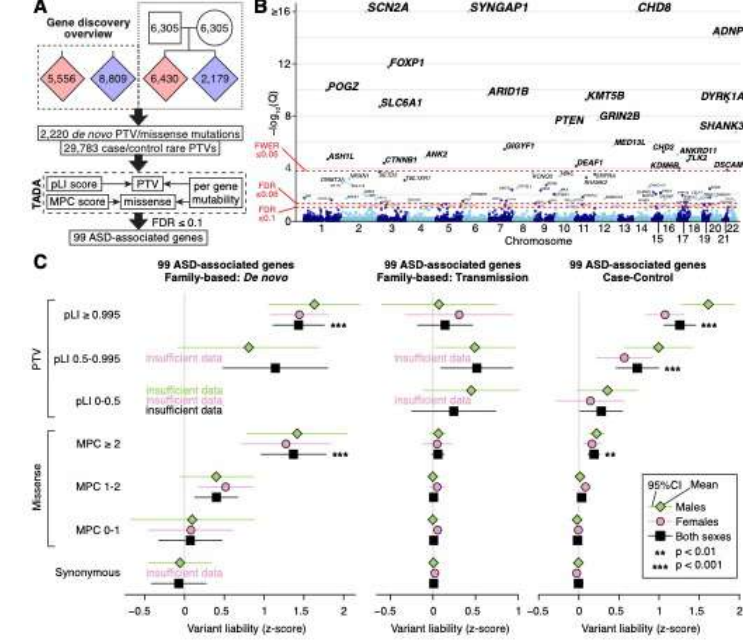
Some examples – Comp. Biology

Data Science for Gene Risk Prediction

It is not enough to collect the data.

What does the data tell us?

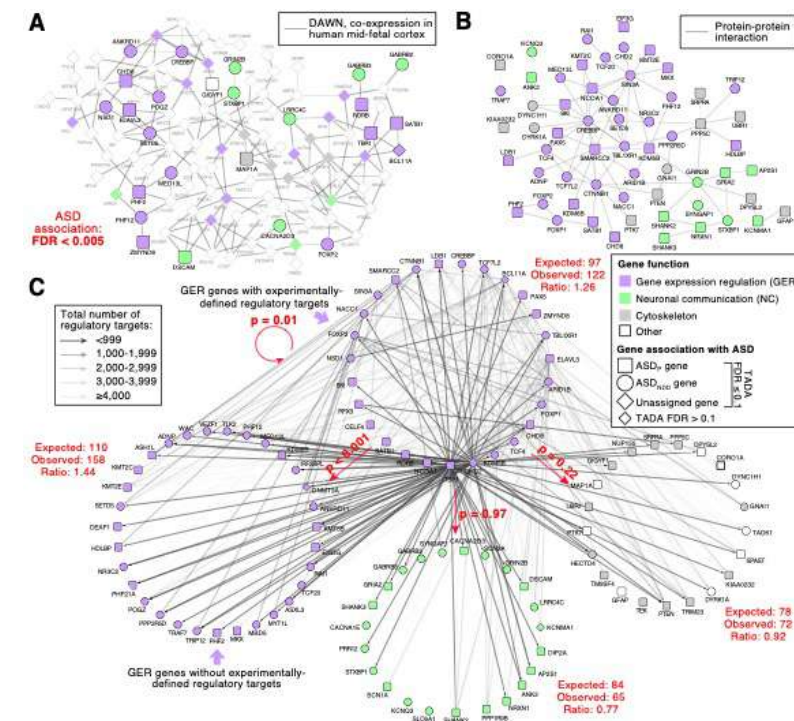
Use methods to analyze the it.



HEALTH • AUTISM

Researchers Find 102 Genes Linked to Autism in One of the Largest Studies of Its Kind to Date

In a **study published Jan. 23 in *Cell***, researchers led by Joseph Buxbaum, director of the Seaver Autism Center for Research and Treatment at Mount Sinai, took advantage of better genetic sequencing technologies and one of the largest databases of DNA samples from people with autism to identify 102 genes associated with autism, including 30 that had never before been connected with the condition. The study also distinguished the genes more closely associated with autism from those that might also contribute to other neurodevelopmental disorders including intellectual and motor disabilities.



Some examples – Comp. Biology

Machine Learning for Gene Risk Prediction

Build algorithms to predict the risk

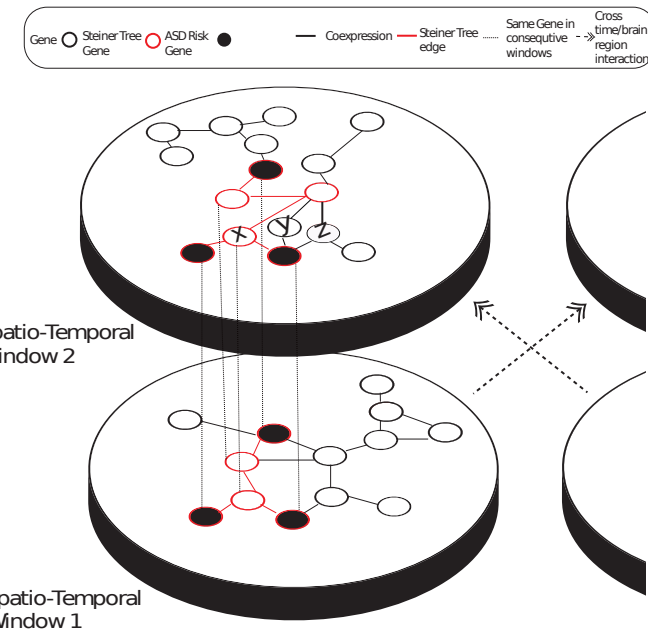


HEALTH • AUTISM

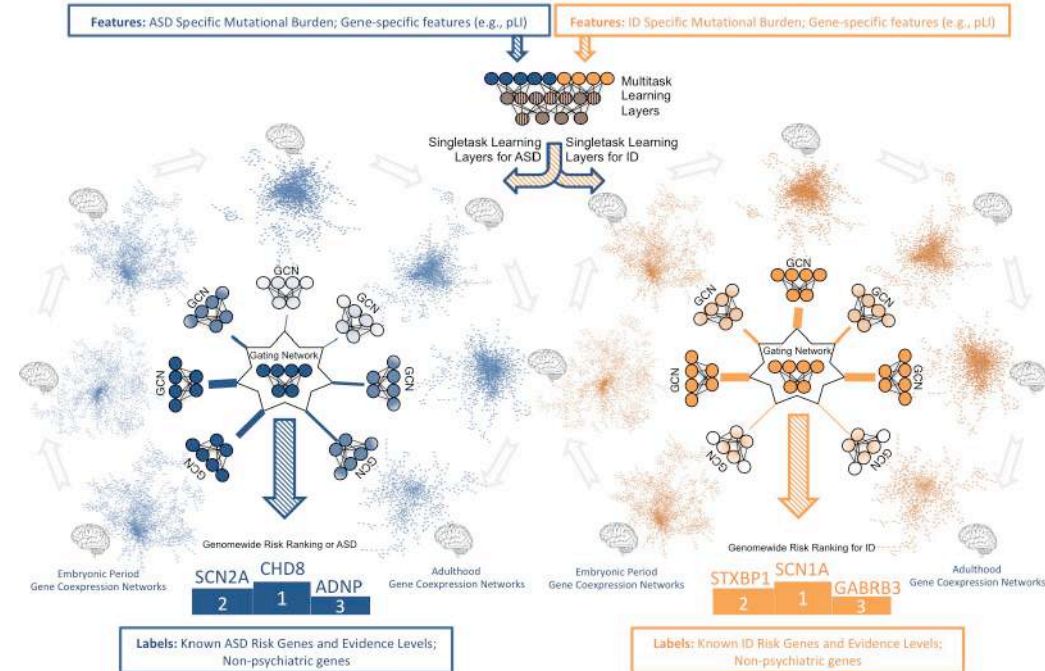
Researchers Find 102 Genes Linked to Autism in One of the Largest Studies of Its Kind to Date

In a **study published Jan. 23 in *Cell***, researchers led by Joseph Buxbaum, director of the Seaver Autism Center for Research and Treatment at Mount Sinai, took advantage of better genetic sequencing technologies and one of the largest databases of DNA samples from people with autism to identify 102 genes associated with autism, including 30 that had never before been connected with the condition. The study also distinguished the genes more closely associated with autism from those that might also contribute to other neurodevelopmental disorders including intellectual and motor disabilities.

Satterstrom *et al.*, CELL 2020



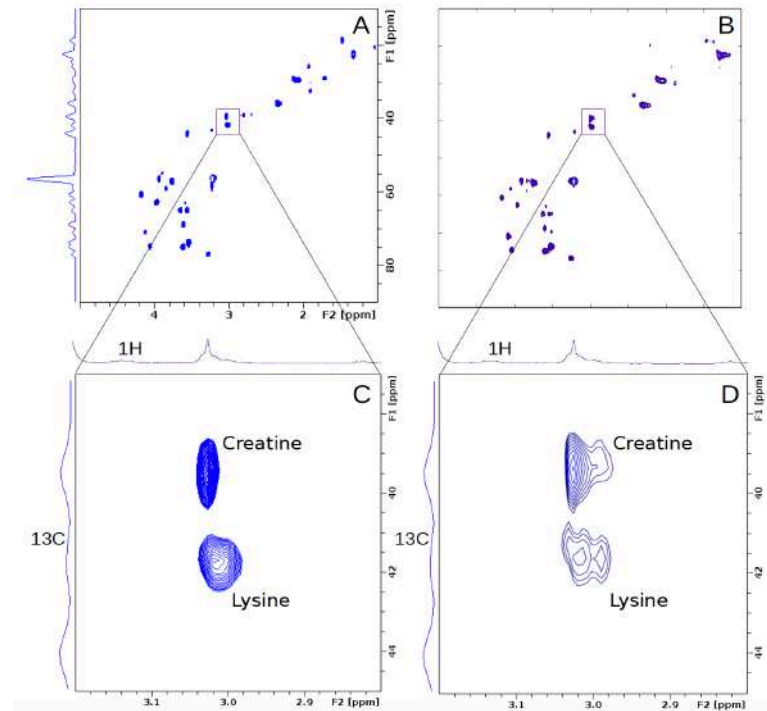
Spatio-temporal Network-based Analysis. Norman and Cicek, Bioinformatics 2019.



Some examples – Comp. Biology

Data Science for Online Feedback to Surgeons

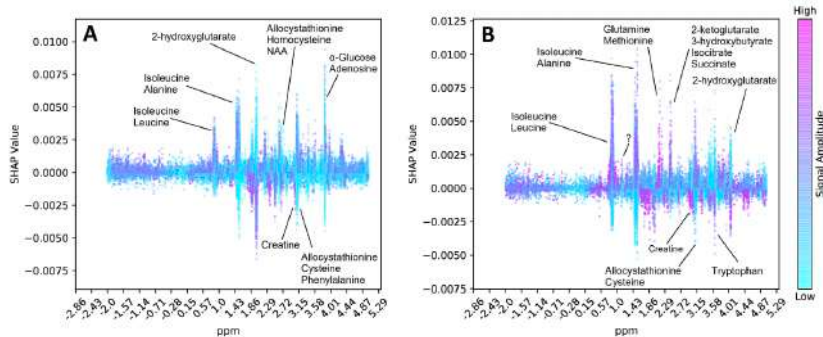
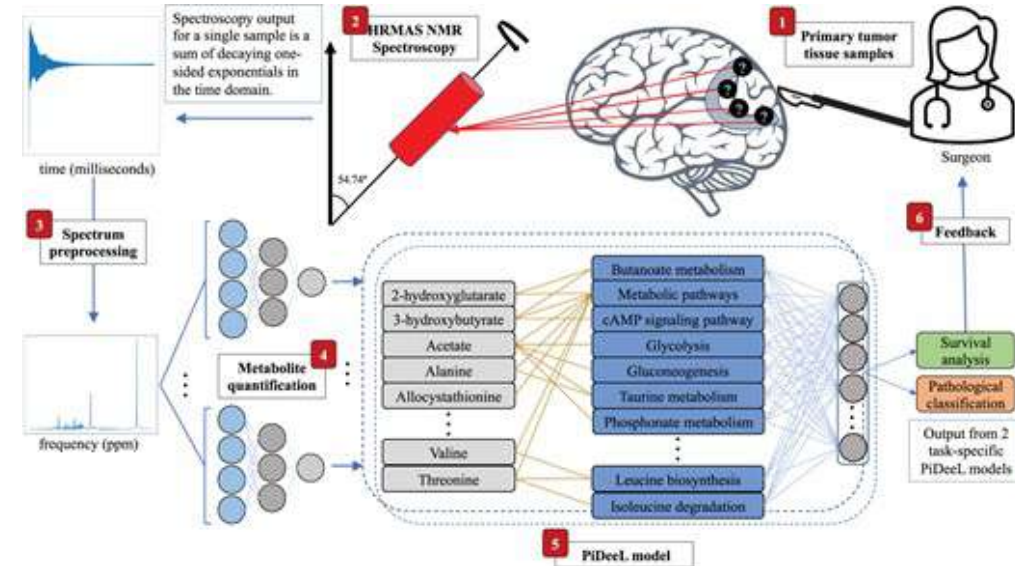
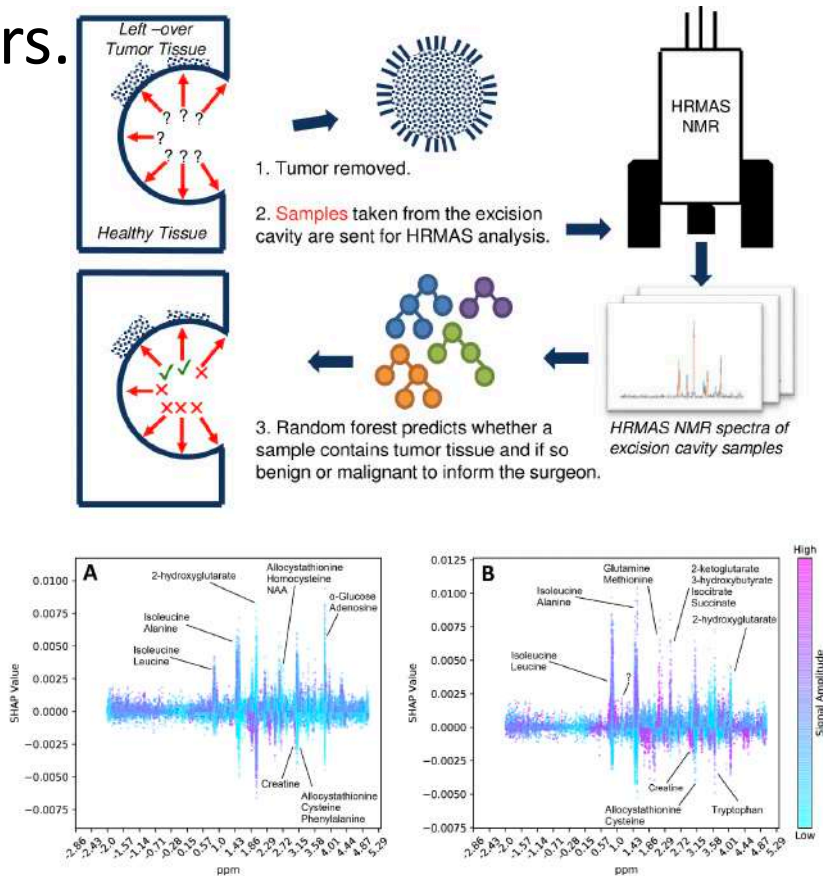
Use Multiple Multivariate Regression to predict the result of a test that is infeasible to perform during surgery due to time requirement.



Some examples – Comp. Biology

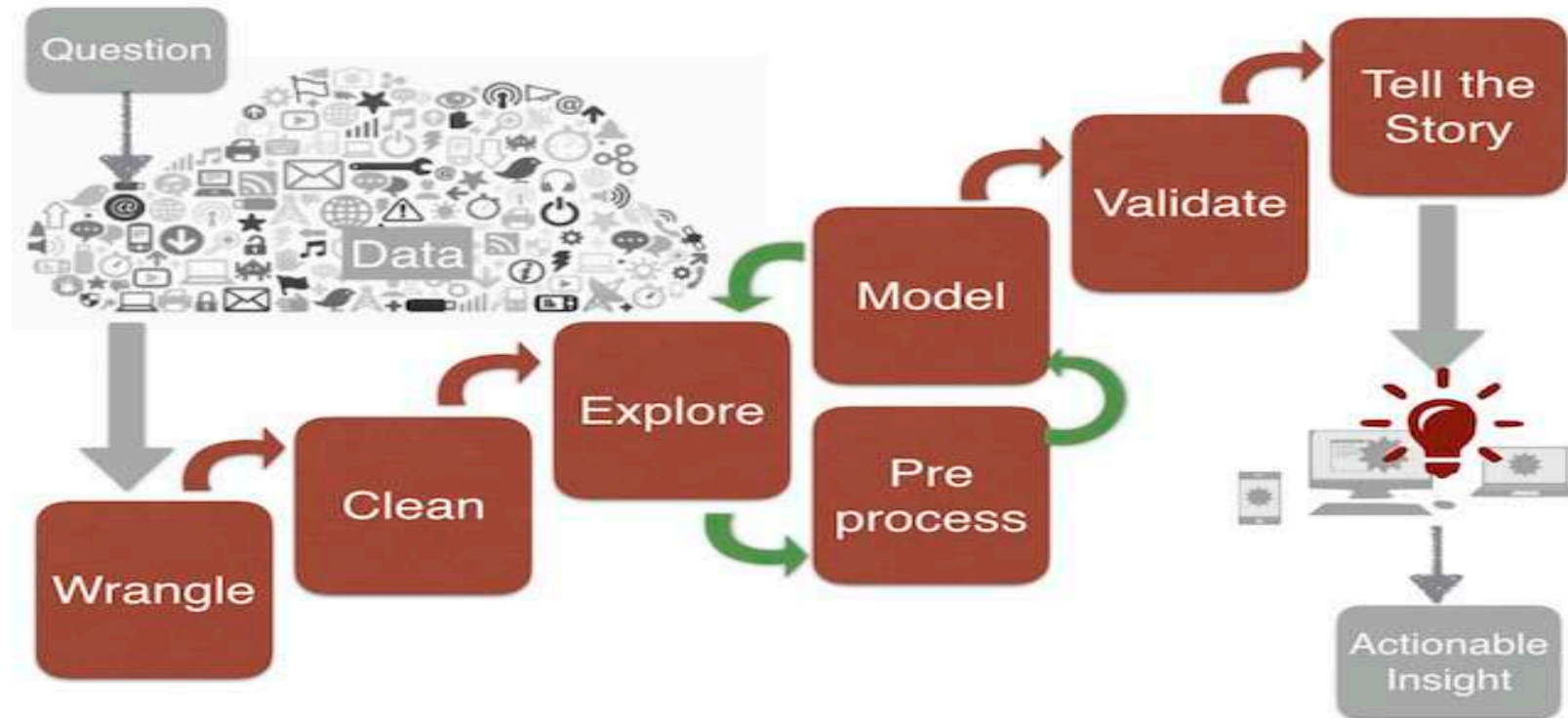
Machine Learning for Online Feedback to Surgeons

Design a neural network that learns important parts of to classify tumors.



Cakmakçı *et al.*, PLoS Computational Biology 2020, 16 (11).
Kaynar *et al.* Bioinformatics 2023, 39 (11), btad684.

Data Science Pipeline



Data Science Pipeline - Data Collection

Many data types, many ways

Sensors

Crowdsourcing, putting humans at work once computers fail:

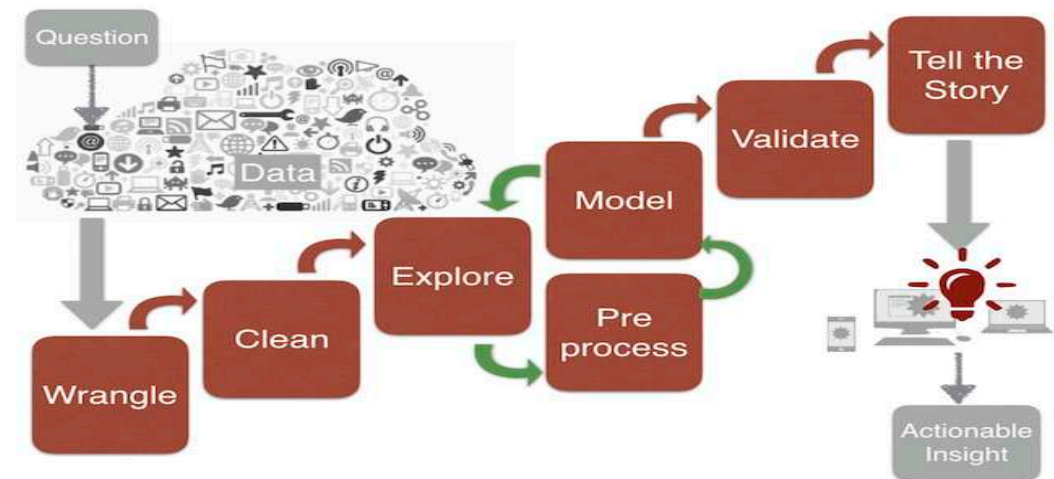
Mechanical Turk

Crawling

Questionnaires..



The Turk

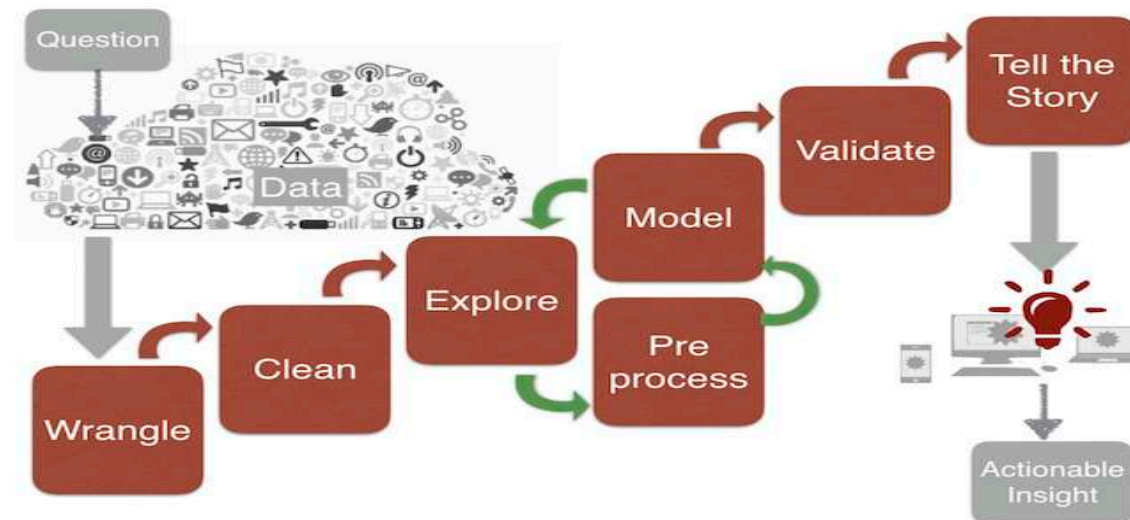


Data Science Pipeline - Data Wrangling

After you obtain the raw data converting it into a more useful format

Gather multiple files into single, standardized format

For example: Unite multiple crawled files into one, get rid of html tags etc.



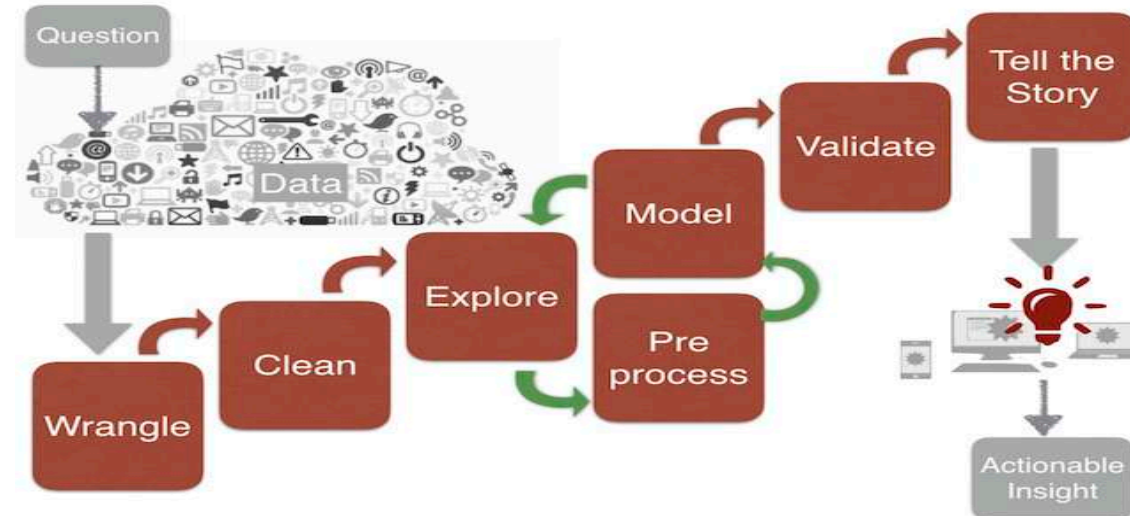
Data Science Pipeline - Data Cleaning

Dig deeper into the data after standardization and detect problems.

Inconsistencies

Outliers

Missing values



Data Science Pipeline

Explore – Preprocess – Model Cycle

1. Explore the structure of the data and decide on the appropriate model to analyze.

For instance: sequence data, maybe LSTM?

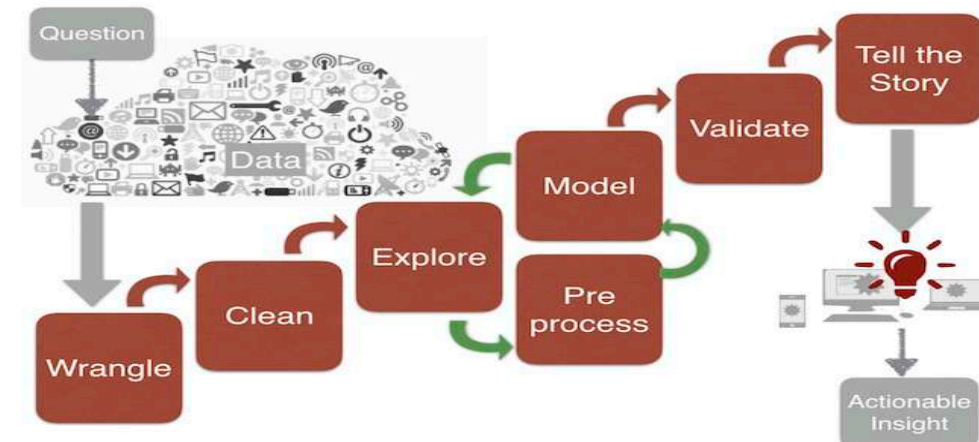
image data, maybe Convolutional Neural Network
transformers for all?

2. Preprocess the data to be fit into the model

For instance, RGB -> Grayscale

3. Apply the model and analyze results

4. Go to 1.



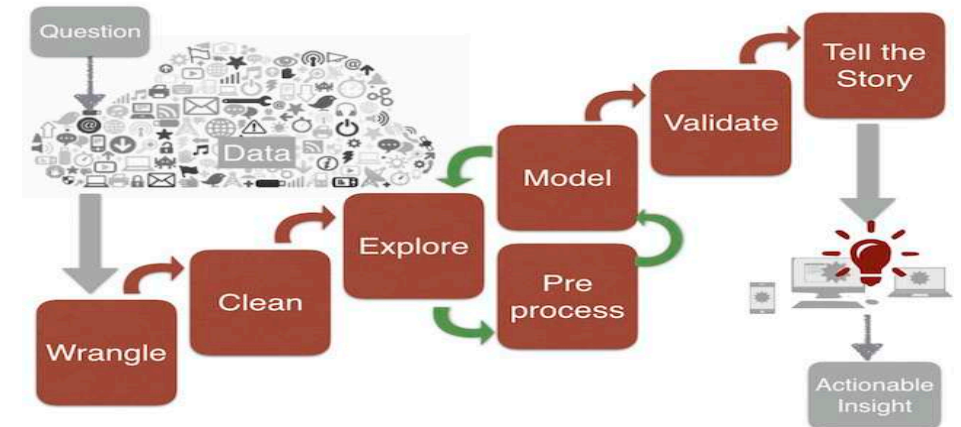
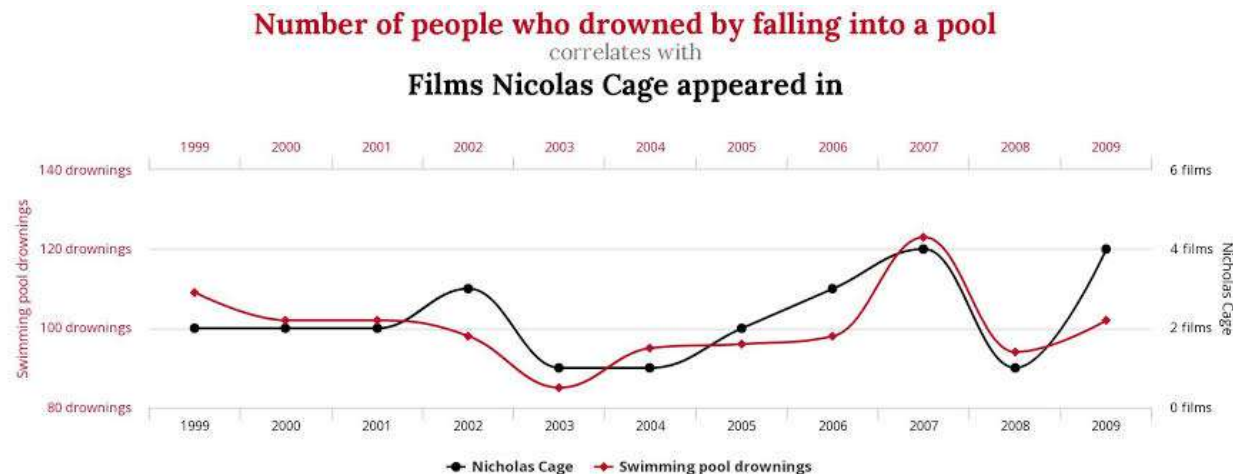
Data Science Pipeline - Validation

After you fine-tuned your model in the previous cycle validate your data on a data that has not been seen by the model.

Validate that your claim is not just random finding.

Multiple hypothesis correction

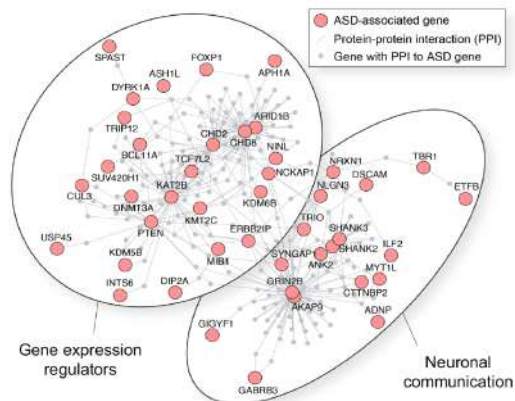
Correlation is not **causation**.



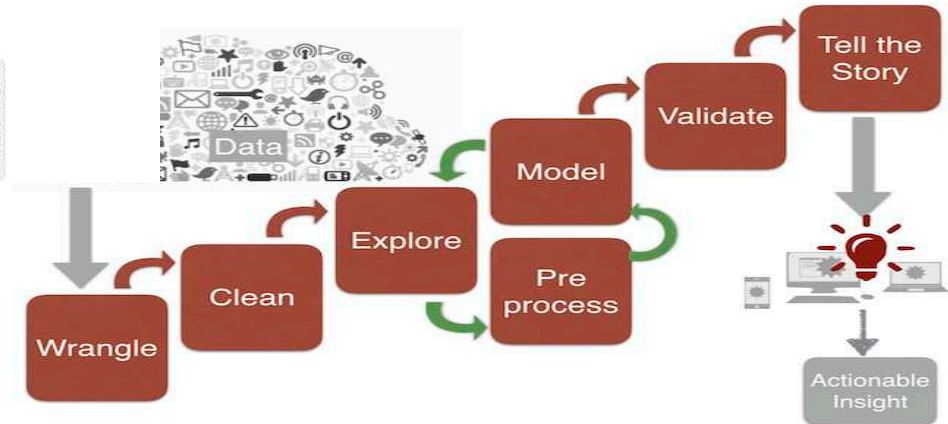
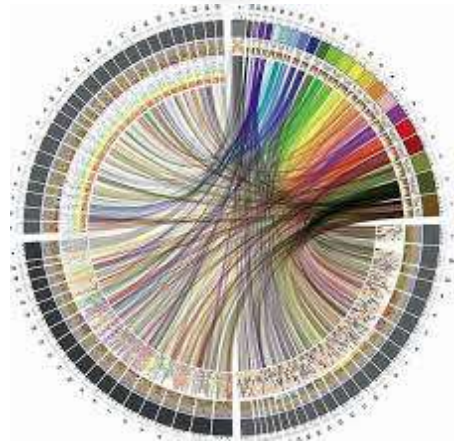
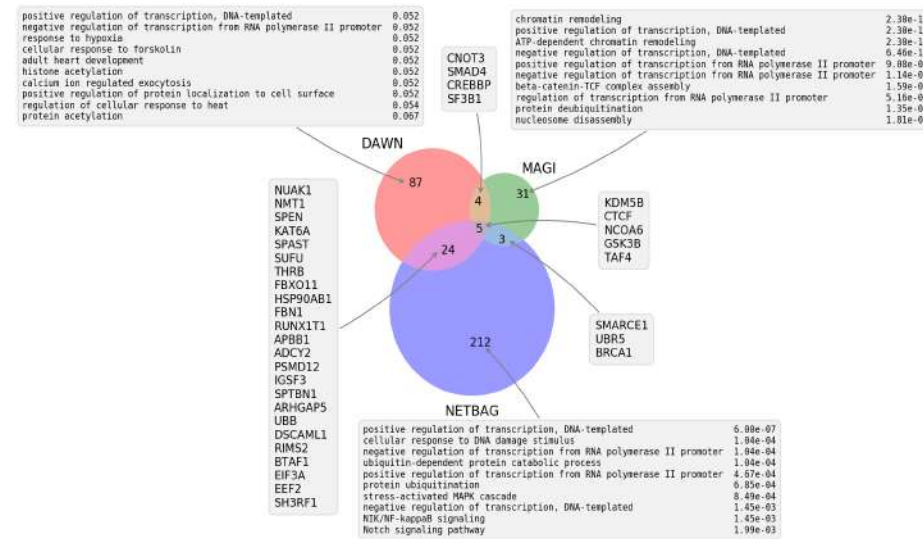
Data Science Pipeline – Story Telling

A data scientist also needs to communicate well.

Infographics and how you convey the story is important.



VS



Data Storage and Cloud

Database Systems

Relational databases, organized around tables, SQL

NoSQL databases for online distributed databases, eventual consistency: Cassandra, Hbase, MangoDB, Neo4j, DynamoDB

Cloud Storage

Ubiquitous computing, data access from everywhere

No worries on losing data

Cloud Computing

Distributed computing on large scale data

Map Reduce, Hadoop

Statistical Modeling

Parametric Models

Family of probability distributions with a finite number of parameters

For example: Binomial distribution has 2 (n, p)

Non-parametric Models

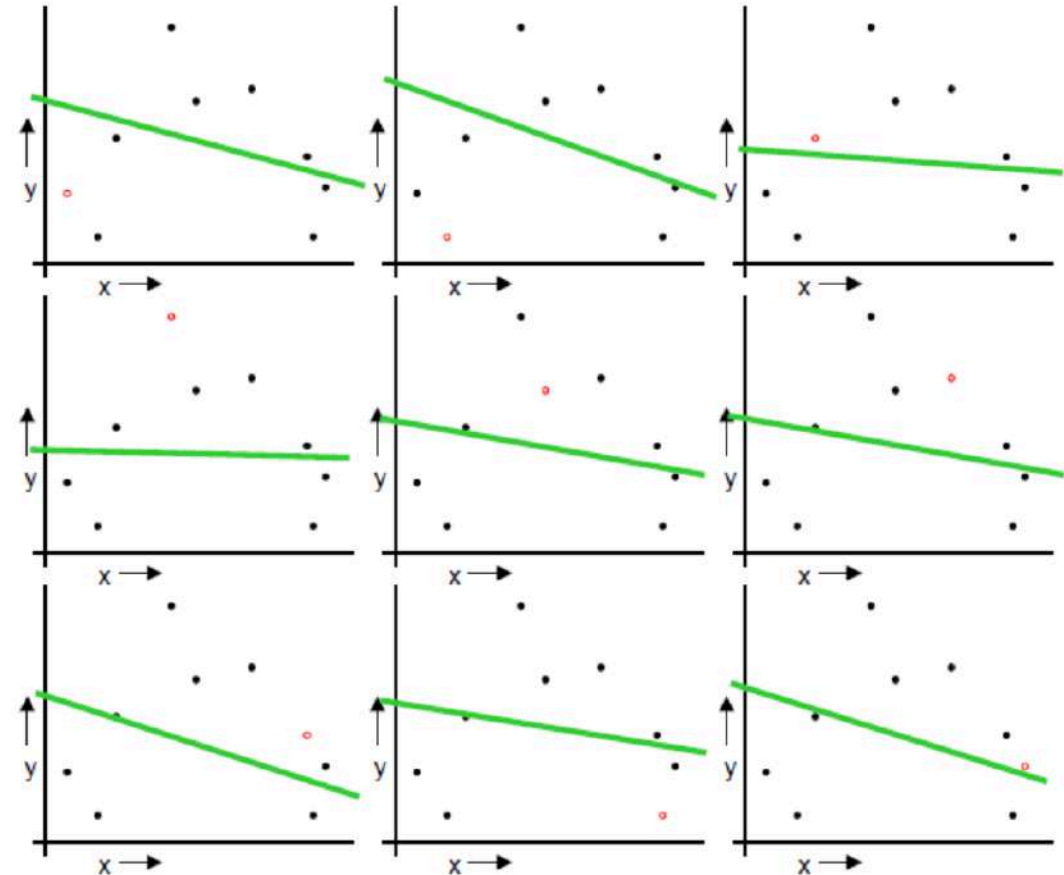
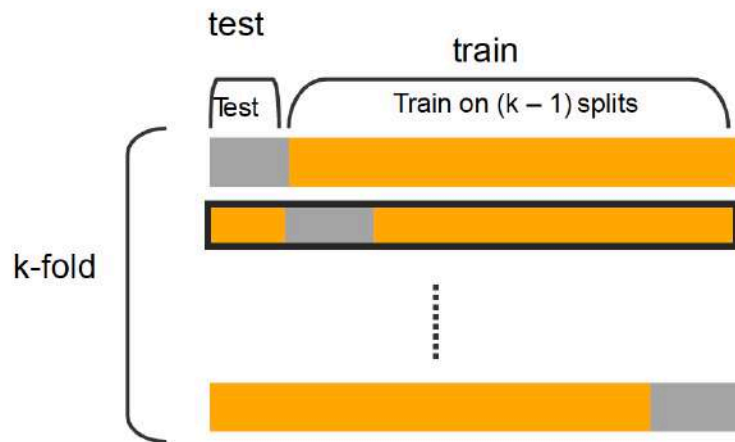
Parameter set is infinite dimensional i.e., grows with the data size. For example: k nearest neighbors classification.

Model Validation

Experimental Design

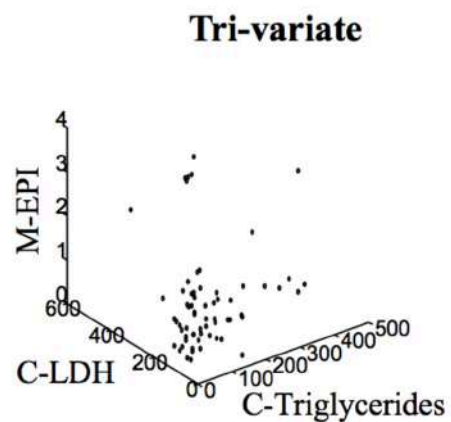
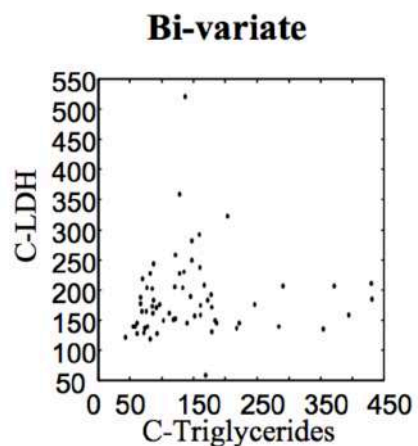
Cross Validation

Statistical Tests for validation



Unsupervised Learning

Feature extraction: Principal Component Analysis, t-SNE etc.



How can we visualize the other variables???

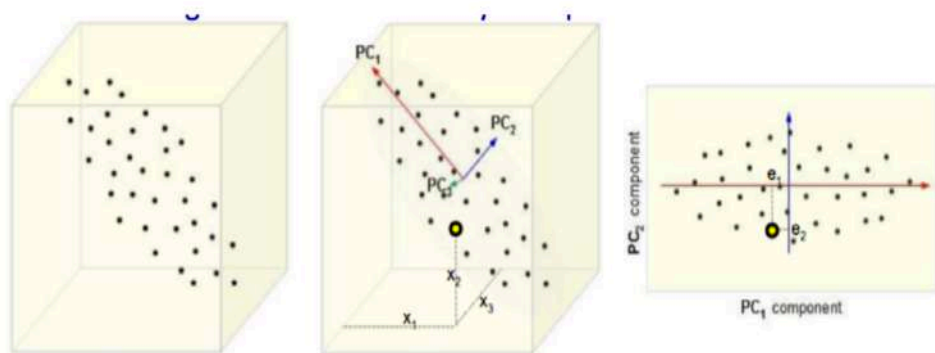
... difficult to see in 4 or higher dimensional spaces...



PC1



PC2

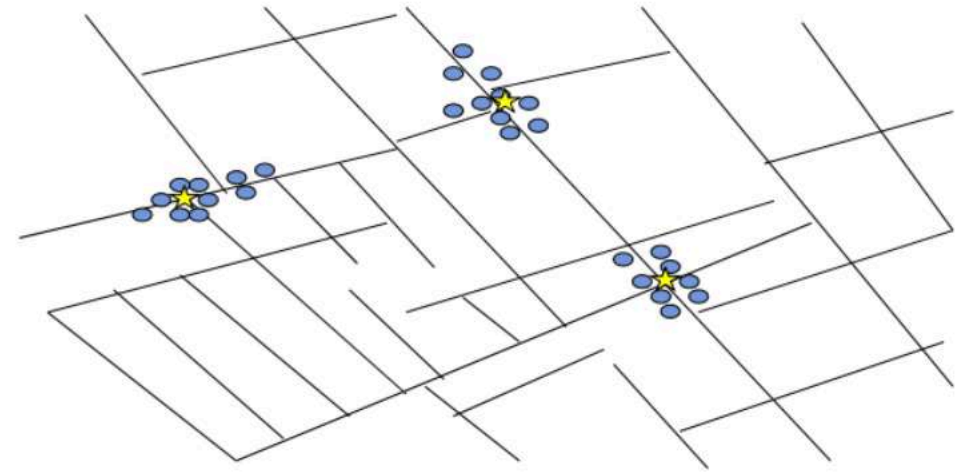


Unsupervised Learning – cont'd

Clustering: Finding groups of data points which are similar to each other.

John Snow, a London physician plotted the location of cholera deaths on a map during an outbreak in the 1850s.

The locations indicated that cases were clustered around certain intersections where there were polluted wells – thus exposing both the problem and the solution



From: Nina Mishra HP Labs

Unsupervised Learning – cont'd

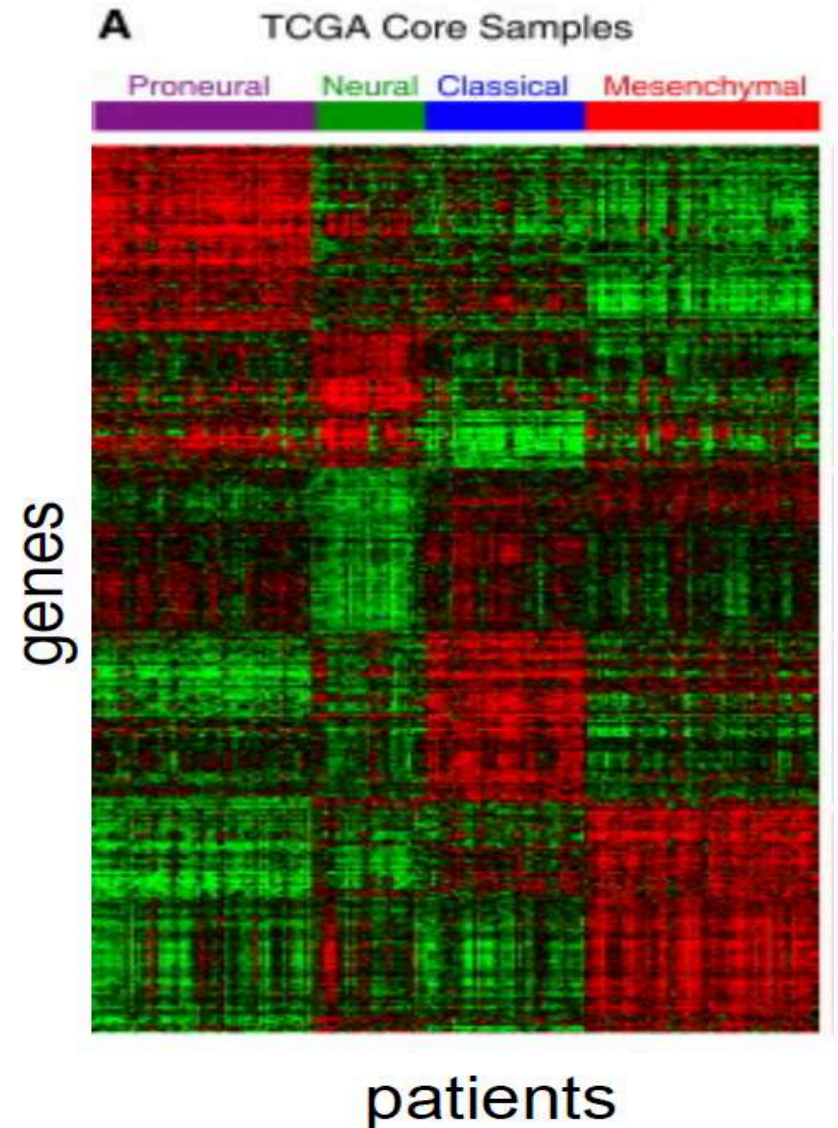
Clustering: Finding groups of data points which are similar to each other.

Given a sample of breast cancer patients and their gene activity level measurements. Can you find subgroups? (e.g., aggressive, mild etc.)

So many other applications:

- Targeted advertising

- LinkedIn contact suggestion



Unsupervised Learning – cont'd

Winner take all rule, competitive learning

Several algorithm examples

k-means

k cluster centers as means of assigned data points

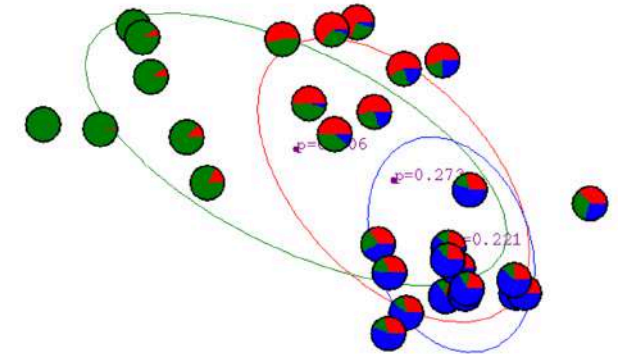
Gaussian Mixture Models

assumes k Gaussian processes generate data

Spectral Clustering

Generate eigenvalues/eigenvectors of the Laplacian of the similarity matrix

Use smallest eigenvalue and corresponding eigenvectors for dimension reduction



GMM example

Supervised Learning

When the data has labels learn a predictive model using features.

Neural Network Architectures

Perceptron

Multi Layer Perceptron

Convolutional Networks

Recurrent Neural Networks

Neural Network Training

Backpropagation

Optimizers

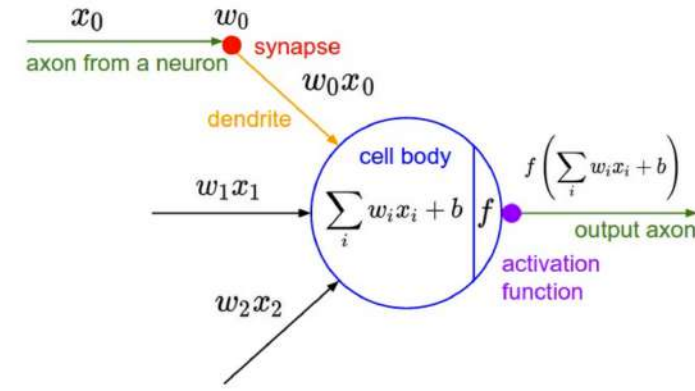
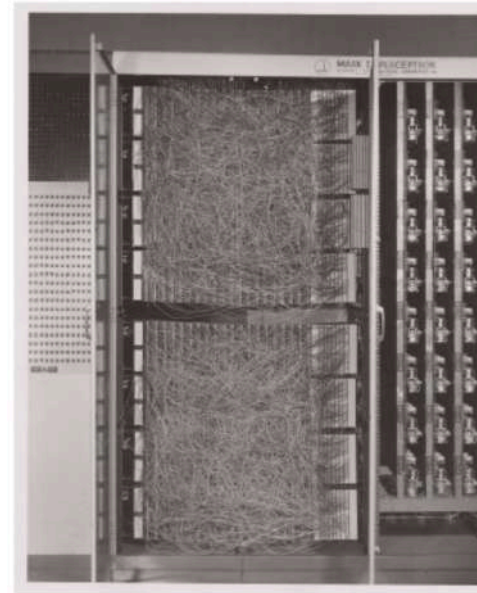
Support Vector Machines

Decision Trees

Ensemble Learning

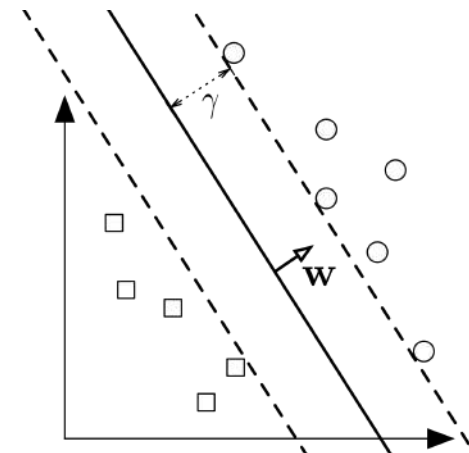
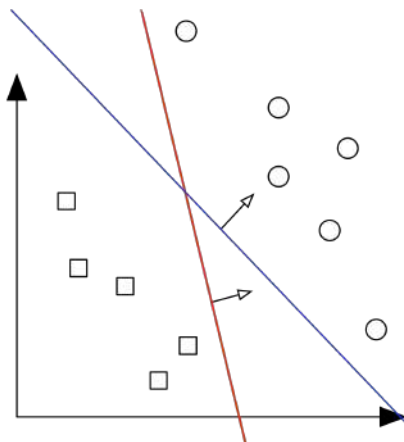
Random Forest

XGBoost, AdaBoost



'Mark I Perceptron at the Cornell Aeronautical Laboratory', hardware implementation of the first Perceptron (Source: Wikipedia / Cornell Library)

Neural Networks



SVM example – image source Cornell cs4780

Reinforcement Learning

Learning a policy by experience, reward, penalty like humans.

Q-Learning

Deep Q-Network



AlphaGo beats a 9-dan (professional) 4-1, gets 9-dan
Later AlphaZero is developed for GO, Shogi and Chess



AlphaZero beats a top professional player. First, time in a RTS game.
Again, by DeepMind.