# GE 461
# Introduction to Data Science

Spring 2022

# Course Website

All course related material will be provided in the course website

http://www.cs.bilkent.edu.tr/~ge461/2022Spring

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:** Tue 10:30 – 12:20 and Thursday 15:30 – 17:30.

**Where:** B-204.

**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. Instructors will decide.

    Projects will be uploaded to Moodle.

    Piazza will be used as the forum to discuss.

**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.

**Harvard Business Review**

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

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

vs

## Data Scientist Salaries

6,606 Salaries    Updated Jan 22, 2020

| Industries ∨ | Company Sizes ∨ | Years of Experience ∨ |

Average Base Pay

**$113,309** /yr

$83K Low   $113K Average   $154K High

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

## Computer Engineer Salaries

256,924 Salaries    Updated Jan 22, 2020

| Industries ∨ | Company Sizes ∨ | Years of Experience ∨ |

Average Base Pay

**$92,046** /yr
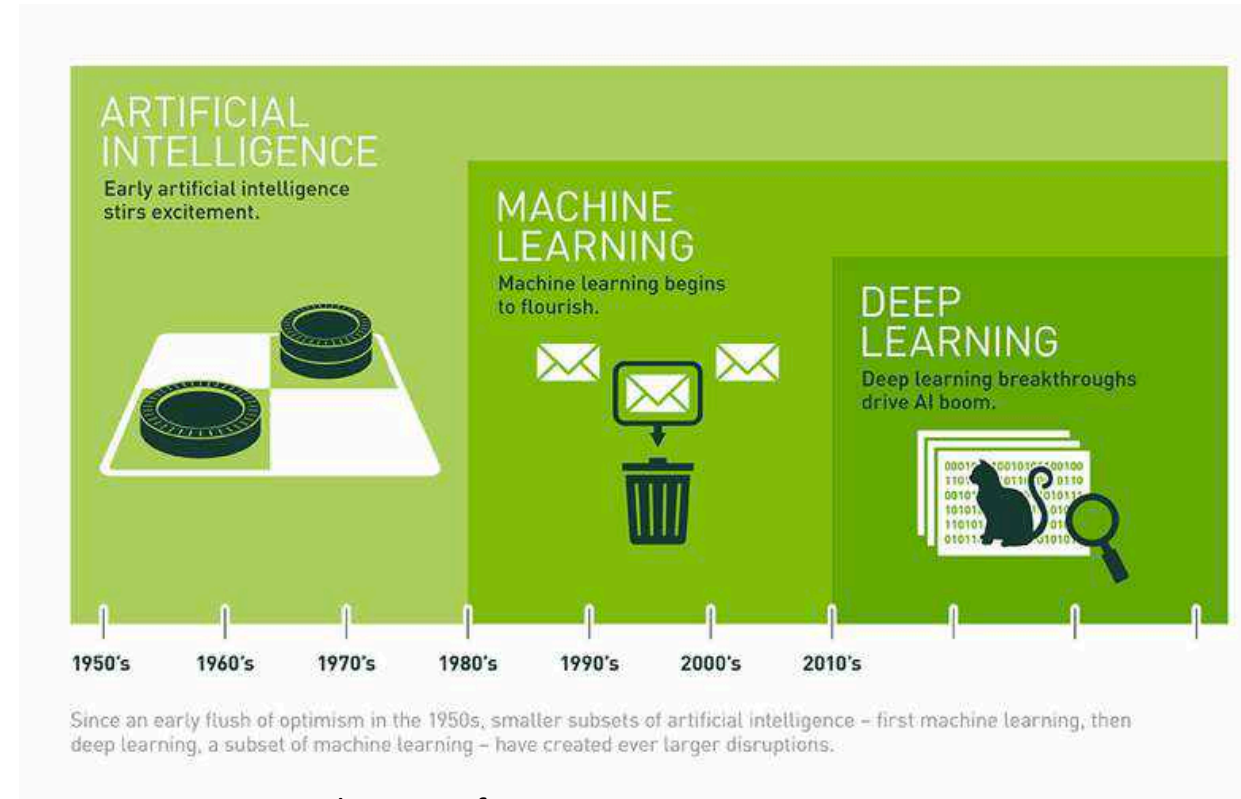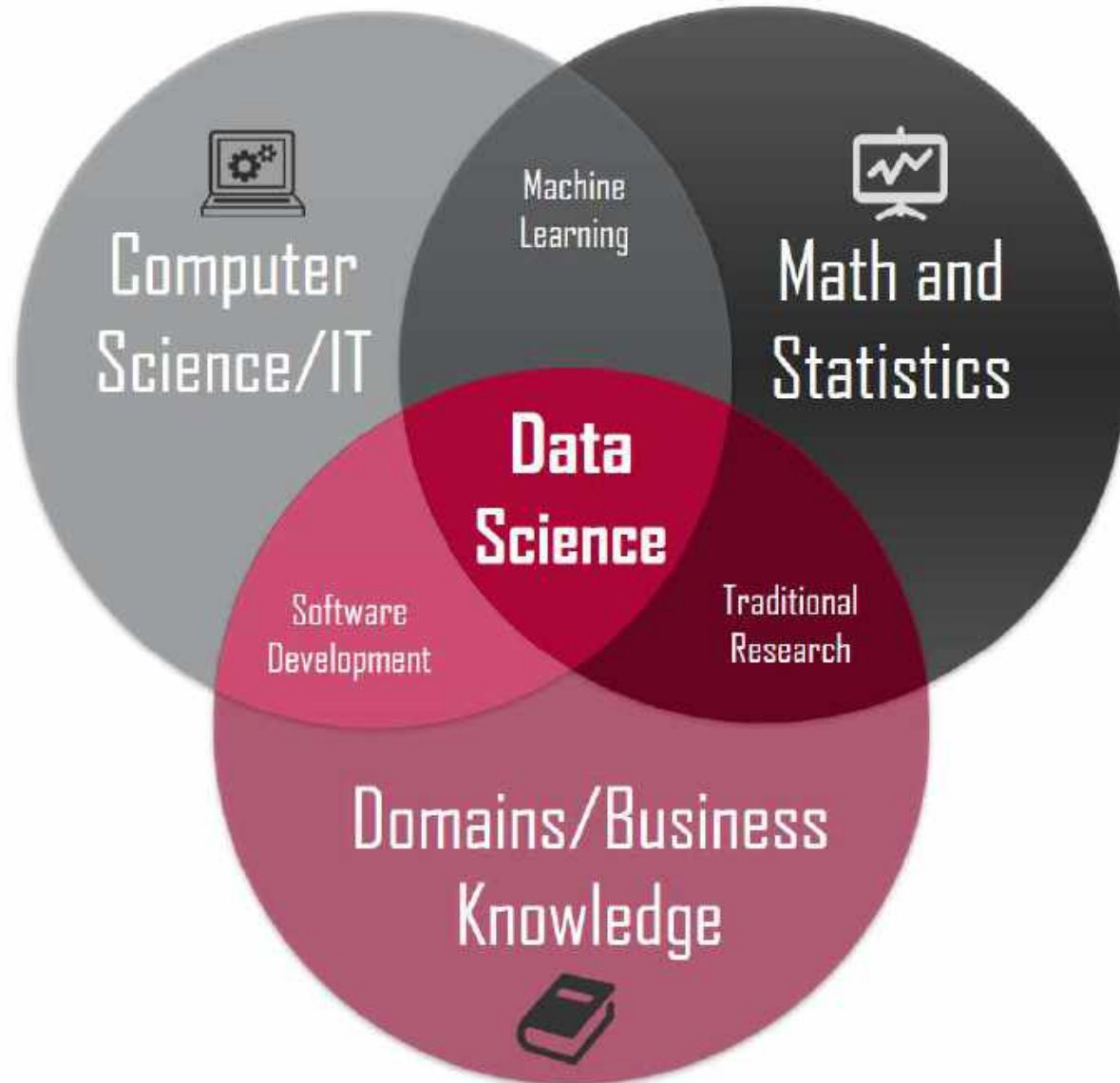
$63K Low   $92K Average   $134K High

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
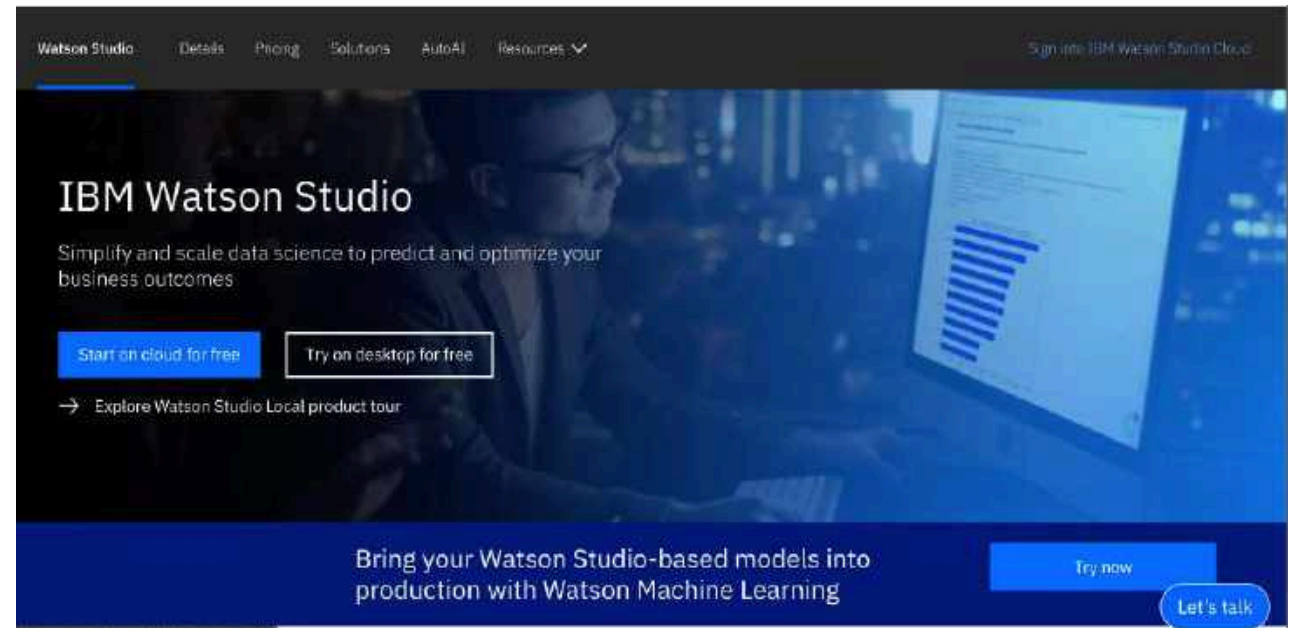
# What is NOT Data Science?



Data Science makes use of AI, ML, DL

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

Image source: Rob Tibshirani, Stanford Stats 101

# 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

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?*

**#LOVE** IS POSTED **23,211** TIMES

**GIPHY** SERVES UP **4,800,000** GIFS

**NETFLIX** USERS STREAM **694,444** HRS OF VIDEO

**GRUBHUB** RECEIVES **8,683** ORDERS

**INSTAGRAM** USERS POST **277,777** STORIES

**TWITCH** USERS VIEW **1,000,000** VIDEOS

**YOUTUBE** USERS WATCH **4,500,000** VIDEOS

**TUMBLR** USERS PUBLISH **92,340** POSTS

**TWITTER** USERS SEND **511,200** TWEETS

**390,030** APPS ARE DOWNLOADED

**188,000,000** EMAILS ARE SENT

**18,100,000** TEXTS ARE SENT

**SKYPE** USERS MAKE **231,840** CALLS

**GOOGLE** CONDUCTS **4,497,420** SEARCHES

**TINDER** USERS SWIPE **1,400,000** TIMES

**VENMO** PROCESSES **$162,037** TRANSACTIONS

**UBER** USERS TAKE **9,772** RIDES

**AIRBNB** BOOKS **1,389** RESERVATIONS

**AMERICANS** USE **4,416,720** GB OF INTERNET DATA

**INSTAGRAM** USERS POST **55,140** PHOTOS

2019 *every* **MINUTE** of the **DAY**

GLOBAL INTERNET POPULATION GROWTH 2013-2018 (IN BILLIONS)

2014 — 3.0 | 2016 — 3.4 | 2017 — 3.8 | 2018 — 4.3

Learn more at domo.com

---

# DATA NEVER SLEEPS 8.0

## How much data is generated *every minute?*

**ZOOM** HOSTS **208,333** PARTICIPANTS IN MEETINGS

**REDDIT** SEES **479,452** PEOPLE ENGAGE WITH CONTENT

**NETFLIX** USERS STREAM **404,444** HOURS OF VIDEO

**DOORDASH** DINERS ORDER **555** MEALS

**INSTAGRAM** USERS POST **347,222** STORIES

**FACEBOOK** USERS UPLOAD **147,000** PHOTOS

**YOUTUBE** USERS UPLOAD **500 HRS** OF VIDEO

**WHATSAPP** USERS SHARE **41,666,667** MESSAGES

**TWITTER** GAINS NEW USERS **319**

**$3,805** IS SPENT ON MOBILE APPS

**1,388,889** PEOPLE MAKE VIDEO/VOICE CALLS

**CONSUMERS SPEND** **$1,000,000** ONLINE

**MICROSOFT TEAMS** CONNECT **52,083** USERS

**AMAZON** SHIPS **6,659** PACKAGES

**INSTAGRAM** BUSINESS PROFILE ADS SEE CLICKS **138,889**

**SPOTIFY** ADDS **28 TRACKS** TO ITS MUSIC LIBRARY

**VENMO** USERS SEND **$239,196**

**TIKTOK** IS INSTALLED **2,704** TIMES

**LINKEDIN** USERS APPLY FOR **69,444** JOBS

**FACEBOOK** USERS SHARE **150,000** MESSAGES

2020 *every* **MINUTE** of the **DAY**

GLOBAL INTERNET POPULATION GROWTH 2014-2020 (IN BILLIONS)

2014 — 3.0 | 2016 — 3.4 | 2018 — 4.3 | 2020 — 4.5

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.

**TWITTER** USERS POST **575k** TWEETS

**TIKTOK** USERS WATCH **167M** VIDEOS

**GOOGLE** CONDUCTS **5.7M** SEARCHES

**DISCORD** USERS SEND **668k** MESSAGES

**12M** PEOPLE SEND AN IMESSAGE

**INSTAGRAM** USERS SHARE **65k** PHOTOS

**CLUBHOUSE** CREATES **208** ROOMS

**FACEBOOK** USERS SHARE **240k** PHOTOS

**SNAPCHAT** USERS SEND **2M** SNAPSHOTS

**FACEBOOK** USERS RECEIVE **44M** VIEWS

**AMAZON** CUSTOMERS SPEND **$283k**

**YOUTUBE** USERS STREAM **694k** HOURS

**6M** PEOPLE SHOP ONLINE

**NETFLIX** USERS STREAM **452k** HOURS

**STRAVA** ATHLETES SHARE **1.5k** ACTIVITIES

**TEAMS** CONNECTS **100k** USERS

**ZOOM** HOSTS **856** MINUTES OF WEBINARS

**SLACK** USERS SEND **148k** MESSAGES

**VENMO** USERS SEND **$304k**

**INSTACART** USERS SPEND **$67k**

EVERY **1 MINUTE** OF THE **DAY** PRESENTED BY DOMO

As of July 2021, the internet reaches 60% 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.

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 it all.
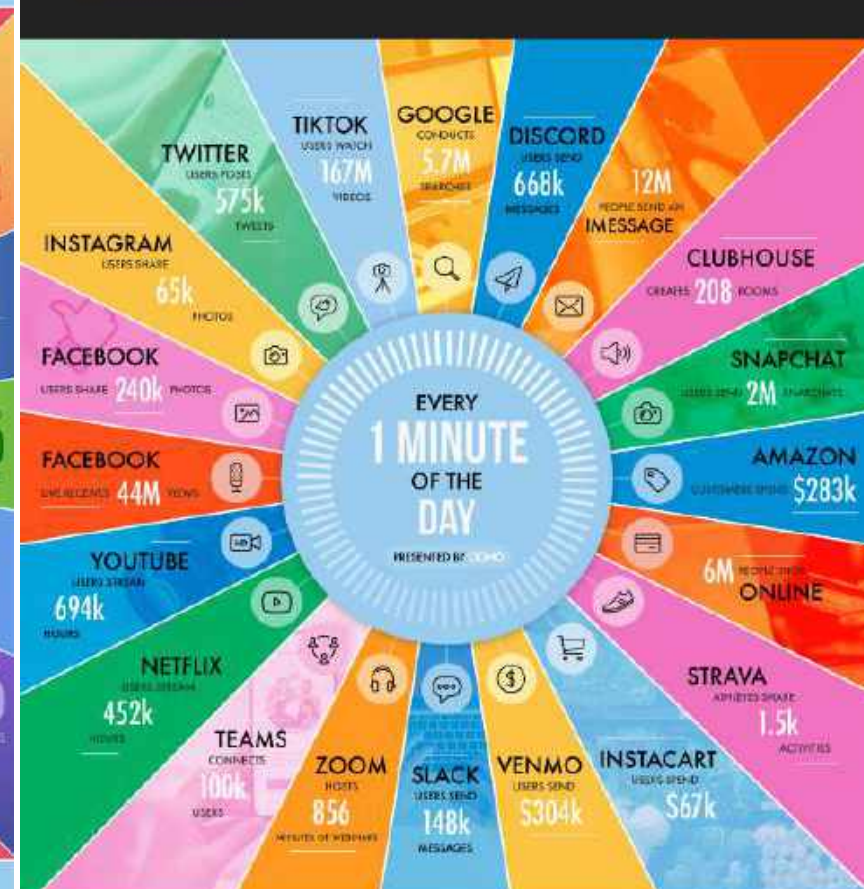
**Global Internet Population Growth** (IN BILLIONS)

2016 — 3.3 | 2018 — 4.3 | 2020 — 4.3 | 2021 — 5.2

Learn more at domo.com

# 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,… |

ACID = Atomicity, Consistency, Isolation and Durability      CAP = Consistency, Availability, Partition Tolerance

Slide by John Canny

# 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

Vertical exchange between levels

Horizontal exchange between columns

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

CIMMS, U of Wisconsin



Conv 1: Edge+Blob        Conv 3: Texture        Conv 5: Object Parts        Fc8: Object Classes

AlexNet/VGG-F visualization from Brown CSCI1430

# 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

# Some examples – Flu Trends

## Google Flu Trends

# 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.





Satterstrom *et al.,* CELL 2020

# Some examples – Comp. Biology

**Machine Learning** for Gene Risk Prediction

Build algorithms to predict the risk



TIME

U.S.  POLITICS  WORLD  TECH  ENTERTAINMENT  SUBSCRIBE

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 Window 2

Spatio-Temporal Window 1

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



Multi-Task Learning for Autism Gene Risk Prediction. Karakahya *et al.*, in prep.
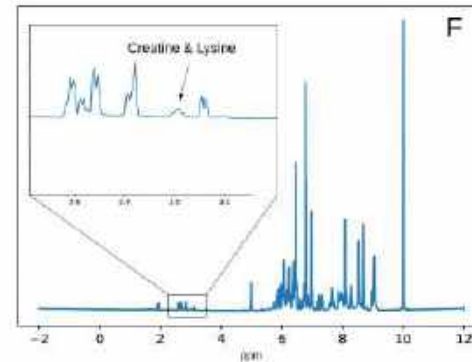
# Some examples – Comp. Biology

**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.
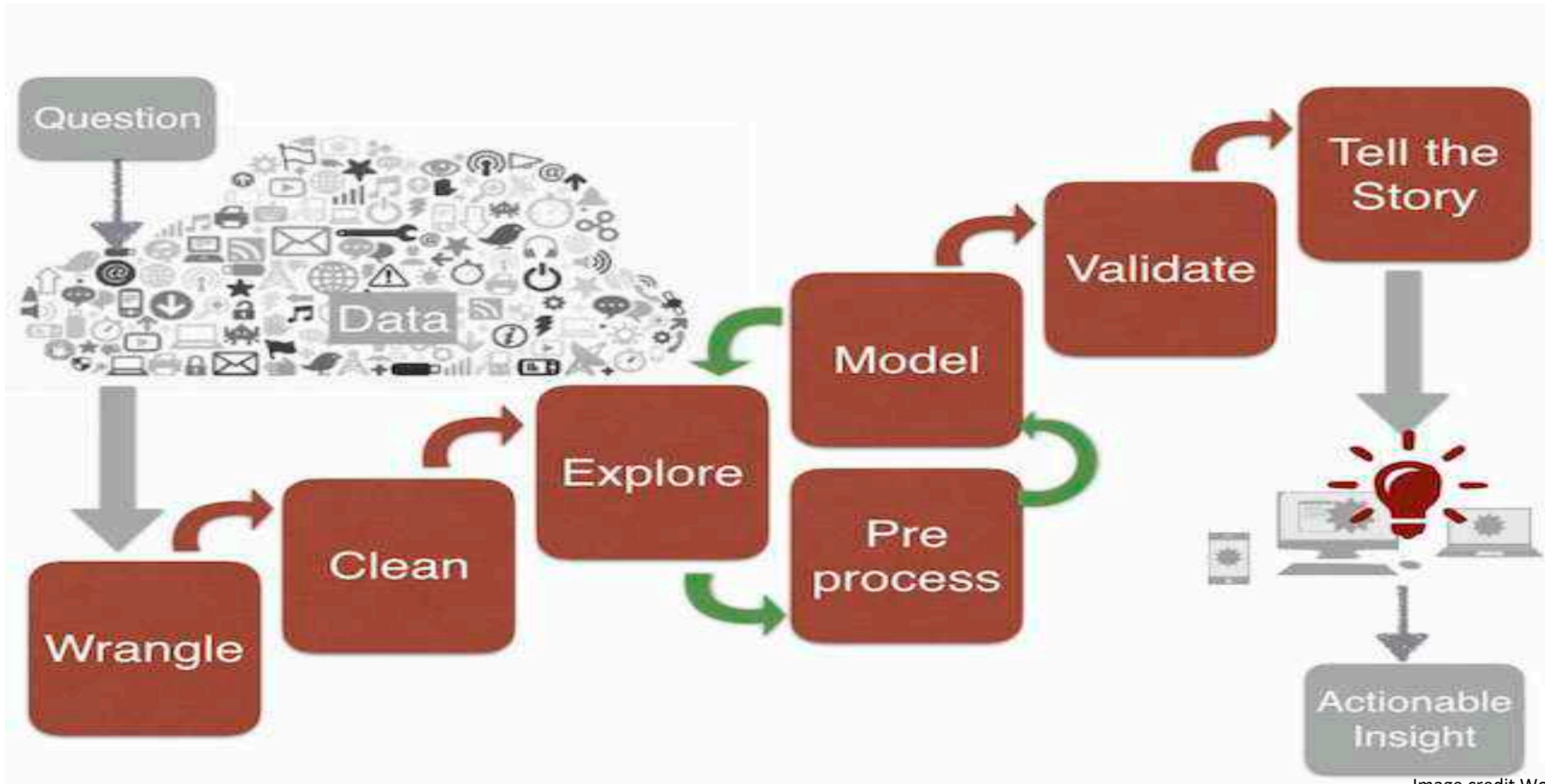


Karakaslar *et al.,* IEEE/ACM TCBB 2019, in press.

# Some examples – Comp. Biology

**Machine Learning** for Online Feedback to Surgeons

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



Cakmakci *et al.,* PLoS Comp. Bio 2020

# Data Science Pipeline



Image credit Wolfram Research

# Data Science Pipeline - Data Collection

Many data types, many ways

    Sensors

    Crowdsourcing, putting humans at work once computers fail:
        <span style="color:red">Mechanical Turk</span>

    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 Networks.

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.

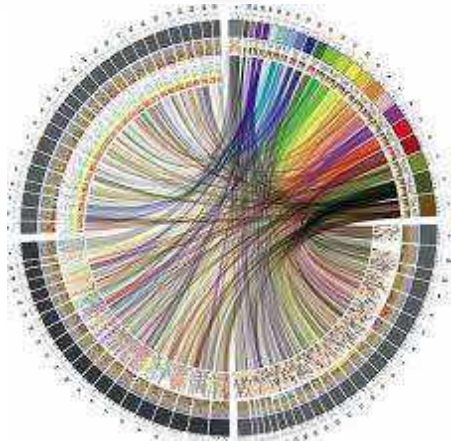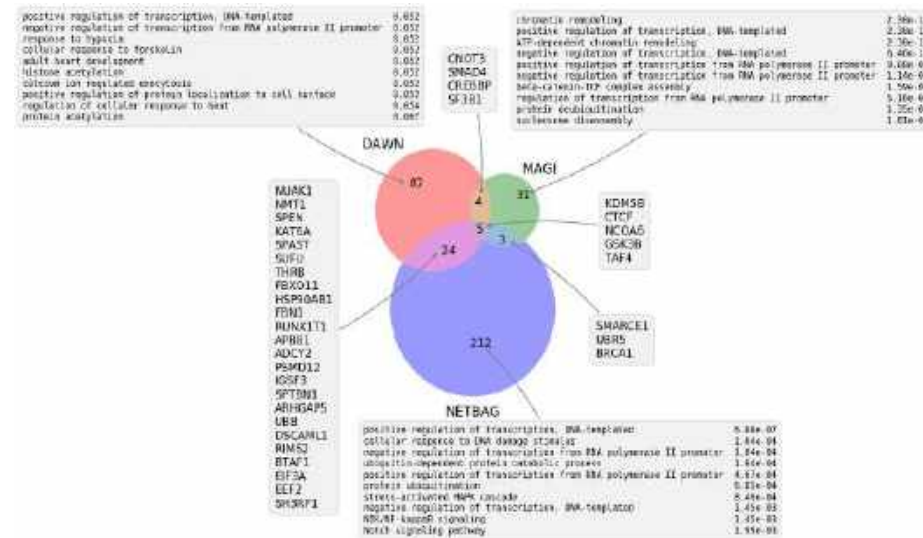Multiple hypothesis correction

Correlation is not causation.



Number of people who drowned by falling into a pool correlates with Films Nicolas Cage appeared in

# 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

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.
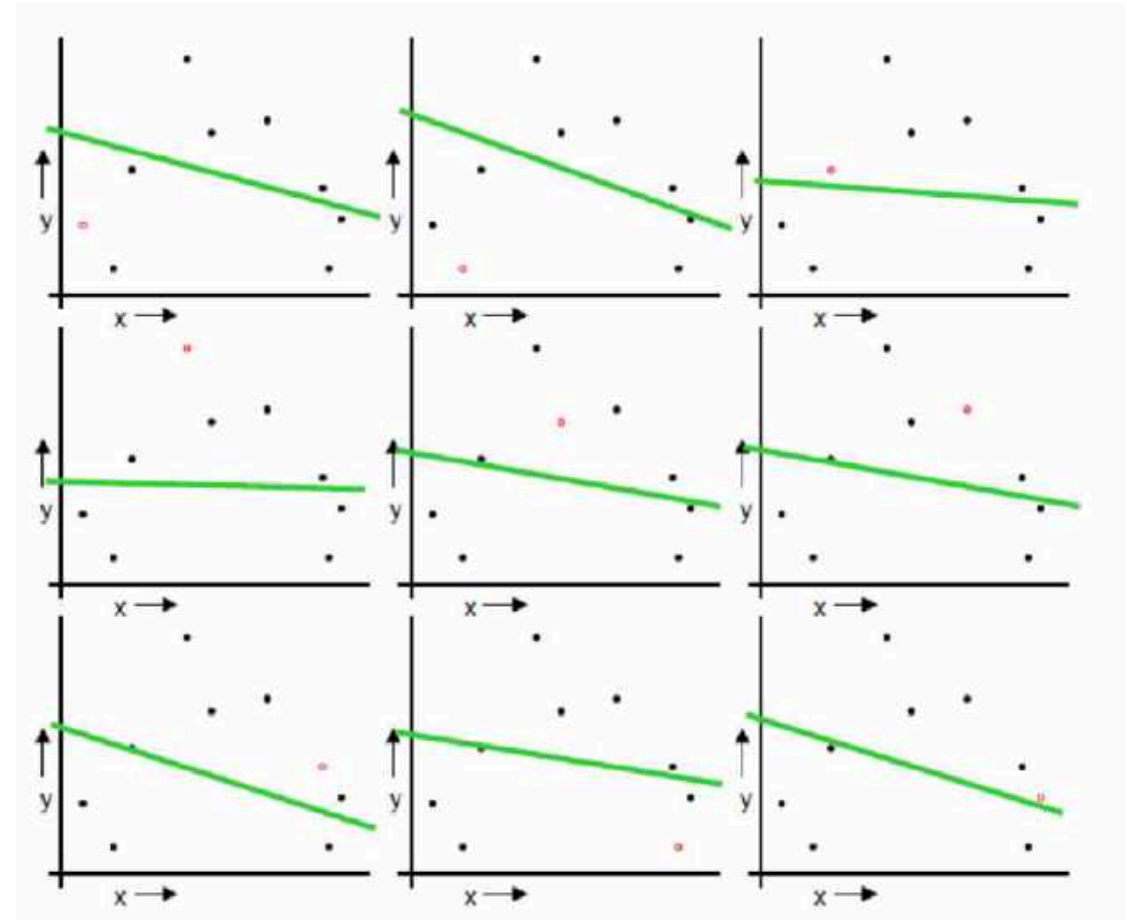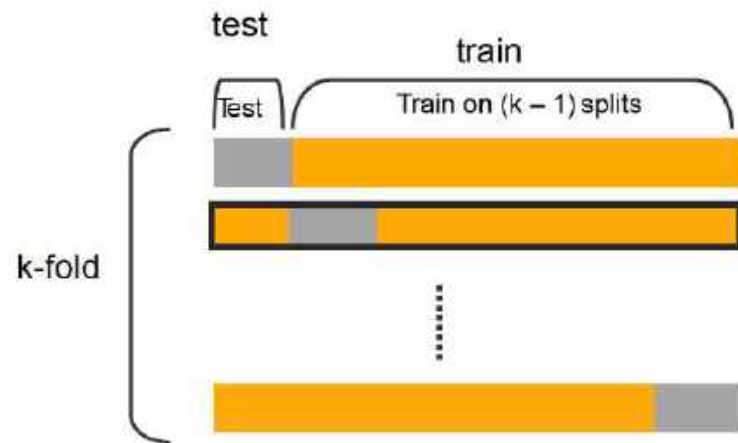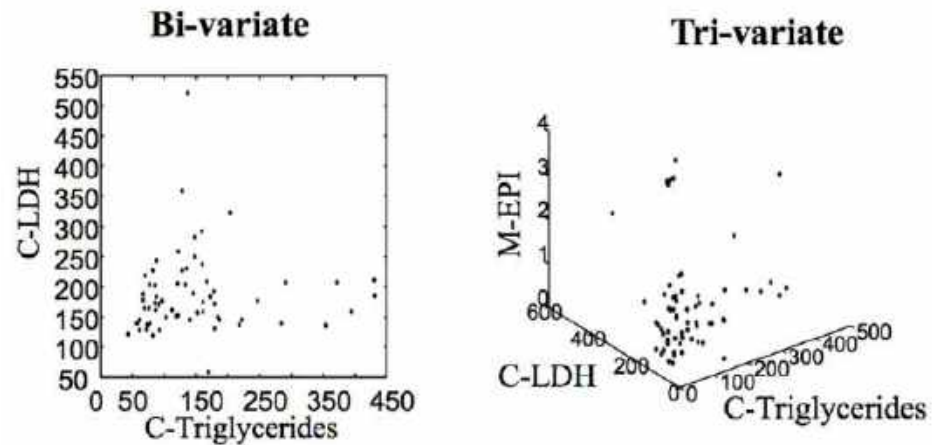
# Model Validation

Experimental Design

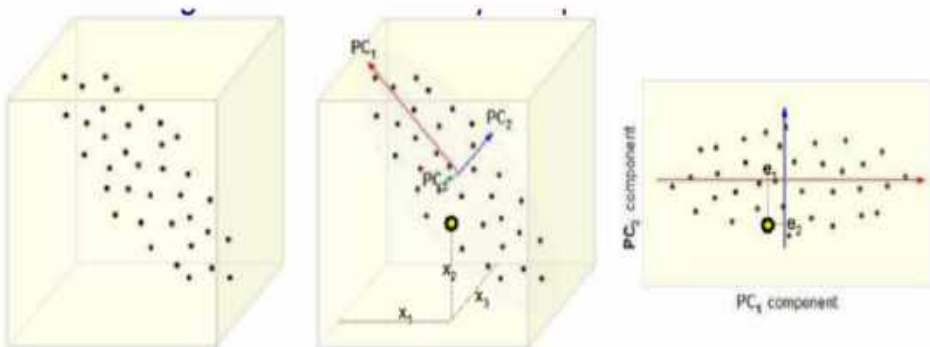Cross Validation

Statistical Tests for validation

# Unsupervised Learning

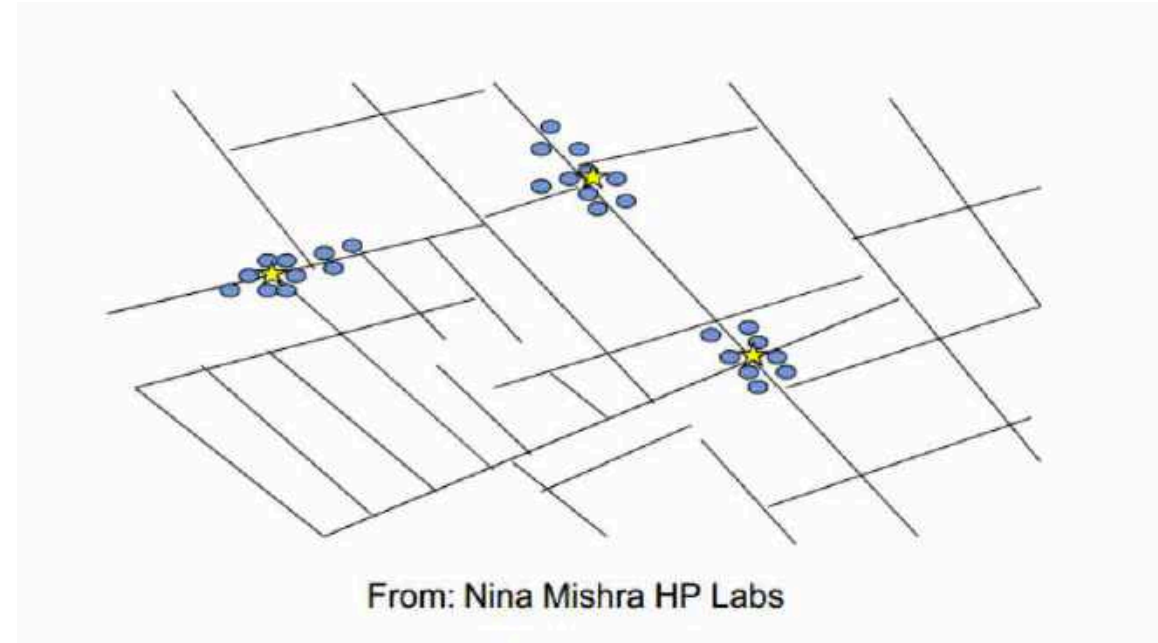Feature extraction: Principal Component Analysis, t-SNE etc.



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
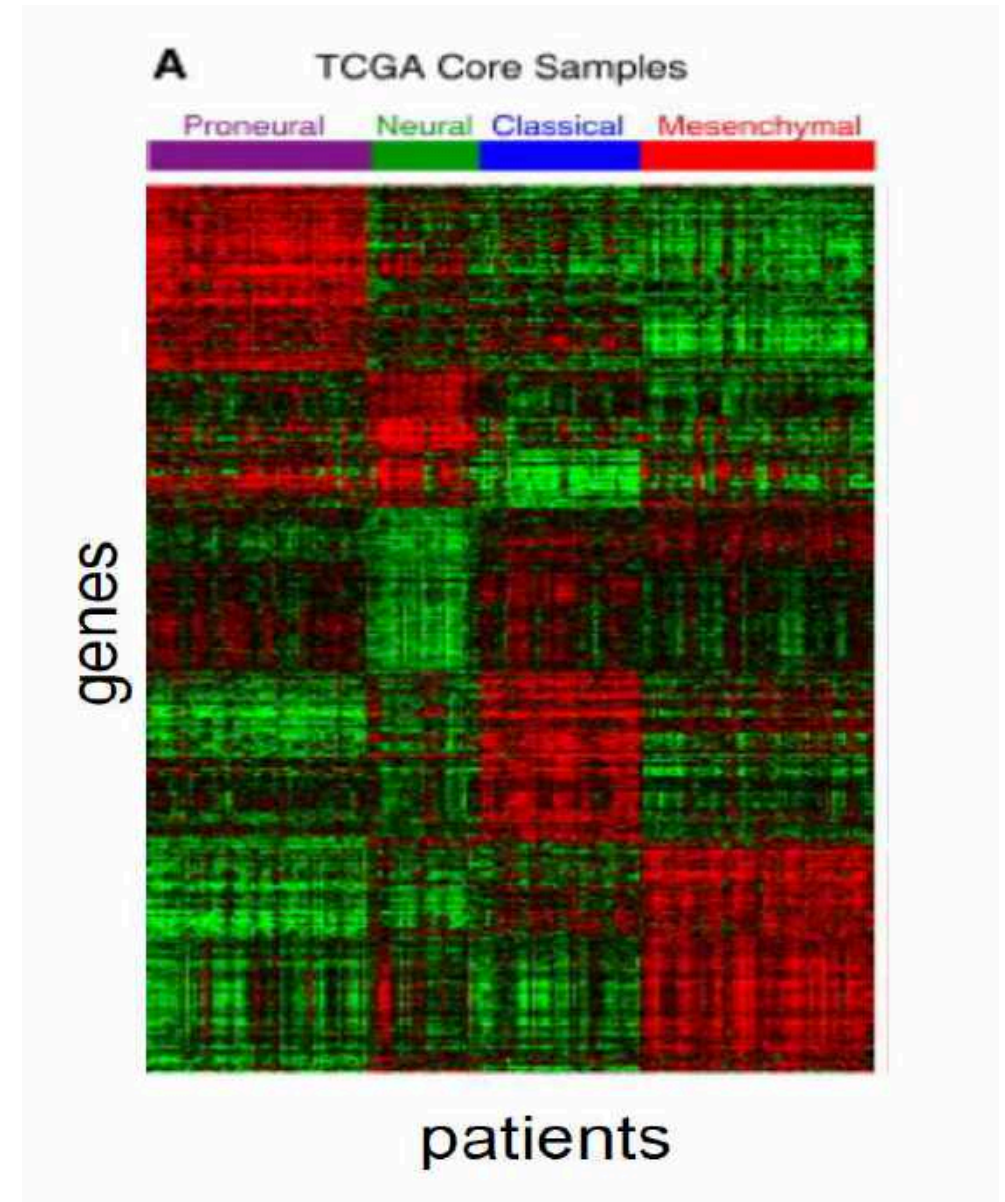
# Unsupervised Learning – cont'd

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



GMM example

Winner take all rule, competitive learning

Several algorithm examples

    k-means

        k cluster centers as means of assigned data points
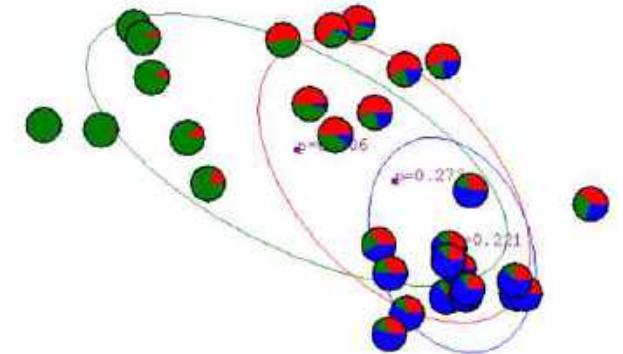
    Gaussian Mixture Models

        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

# 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

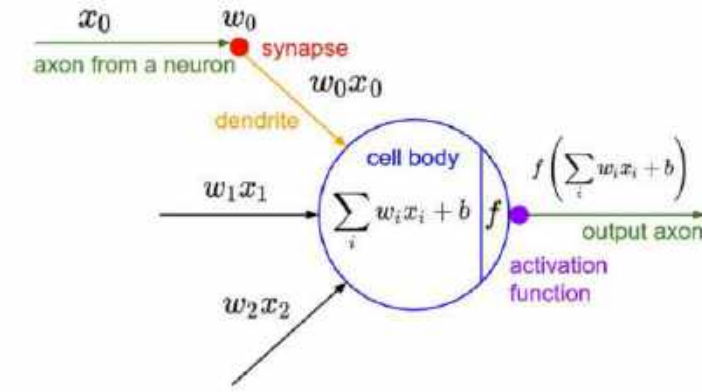    Neural Network Training

        Backpropagation

        Optimizers
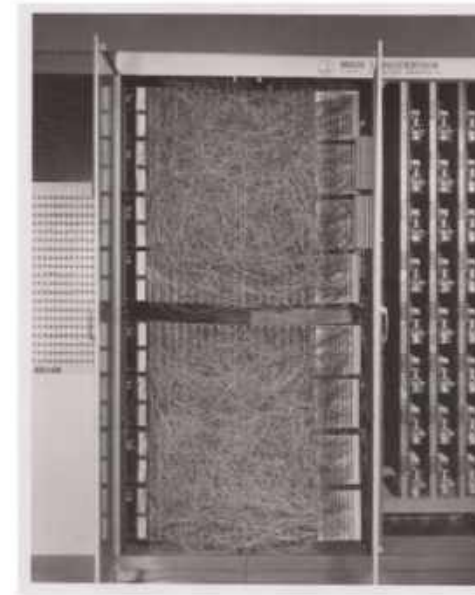
    Support Vector Machines

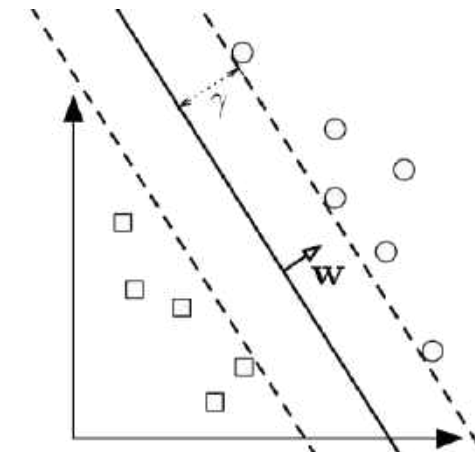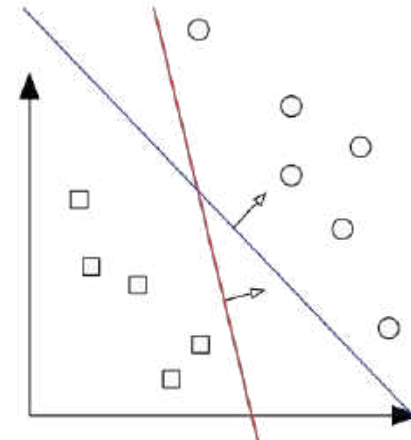    Decision Trees

    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.