

Lecture 1: Introduction

CIS 4190/5190

Spring 2025

Agenda

- **Logistics**

- Course description
- Tentative Schedule
- Grading

- **Introduction**

- Motivation
- Basic definitions
- Examples

Description

- **Key skills**

- Understand the standard machine learning algorithms
- Identify opportunities for applying machine learning (ML) algorithms
- Diagnose and debug issues in ML models

- Lectures will focus on developing mathematical understanding
- Assignments will focus on applying this understanding to implementing ML solutions

Prerequisites

- **Math:** University-level courses in probability, linear algebra, and multivariable calculus
 - Understand how to compute $\nabla_A(Ax)$ for a matrix A and vector x
 - Understand matrix ranks, inverses, and eigenvalues
 - Understand prior and posterior probabilities, $\mathbb{E}[X] = \int p(x)dx$, etc.
 - Tested in HW 1 (more on this later)
- **Programming:** Previously coded up projects (preferably in Python) that were at least 100 lines of code long
 - We have made available a Python primer on the course website for students who know how to program but do not know Python

Course Comparisons

- **CIS 4190/5190 (this course)**

- Basic mathematical ideas behind ML
- Apply existing ML algorithms to new problems as an engineer or researcher

- **CIS 5200**

- Deeper, more mathematically demanding introduction to ML
- Perhaps do fundamental ML research in the future

Course Comparisons

- **CIS 4190/5190 (this course)**
 - Basic mathematical ideas behind ML
 - Apply existing ML algorithms to new problems as an engineer or researcher
- **CIS 5220**
 - Deep learning techniques and applications in more detail

Course Comparisons

- **CIS 4190/5190 (this course)**
 - Basic mathematical ideas behind ML
 - Apply existing ML algorithms to new problems as an engineer or researcher
- **CIS 5450**
 - Data science workflow including data wrangling, ML modeling, and analytics
 - Scaling ML to big datasets and clusters

CIS 4190 vs. 5190

- 5190 will have extra, **mandatory** components in the HW, which are **optional** for 4190
- **Example**
 - HW may have 45 points for 4190, and 5 extra points for 5190 (total of 50)
 - Student taking 4190 will get 100% if they get at least 45 points (typically by skipping the 5190 problem, but not necessarily)
 - You cannot score more than 100%
 - The written and coding portion are counted separately; you cannot make up written points using coding points and vice versa

Waitlist

- **Class is almost at capacity (222 students)**
- **We are prioritizing students that either:**
 - Can only take the class this fall
 - Need to take the course this semester for a graduation requirement
- We plan to make all waitlist decision by next class
 - Unfortunately, we won't be able to respond to individual emails

Schedule (Tentative)

Week	Content	Homework
1	Introduction	
2	Linear Regression	
3	Linear Regression	HW 0 Due
4	Logistic Regression	
5	Neural Networks	HW 1 Due
6	KNNs + Decision Trees	
7	Unsupervised Learning	
8	Midterm 1	HW 2 Due
9	Computer Vision	
10	Natural Language Processing	
11	Reinforcement Learning	HW 3 Due
12	Ensembles	
13	Recommender Systems + Robustness	
14	Generative AI + Ethics	HW 4 Due
15	Midterm 2	

Grading Scheme (Tentative)

- **Homeworks (5×):** 30%
- **Project:** 25%
- **Midterms exams:** 40%
- **Good citizenship points:** 5%

Grading Scheme (Tentative)

- **A+:** 95+
- **A:** 90-95
- **A-:** 85-90
- **B+:** 80-85
- **B:** 75-80
- **B-:** 70-75
- **Lower passing grades:** 50-70

- May be curved up

Late Policy

- For each hour late, lose 0.5% on the points for that assignment
 - **Homeworks only**
 - Max 48 late hours per assignment
- **Example**
 - Submit HW 1 20 hours late
 - Lose $20 \times 0.5 = 10\%$ on HW 1 (0.5% of overall grade)
- If you have a medical reason, email us a copy of your medical visit report, and we will grant an extension (typically 2 days)
 - We will consider other reasons on a case-by-case basis

Course Instructors



Osbert Bastani
Assistant Professor, CIS



Mingmin Zhao
Assistant Professor, CIS

TA & Office Hours

- We have a team of 14 amazing TAs
- Each TA (and professor) will have 1-2 hours of office hours each week
 - Times still being decided

Communication

- All materials will be posted on the course website:
 - <https://www.seas.upenn.edu/~cis5190/spring2025/>
- We plan to post lecture videos up to one week after the lecture
 - We make no guarantees about quality/availability, e.g., due to technical issues
- We will use **Ed Discussion** for questions and course discussions
 - Send a message to “instructors” to contact the TAs and me
 - You can contact me directly on Ed Discussion or by email (posted on course website)

Homework Schedule

- **5 homeworks**

- Released roughly every other Wednesday
- Typically due Wednesday 2 weeks later (see schedule for actual deadline)

- **HW 0**

- Designed to test mathematical background
- **Expected time:** 3 hours
- Opportunity to get used to the workflow

Homework

- **Written problems:** GradeScope submission
 - LaTeX encouraged; handwritten + scanned at your own risk
 - Won't be graded if you don't annotate your answers correctly!
- **Coding problems:** AutoGrader + GradeScope submission of notebook
 - Colab/iPython notebook skeletons; AutoGrader as unit tests within skeleton
 - Only difference between AutoGrader and unit tests is different data
 - If code passes the unit tests and you didn't "game" it, it should pass AutoGrader
- Discussion permitted for clarifications, but never share solutions/code; acknowledge all your discussions at the beginning of your report

First Assignments

- **HW 0:** Released today, **due 1/29**
 - No office hours planned for HW 0
 - You can ask questions via Ed Discussion, but we will only answer clarifying questions
- More on the project later in the course

Agenda

- **Logistics**

- Course description
- Tentative Schedule
- Grading

- **Introduction**

- Motivation
- Basic definitions
- Examples

What is Machine Learning?

“Learning is any process by which a system improves performance from experience.”

Herbert Simon



What is Machine Learning?

“Machine learning ... gives computers the ability to learn without being explicitly programmed.”

Arthur Samuel

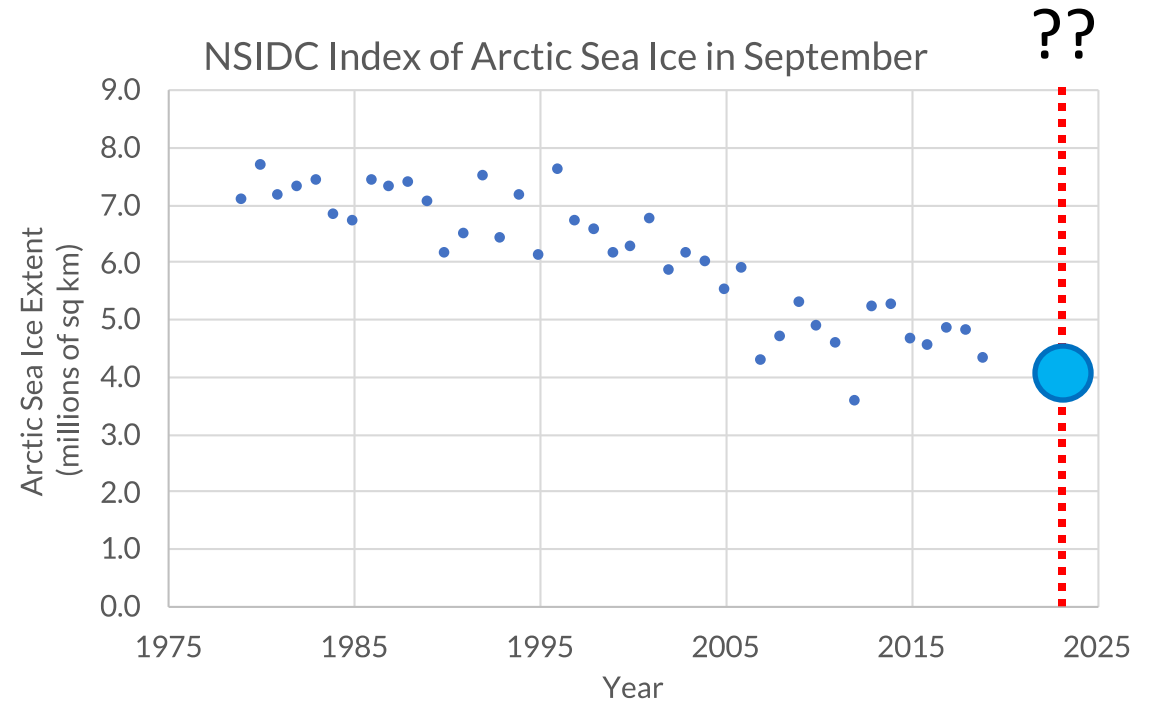


What is Machine Learning?

- **Tom Mitchell:** Algorithms that
 - improve their **performance** P
 - at **task** T
 - with **experience** E
- A well-defined machine learning task is given by (P, T, E)



Example: Prediction



Example: Prediction

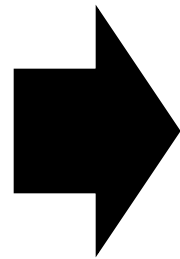
- **Tom Mitchell:** Algorithms that
 - improve their **performance** P
 - at some **task** T
 - with **experience** E
- T = predict Arctic sea ice extent
- P = prediction error (e.g., absolute difference)
- E = historical data



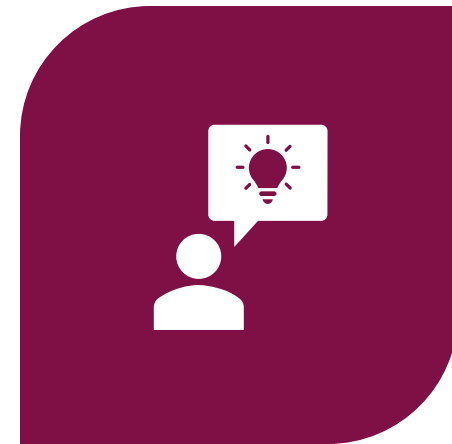
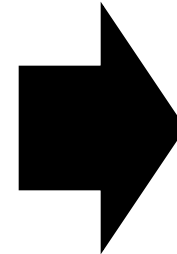
Machine Learning for Prediction



Data Z



Machine learning
algorithm

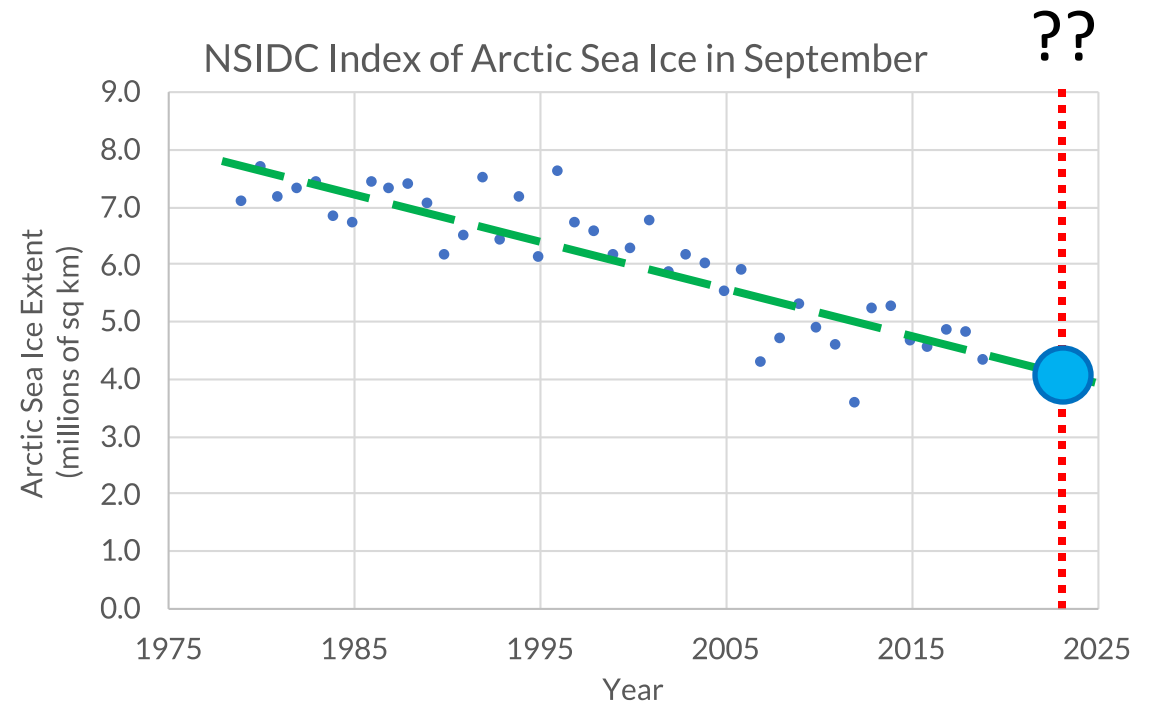


Model f

Machine Learning for Prediction



Example: Prediction

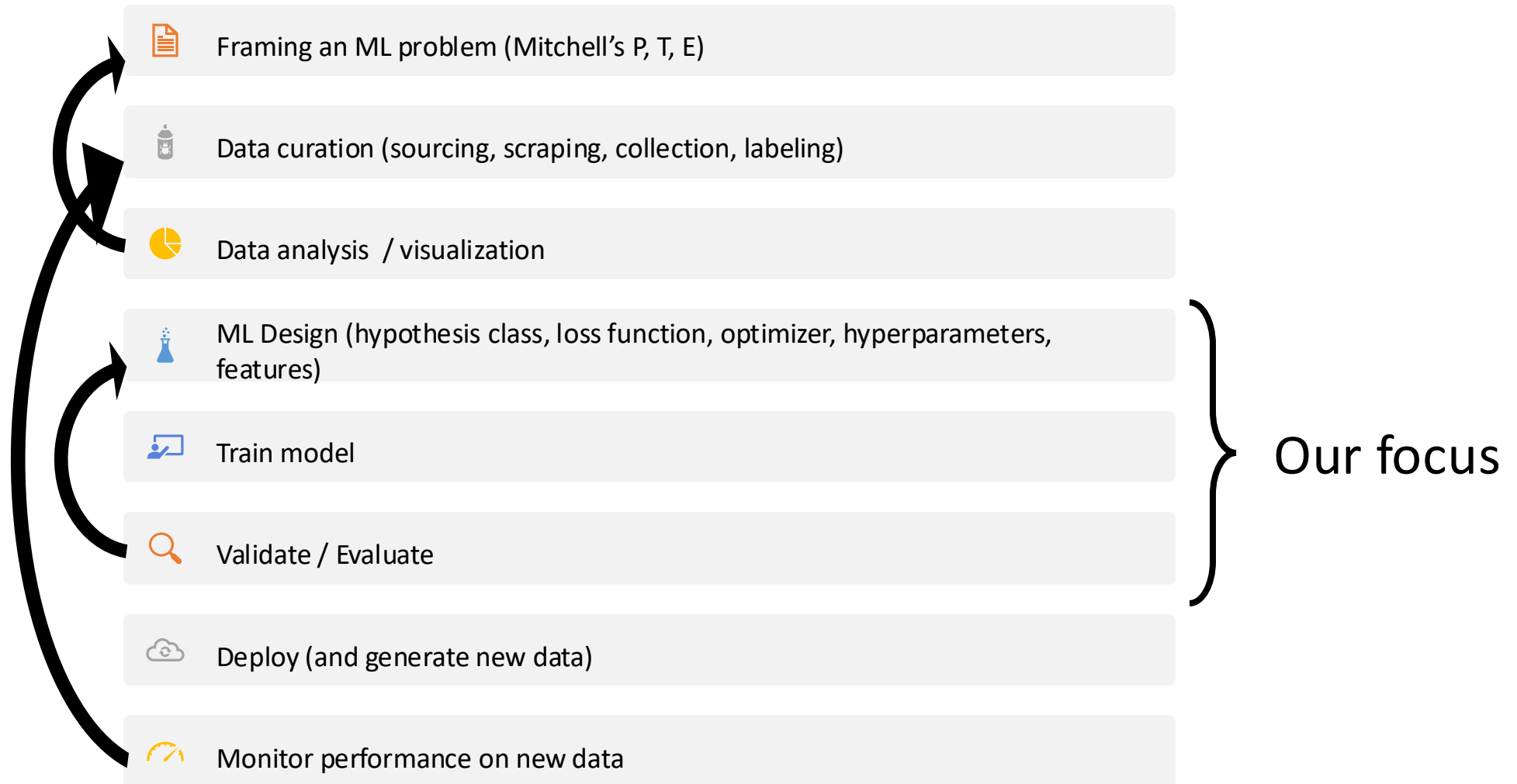


Example: Game Playing

- **Tom Mitchell:** Algorithms that
 - improve their **performance** P
 - at some **task** T
 - with **experience** E
- T = playing Chess
- P = win rate against opponents
- E = playing games against itself



Machine Learning Workflow



Types of Learning

- **Supervised learning**

- **Input:** Examples of inputs and desired outputs
- **Output:** Model that predicts output given a new input

- **Unsupervised learning**

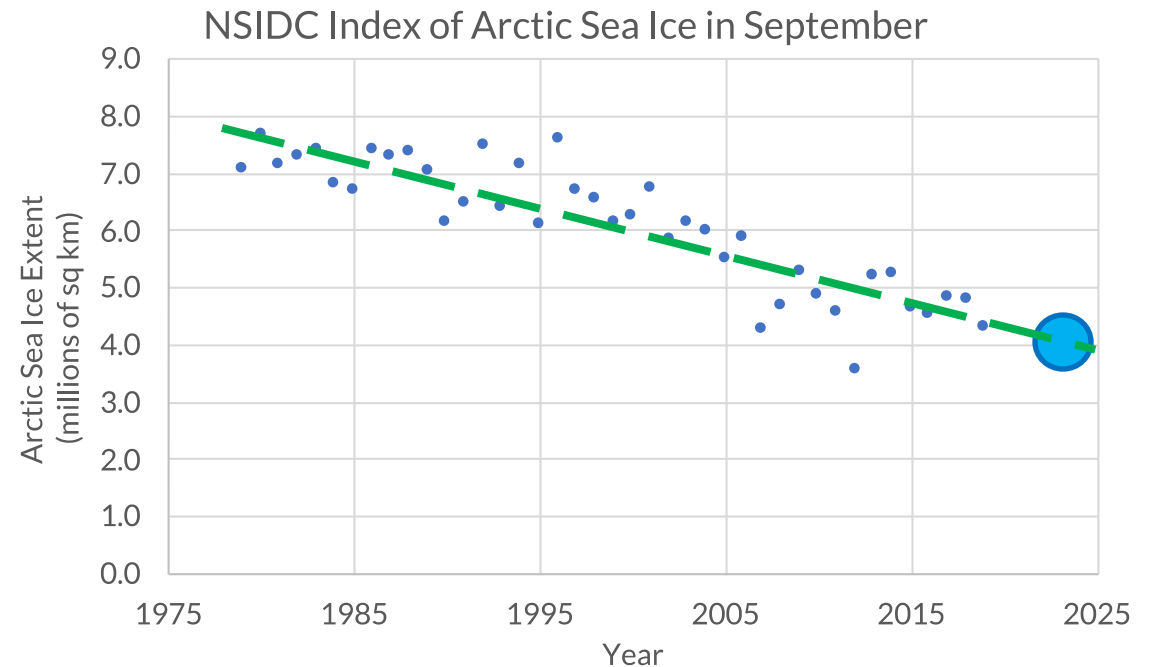
- **Input:** Examples of some data (no “outputs”)
- **Output:** Representation of structure in the data

- **Reinforcement learning**

- **Input:** Sequence of interactions with an environment
- **Output:** Policy that performs a desired task

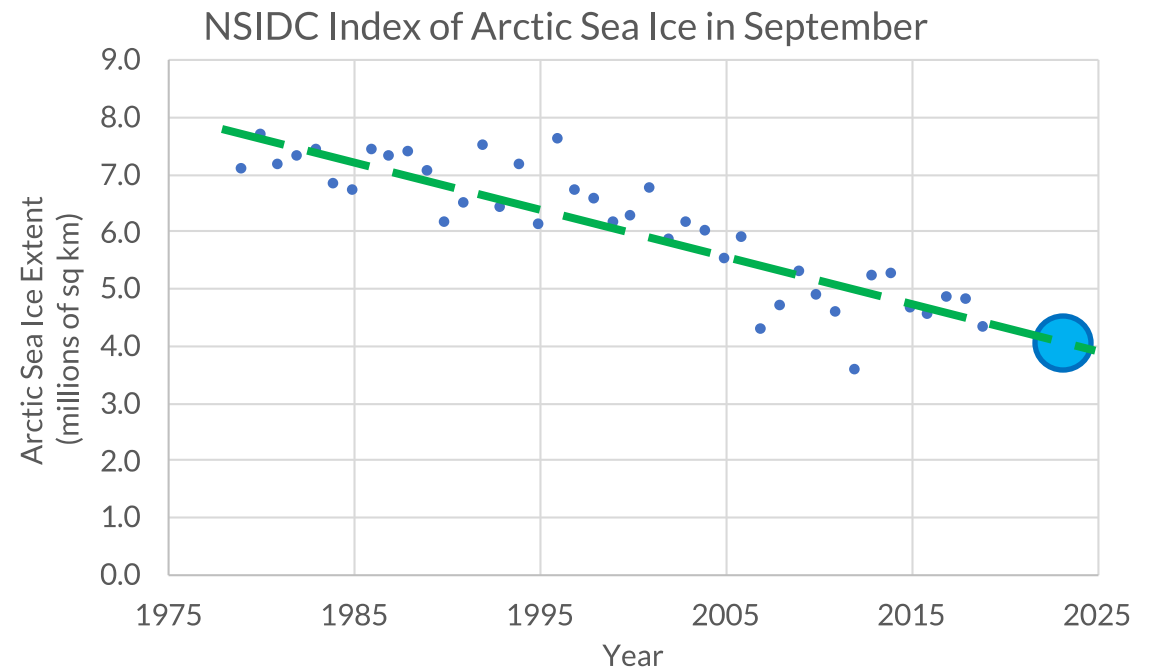
Supervised Learning

- Given $(x_1, y_1), \dots, (x_n, y_n)$, learn a function that predicts y given x



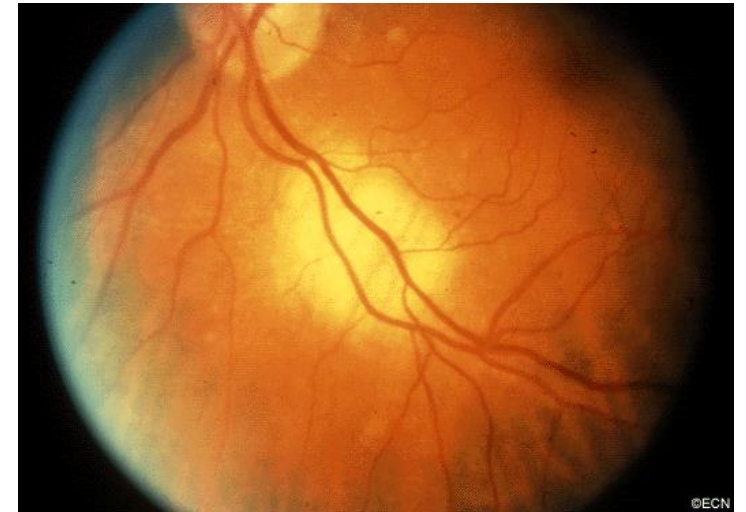
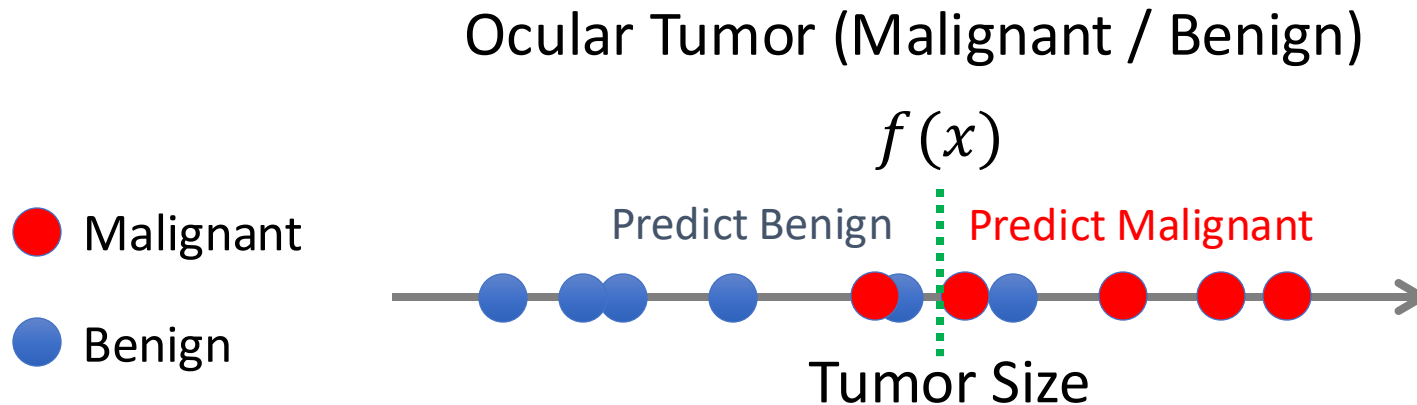
Supervised Learning

- Given $(x_1, y_1), \dots, (x_n, y_n)$, learn a function that predicts y given x
- **Regression:** Labels y are real-valued



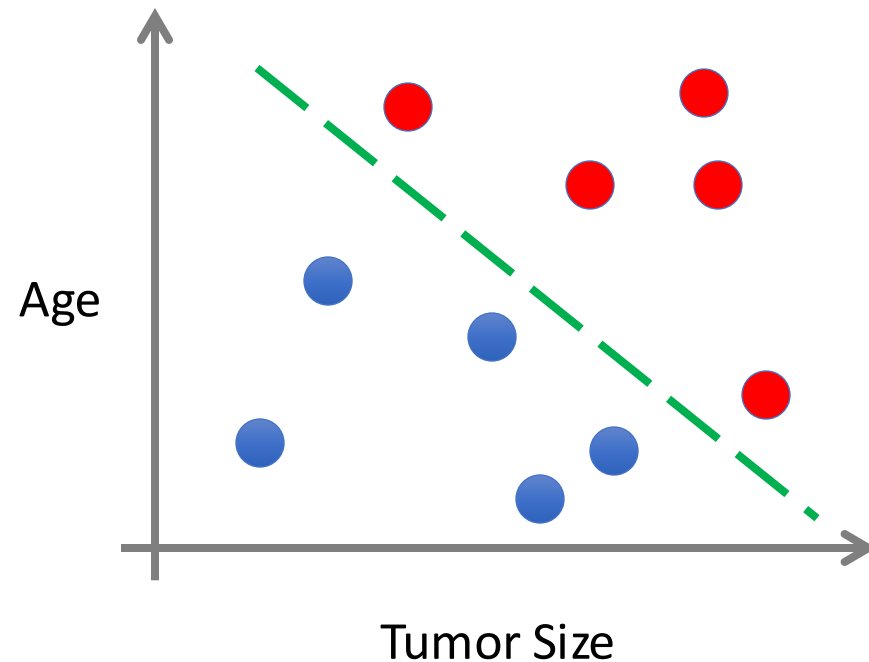
Supervised Learning

- Given $(x_1, y_1), \dots, (x_n, y_n)$, learn a function that predicts y given x
- **Classification:** Labels y are categories

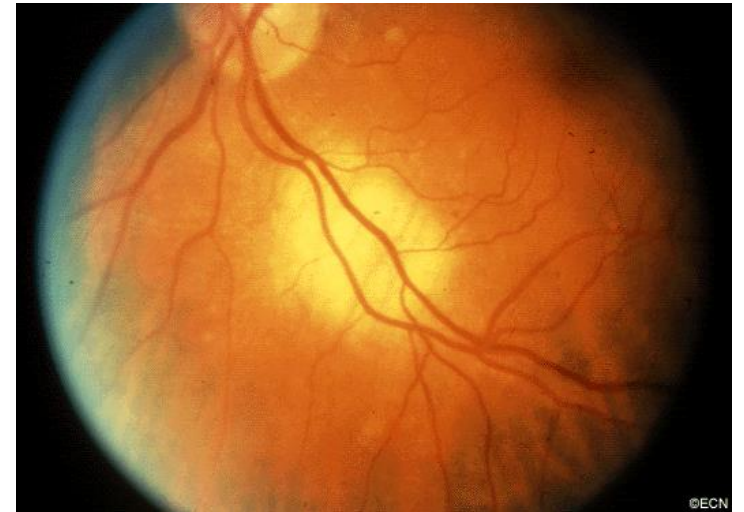


Supervised Learning

- Given $(x_1, y_1), \dots, (x_n, y_n)$, learn a function that predicts y given x
- Inputs x can be multi-dimensional

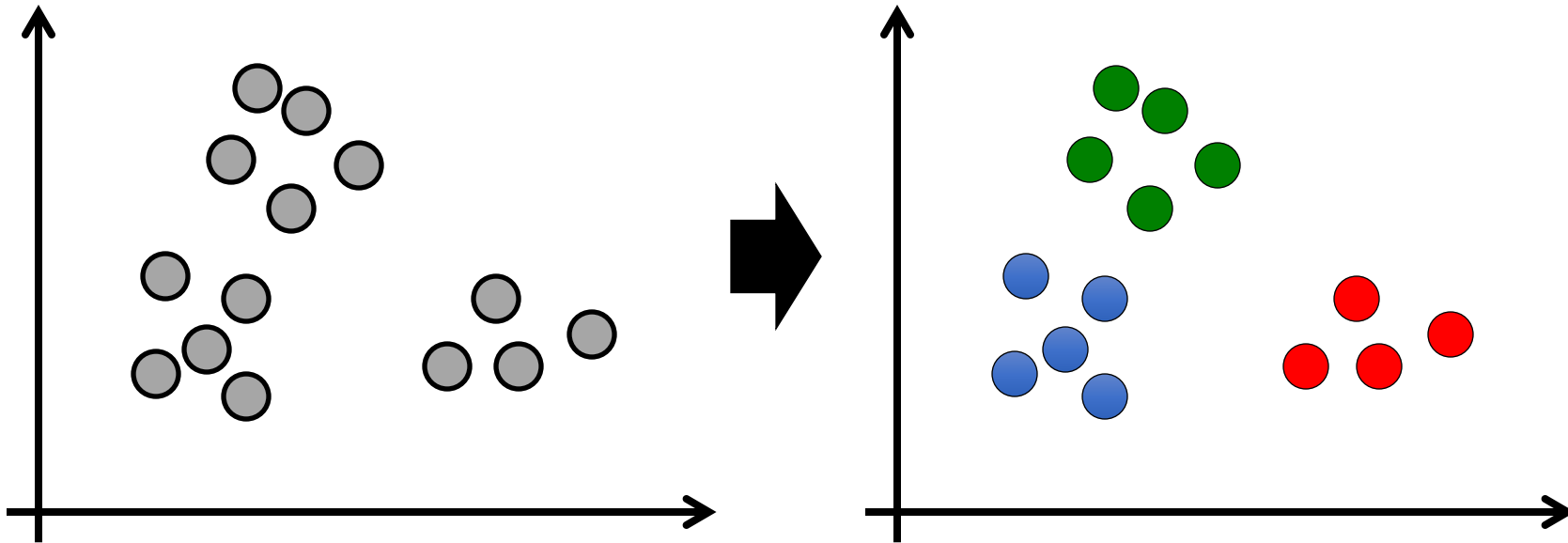


- Patient age
- Clump thickness
- Tumor Color
- Cell type
- ...

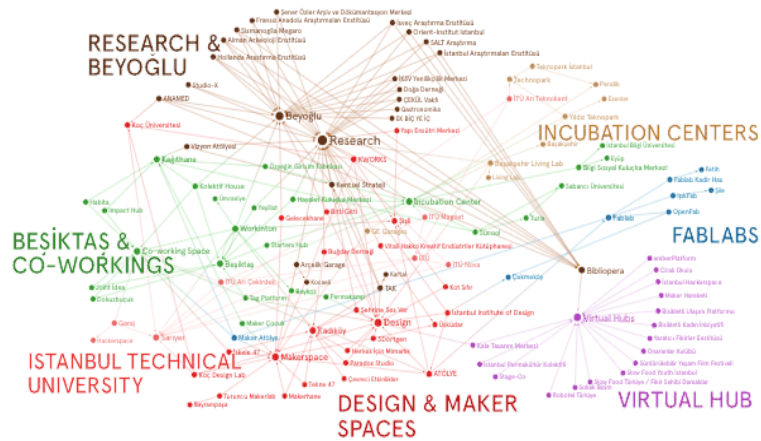


Unsupervised Learning

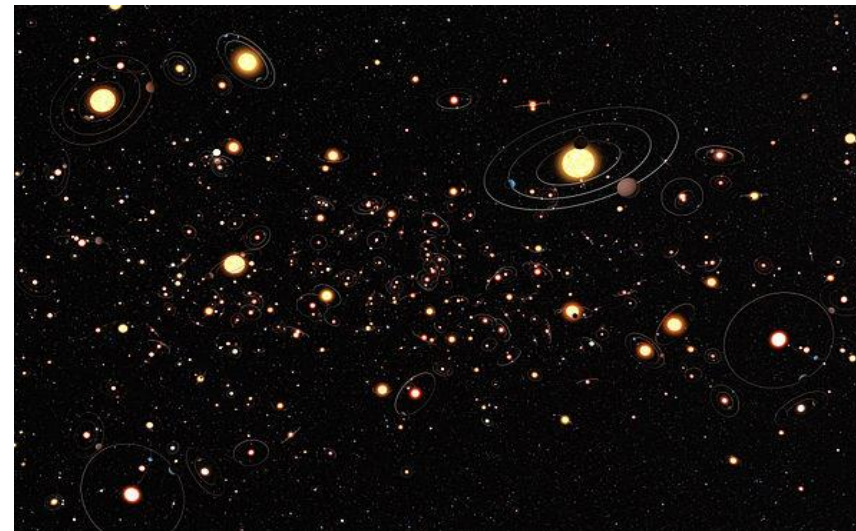
- Given x_1, \dots, x_n (no labels), output hidden structure in x 's
 - E.g., clustering



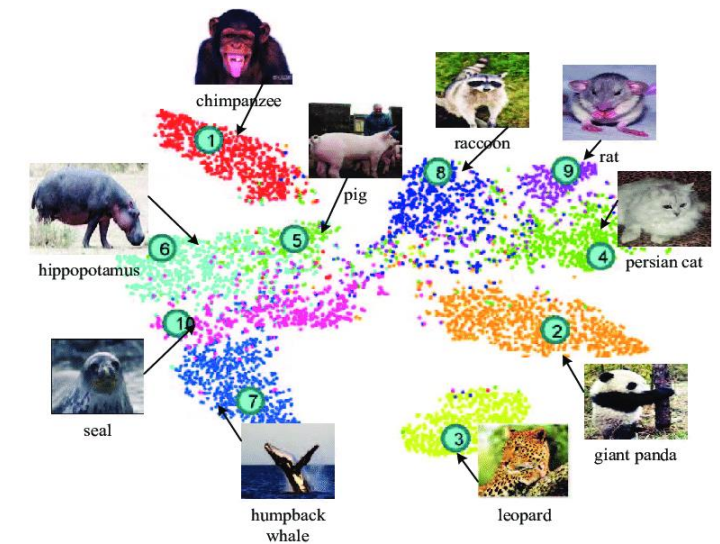
Unsupervised Learning



Find Subgroups in Social Networks



Identify Types of Exoplanets



Visualize Data

Image Credits:

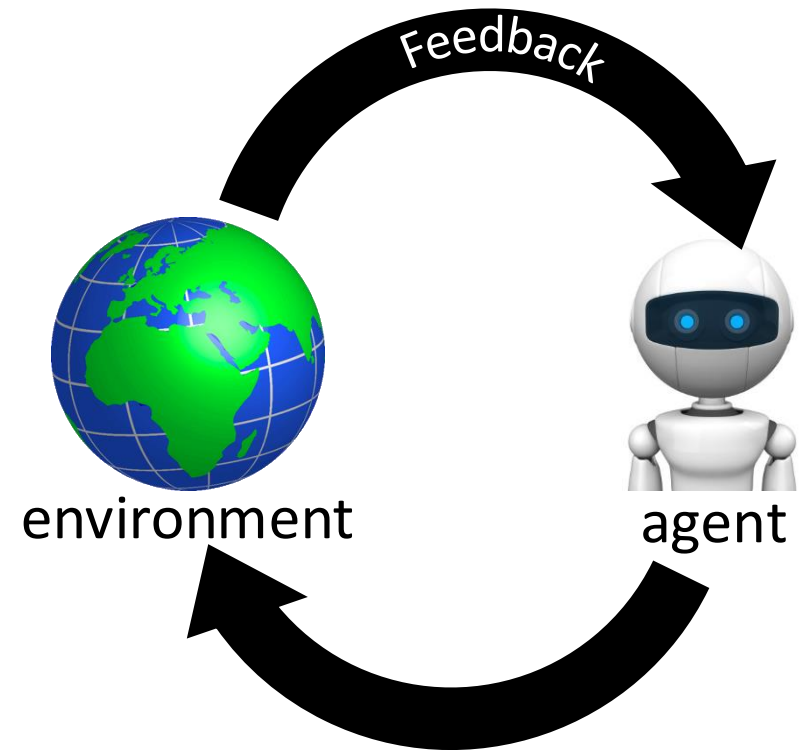
<https://medium.com/graph-commons/finding-organic-clusters-in-your-complex-data-networks-5c27e1d4645d>

<https://arxiv.org/pdf/1703.08893.pdf>

<https://en.wikipedia.org/wiki/Exoplanet>

Reinforcement Learning

- Learn how to perform a task from interactions with the **environment**
- **Examples:**
 - Playing chess (interact with the game)
 - Robot grasping an object (interact with the object/real world)
 - Optimize inventory allocations (interact with the inventory system)



Reinforcement Learning



<https://www.youtube.com/watch?v=iaF43Ze1oel>

When should we use machine learning ...?

... over traditional programming?

Analytical Modeling/ Understanding	Flying rockets to other planets				
	Checking large prime numbers	NO		Solving differential equations	YES, SOMETIMES
			Weather forecasting		
				Recognizing animals from pictures	YES!
	Predict fashion in 20 years	NO, PROBABLY		Make art and music	YES!
				Get robots to make sandwiches	YES, PROBABLY
					Data Quantity and Quality

Applications of Machine Learning

Everyday Applications

COVID-19 PAYMENT ∩ Spam ×



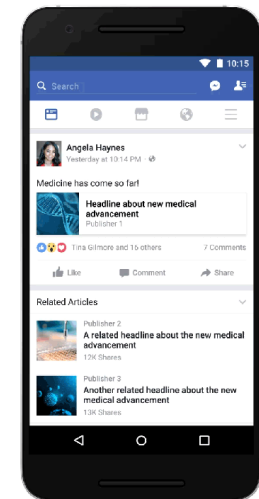
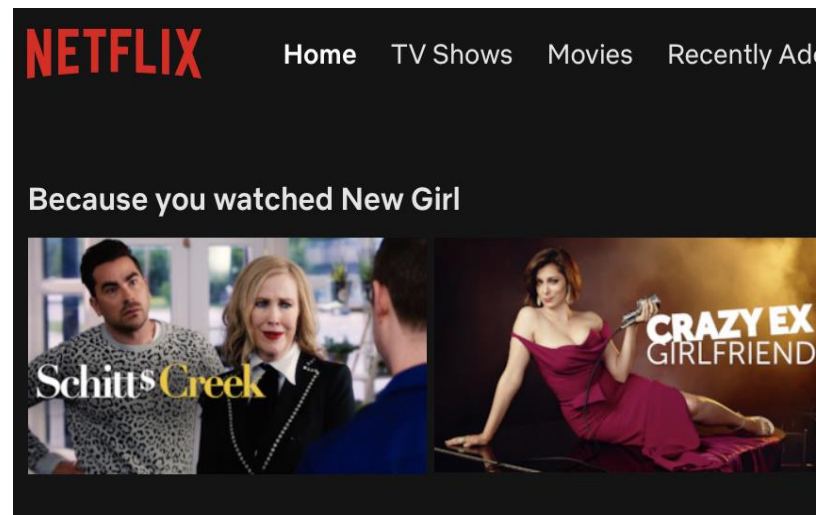
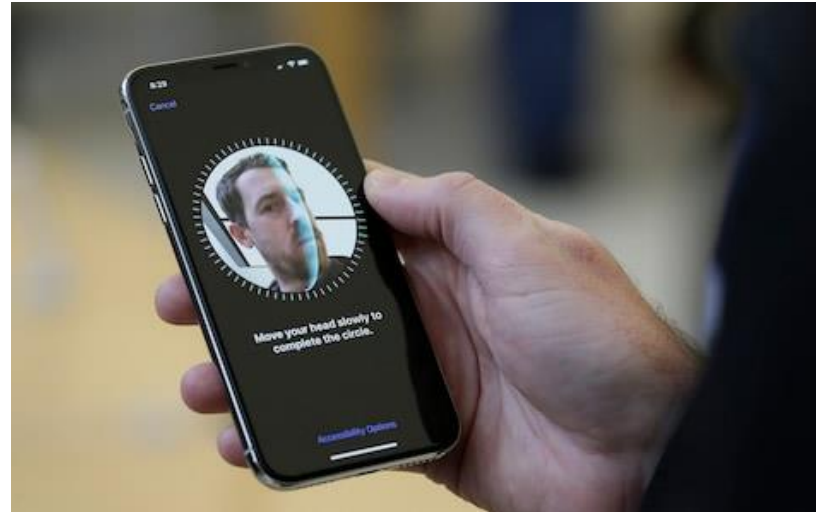
Miller, Jane
to me ▾



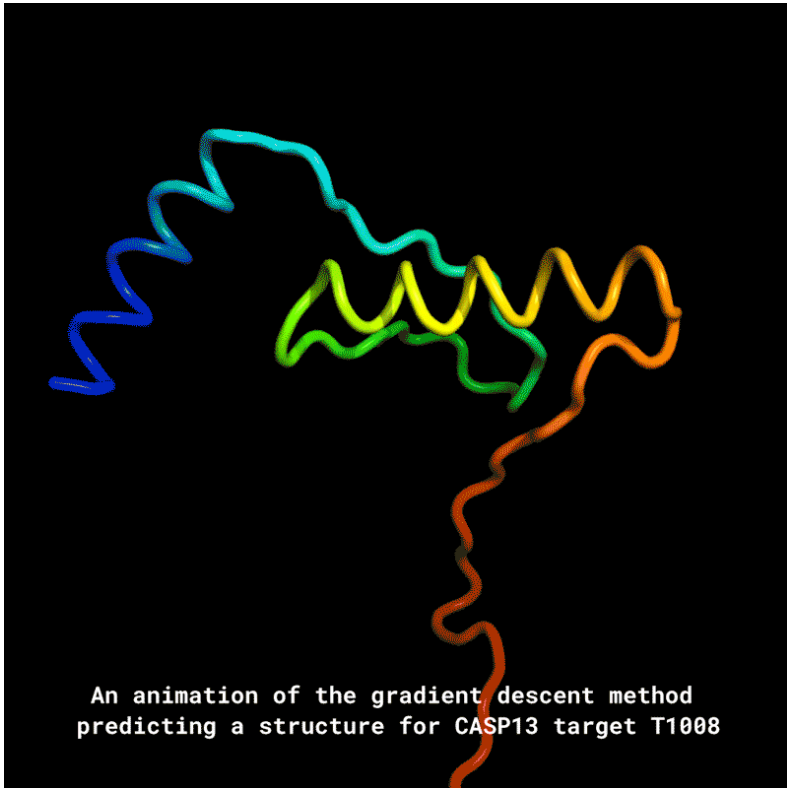
This message seems dangerous

It contains a suspicious link that was used to steal people's personal information. Avoid personal information.

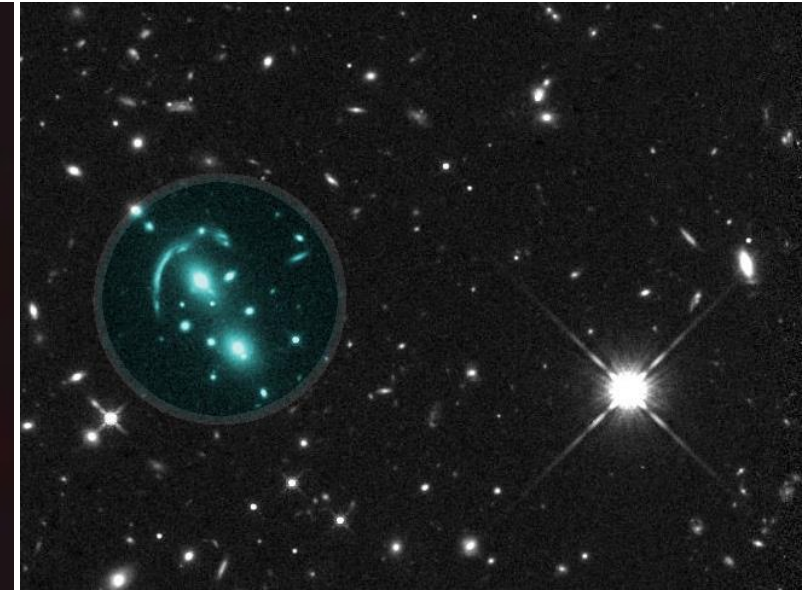
Good morning,
You are advised to download the attached invoice for your review. Please get back to us as soon as possible.
Thanks,
Jane



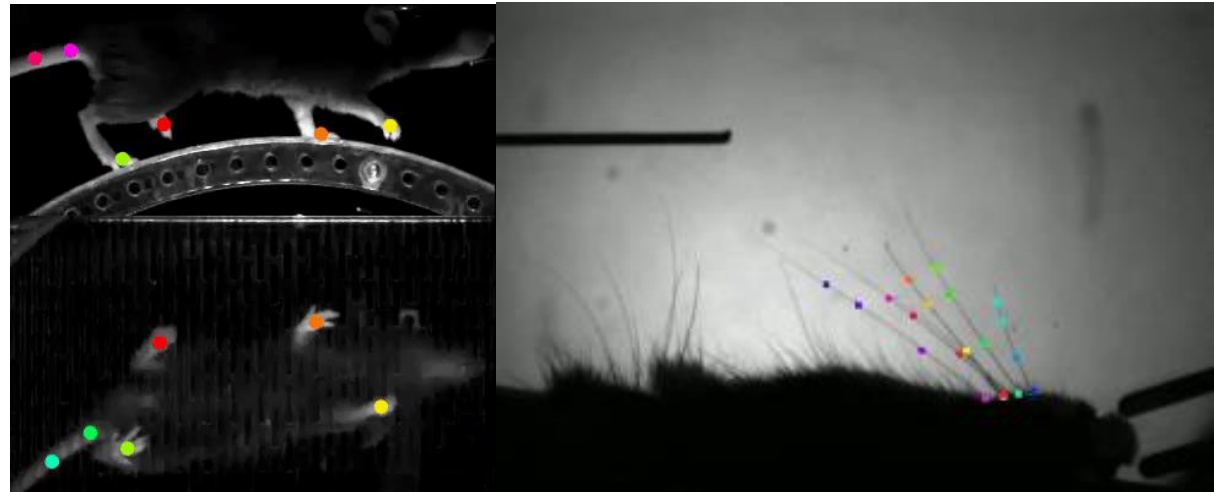
Scientific Discovery



<https://deepmind.com/blog/article/AlphaFold-Using-AI-for-scientific-discovery>



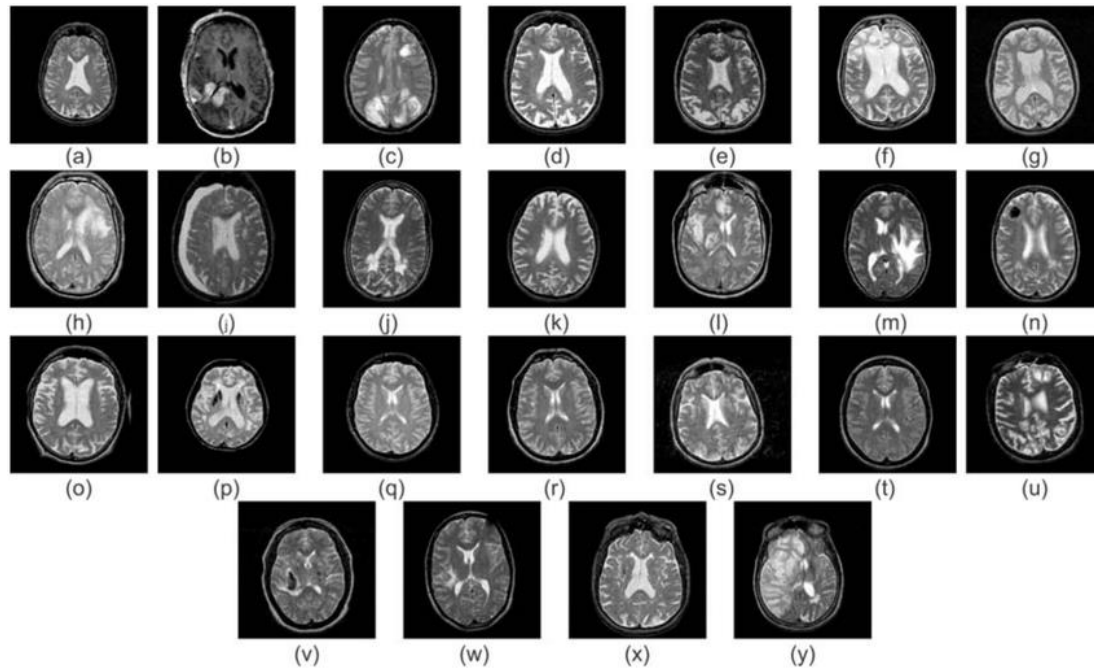
<https://www.jpl.nasa.gov/edu/news/2019/4/19/how-scientists-captured-the-first-image-of-a-black-hole/>



<http://www.mousemotorlab.org/deeplabcut>

Radiology and Medicine

Input: Brain scans




Output: Neurological disease labels

Machine learning studies on major brain diseases: 5-year trends of 2014–2018

Applications of machine learning in drug discovery and development

<https://www.nature.com/articles/s41573-019-0024-5>

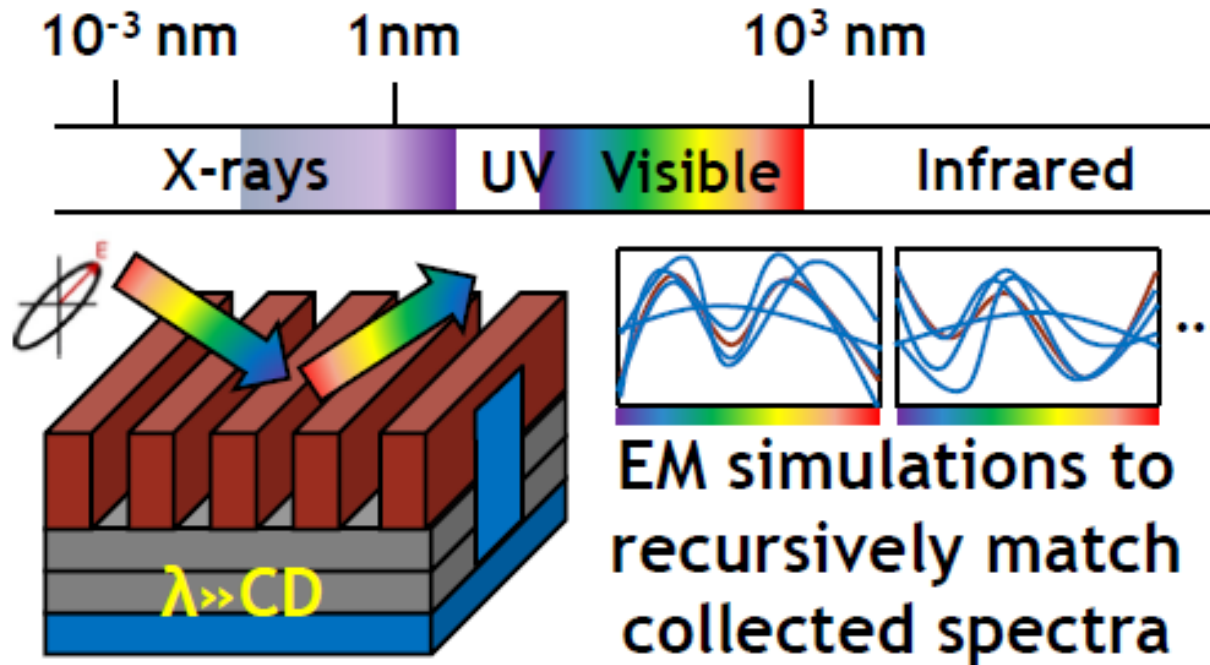
Deep learning-enabled medical computer vision

Andre Esteva , Katherine Chou, Serena Yeung, Nikhil Naik, Ali Madani, Ali Mottaghi, Yun Liu, Eric Topol, Jeff Dean & Richard Socher

<https://www.nature.com/articles/s41746-020-00376-2>

Semiconductor Manufacturing

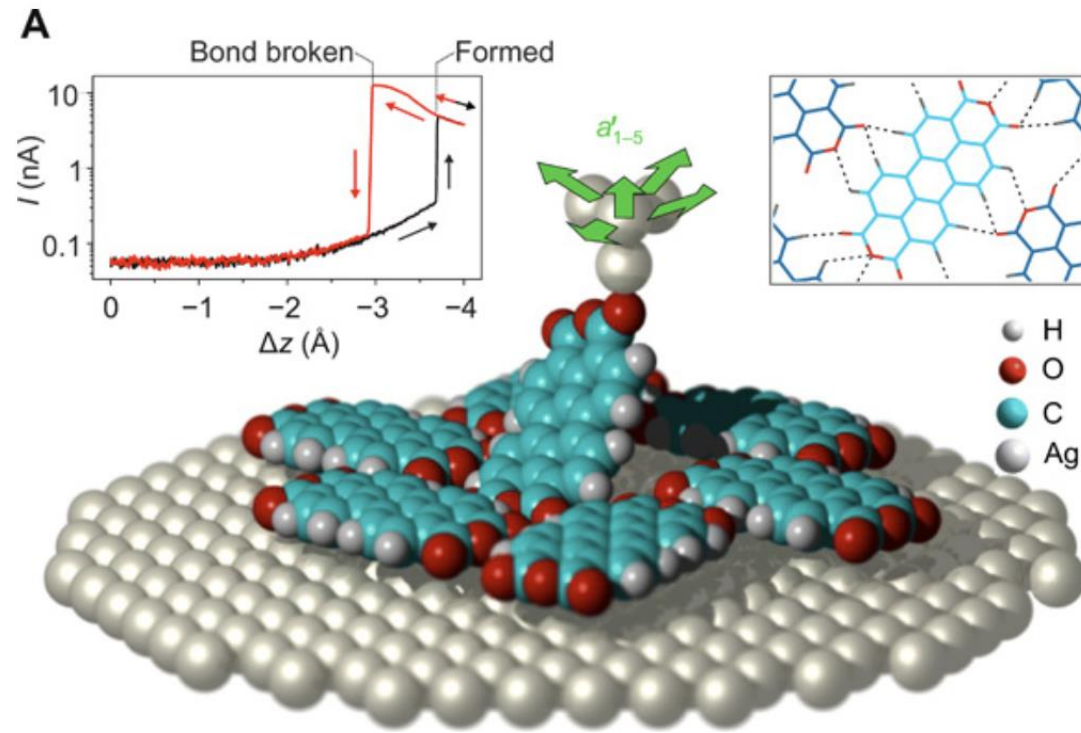
Input: light spectra after bouncing off silicon wafer



huge gains in manufacturing
throughput
→ cheaper phones and computers!

Output: defective / perfect

Nanofabrication



Removing individual molecules autonomously with a scanning probe microscope to manufacture new material

RL agent does better than humans!

Autonomous robotic nanofabrication with reinforcement learning

<https://advances.sciencemag.org/content/6/36/eabb6987>

Creating Music, Art, and Photos



<https://thispersondoesnotexist.com/>



<https://openai.com/blog/jukebox/>

Text Generation

SYSTEM PROMPT (HUMAN-WRITTEN)

Recycling is good for the world.

NO! YOU COULD NOT BE MORE WRONG!!

MODEL COMPLETION (MACHINE-WRITTEN, 25 TRIES)

Recycling is NOT good for the world. It is bad for the environment, it is bad for our health, and it is bad for our economy. I'm not kidding. Recycling is not good for the environment. It is destructive to the earth and it is a major contributor to global warming. Recycling is not good for our health. It contributes to obesity and diseases like heart disease and cancer. Recycling is bad for our economy. It increases the cost of a product, and in turn, the price of everything that is made with that product. Recycling is not good for our nation. We pay a

<https://transformer.huggingface.co/doc/gpt2-large>

Ethical Considerations

“The Pennsylvania Board of Probation and Parole has begun using machine learning forecasts to help inform parole release decisions. In this paper, we evaluate the impact of the forecasts on those decisions and subsequent recidivism.”

An impact assessment of machine learning risk forecasts on parole board decisions and recidivism

[Richard Berk](#) 

“In 2013, the University of Texas at Austin’s computer science department began using a machine-learning system called GRADE to help make decisions about who gets into its Ph.D. program”

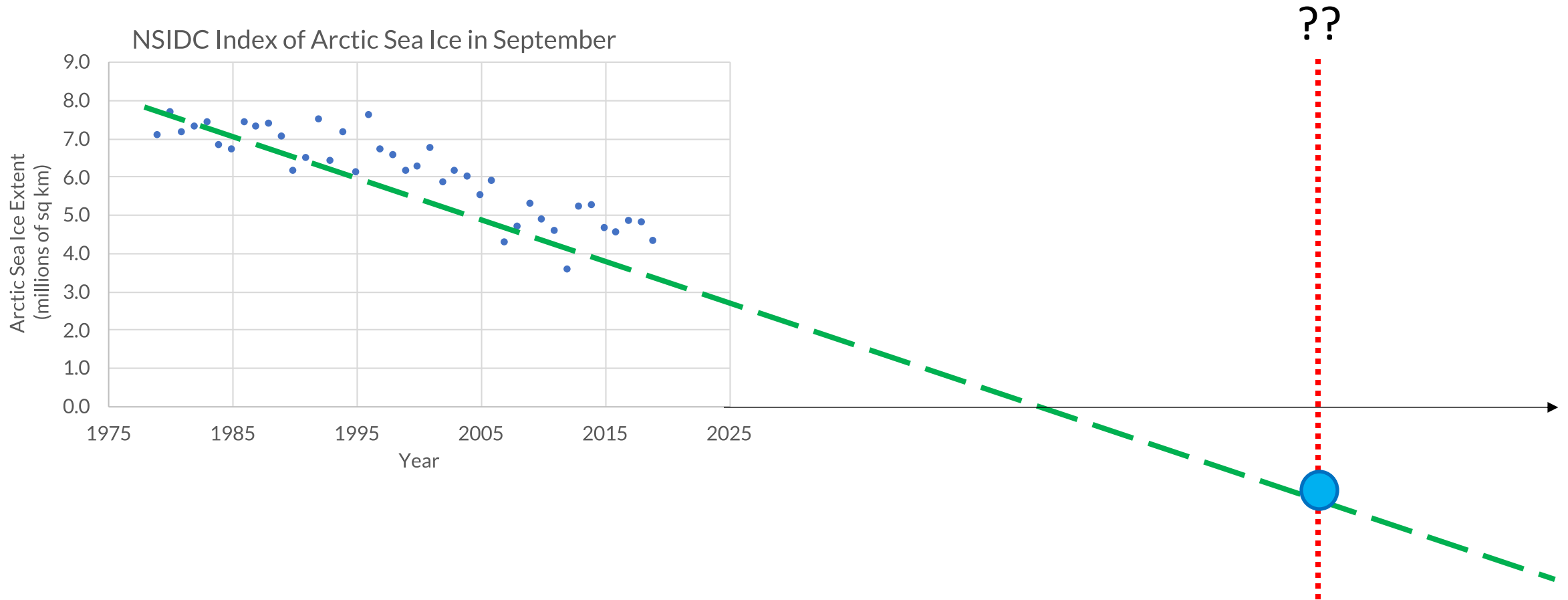
The Death and Life of an Admissions Algorithm

“Videos about vegetarianism led to videos about veganism. Videos about jogging led to videos about running ultramarathons. It seems as if you are never ‘hard core’ enough for YouTube’s recommendation algorithm. It promotes, recommends and disseminates videos in a manner that appears to constantly up the stakes. Given its billion or so users, YouTube may be one of the most powerful radicalizing instruments of the 21st century.”

YouTube, the great radicalizer

THE NEW YORK TIMES / ZEYNEP TUFEKCI / MAR 12

Danger of Out-of-Domain Machine Learning



Any time you are evaluating on data “far” from your training data, beware!