

PH240C: Scope and Introduction

Jingshen Wang

09/01/2021



Class Policy

Homework assignments (65%):

- Biweekly
- The lowest score will be dropped in the final grade, and one late homework (24 hour) is allowed.
- It is encouraged to discuss the problem sets with others, but everyone needs to turn in a unique personal write-up.

Final project write-up and presentation(35%): take home, open books, open notes.

Labs

Please fill in the lab
time change pool!

Supervised Learning (classical approaches)

1. GLM/SVM (09/08)
2. Kernel-based Methods (09/15)
3. Metric Learning (09/22)
4. Tree-based Methods (09/29, 10/06)

3 Labs

Semi-supervised Learning (10/13)

1. Neural Networks (10/20)
2. Deep Neural Networks (10/27)

1 Labs

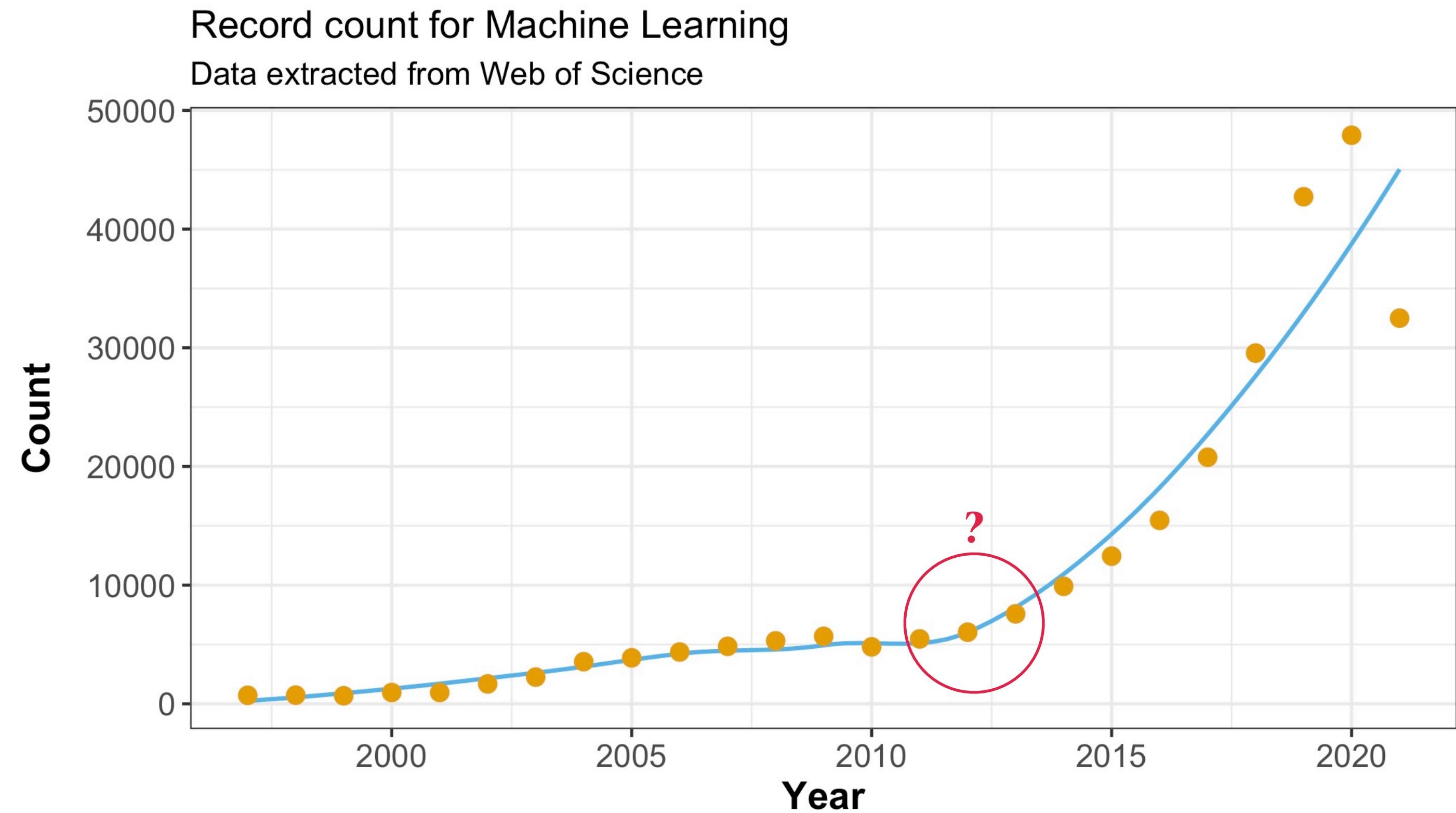
Causal Inference and Clinical Trials

1. Nature's Experiments: Mendelian Randomization (11/10)
2. Bayesian Inference and Design of Experiments (11/17)
3. Adaptive Clinical Trial and Reinforcement Learning (12/01)


2 Labs

What is Machine Learning?

- ▶ The concept of ML is not new
- ▶ Arthur Samuel (1959):
Machine Learning is the field of study that gives the computer the ability to learn without being explicitly programmed



Change Points of ML — Image Recognition

[\[PDF\] Imagenet classification with deep convolutional neural networks](#)

[A Krizhevsky, I Sutskever...](#) - *Advances in neural ...*, 2012 - [proceedings.neurips.cc](#)

We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000 ...

☆  Cited by 85795 [Related articles](#) [All 121 versions](#) 

Showing the best result for this search. [See all results](#)

The trained deep convolutional neural network provides much more accurate prediction than the previous methods

What does it mean for medicine and healthcare?

Studies relied on Deep Learning have showcased its ability to

- ▶ Diagnose some types of skin cancer
- ▶ Identify specific heart-rhythm abnormality like cardiologists
- ▶ Interpret medical scans or pathology slides like highly qualified radiologists
- ▶ Diagnose various eye disease as well as ophthalmologist
- ▶

Nevertheless, machine learning methods need to be powered by big data

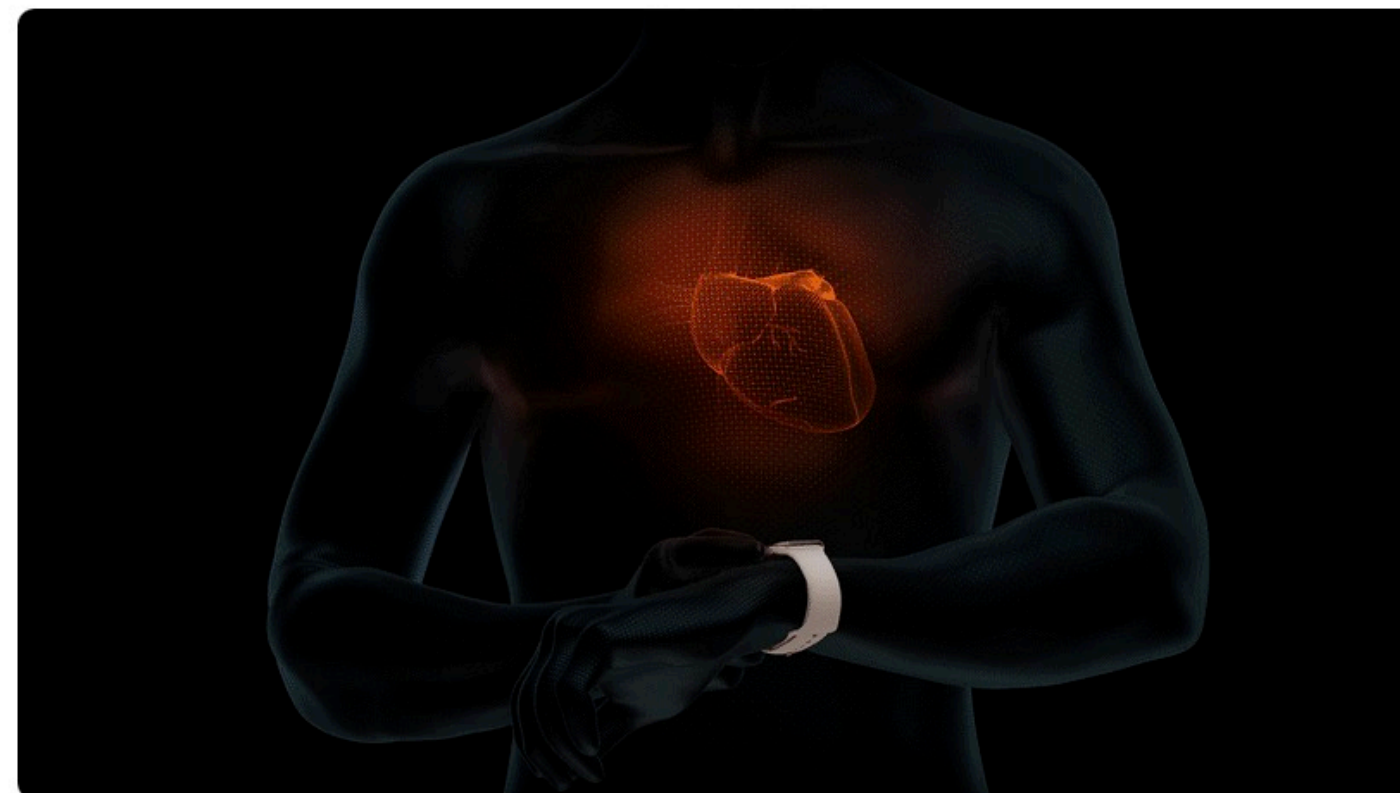
One example

Apple Newsroom needs your permission to [enable desktop notifications](#) when new articles are published

UPDATE
December 6, 2018

ECG app and irregular heart rhythm notification available today on Apple Watch

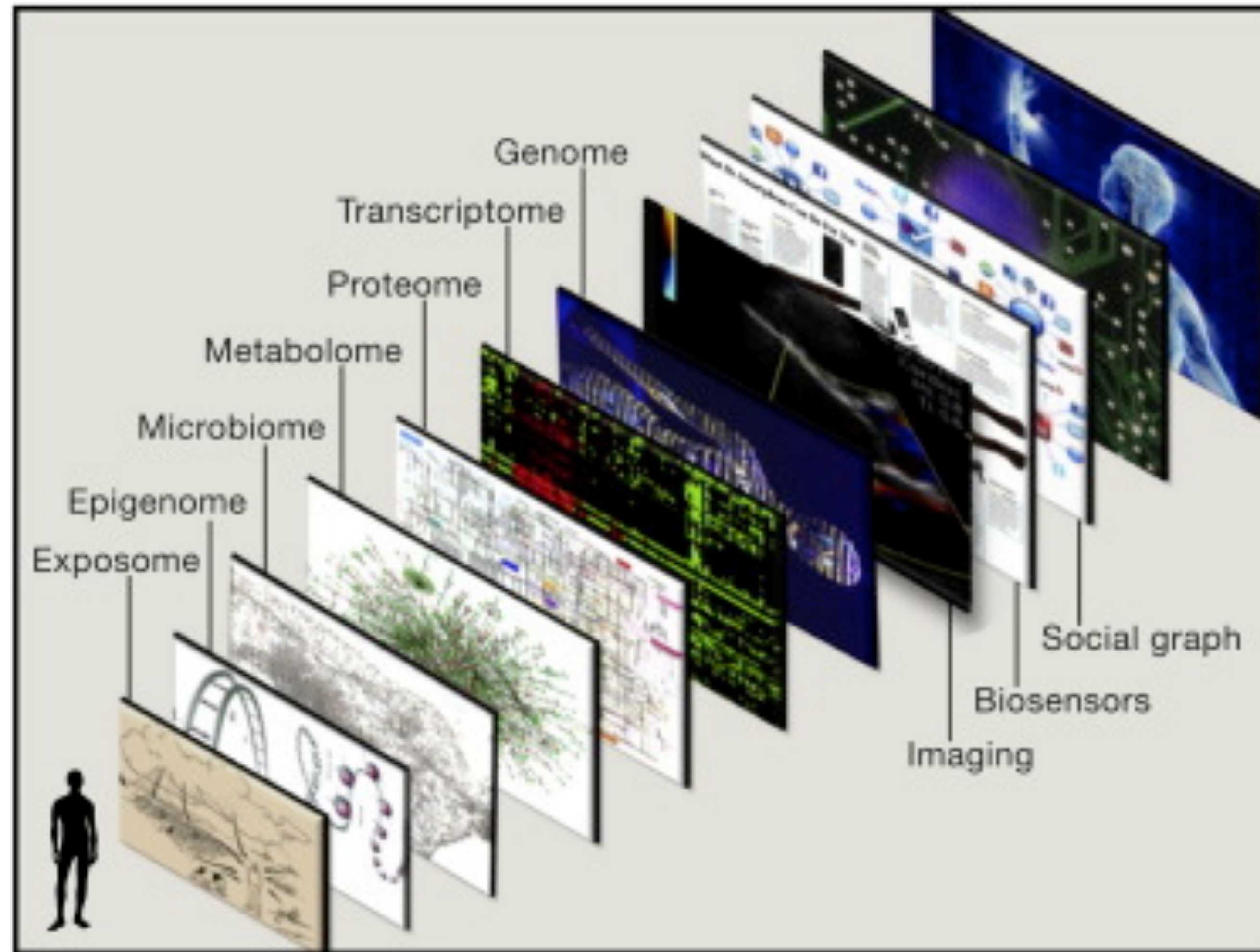
Facebook Twitter Email Link



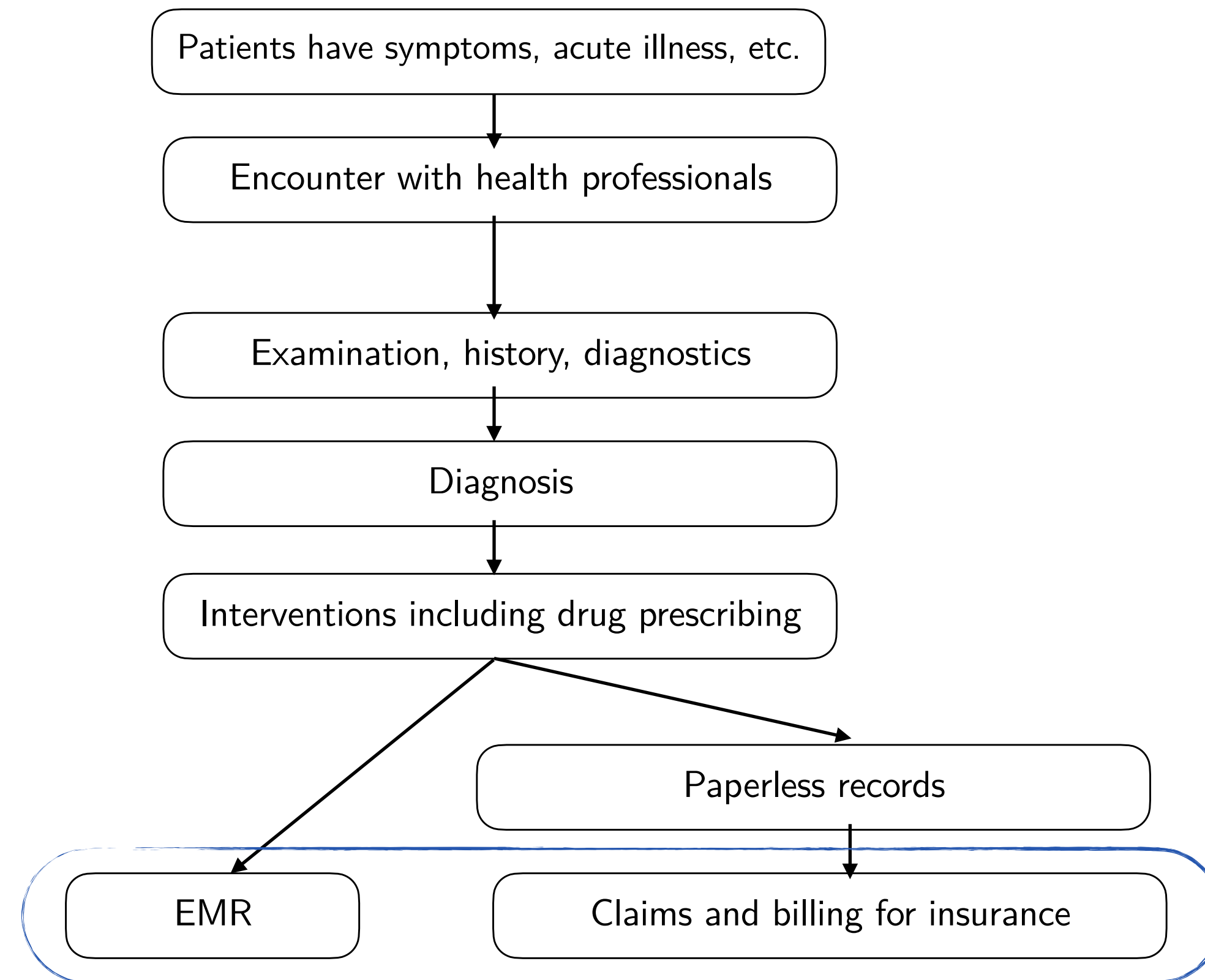
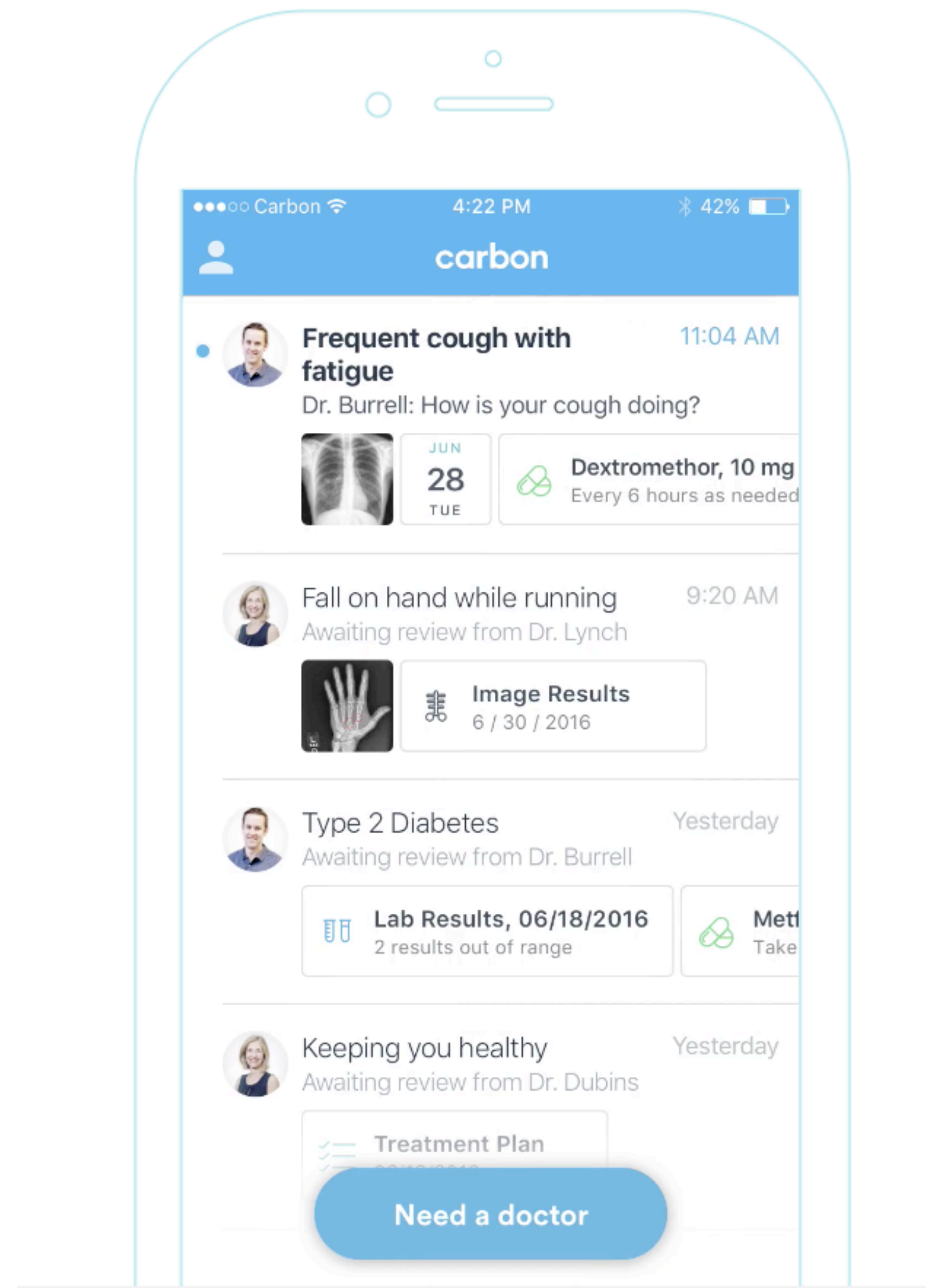
New electrodes in Apple Watch Series 4 now enable customers to take an ECG directly from the wrist.



Geographic Information System of a Human Being



Where do Data Come from?



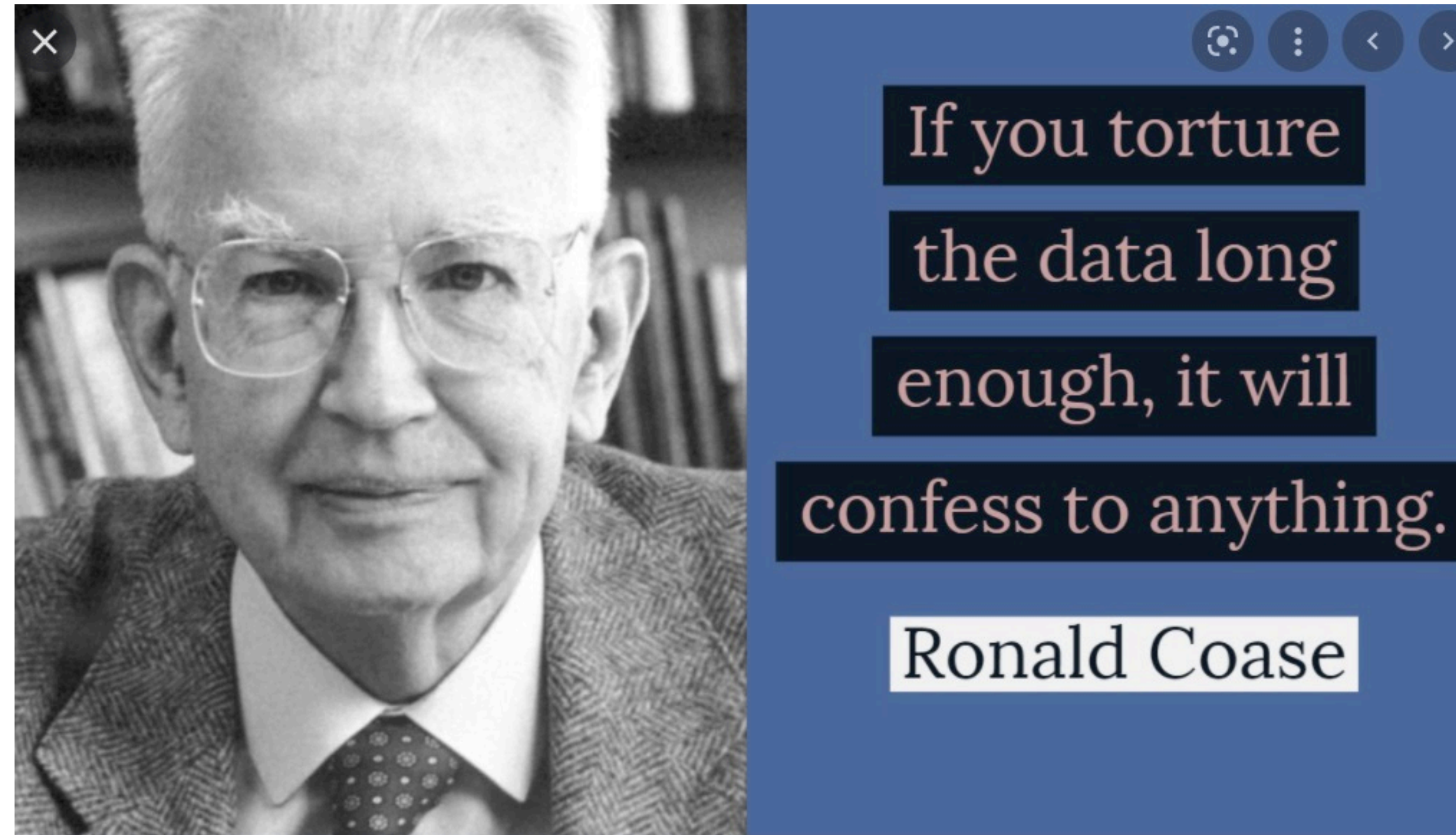
Research Data Base: Electronic Health Record Data

Electronic Health Record Data

EHR stands for electronic health record. According to Wikipedia:

- ▶ EHR is the systematized collection of patient and population electronically-stored health information in a digital format.
- ▶ These records can be shared across different health care settings. Records are shared through network-connected, enterprise-wide information systems or other information networks and exchanges.

What do we want to learn/gain from data?



What do we want to learn/gain from data?

Learning objects in this class:

1. Disease risk (early) prediction/diagnosis — Supervised Learning
2. Utilize massive undiagnosed patient information to improve prediction — Semi-supervised Learning
3. Guide future clinical decisions — causal inference and clinical trial design

We need tools to achieve these goals

Supervised Learning (classical approaches)

1. GLM/SVM
2. Kernel-based Methods
3. Metric Learning
4. Tree-based Methods

Example 1. The prime example: house price prediction

Suppose we have data about

- (1) Features X — square footage, number of rooms, features, whether a house has a garden or not
- (2) Labels/Outcomes Y — the prices of these houses

By leveraging data coming from thousands of houses, we can train a supervised machine learning model to predict a new house's price based on the examples observed by the model.

Example: Text as data

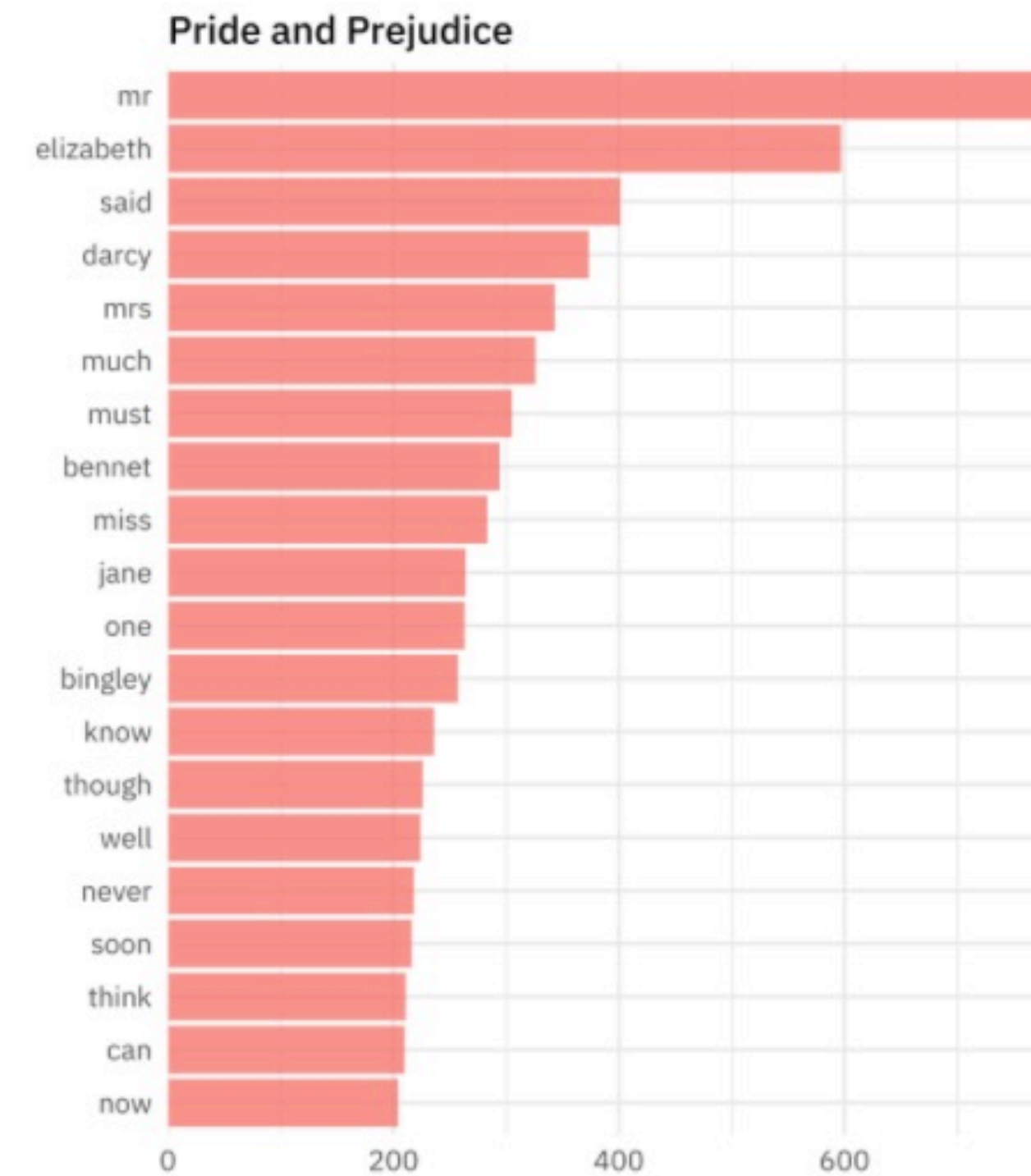
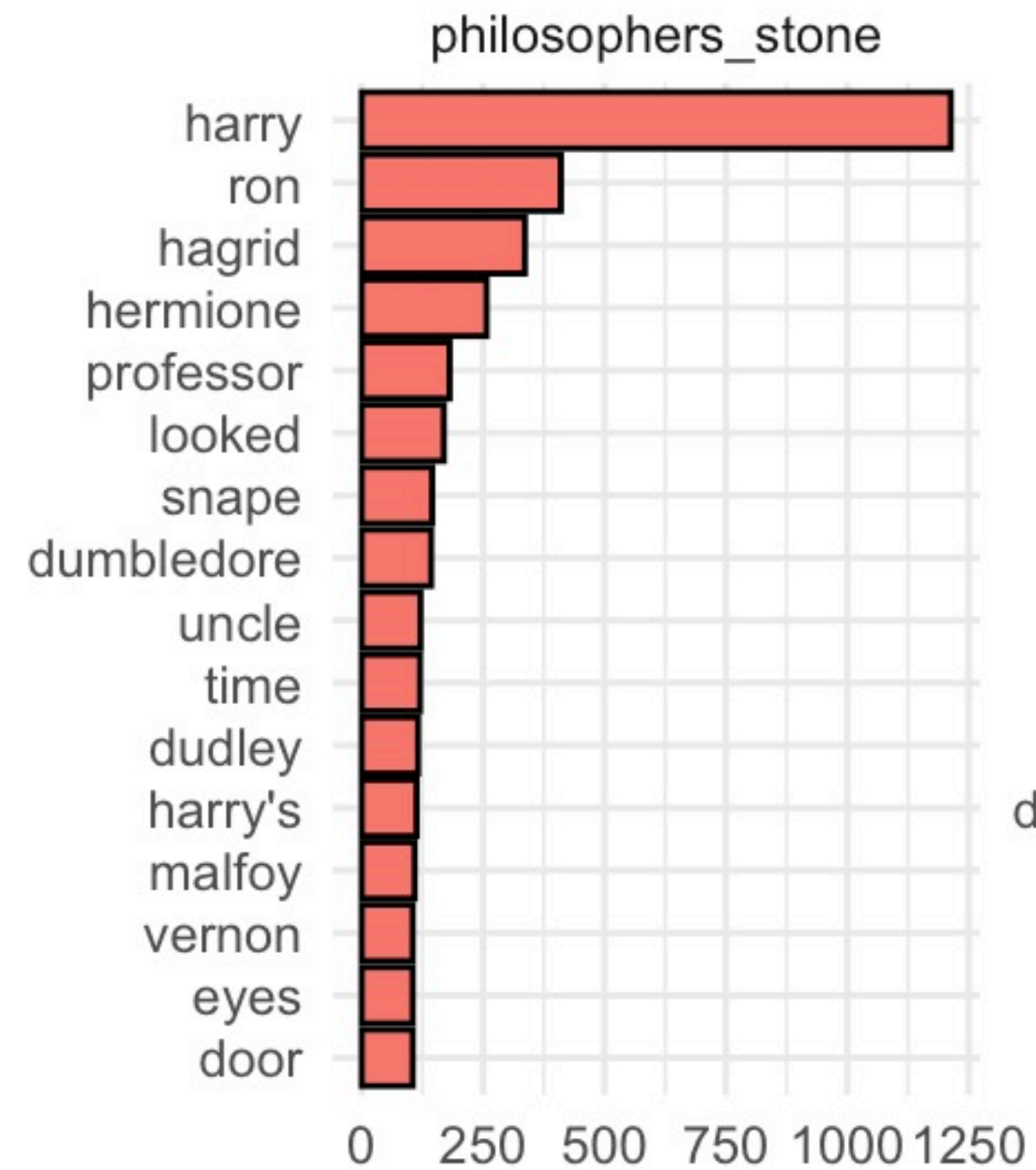
Example 2. Text as data

Can you tell who wrote the following sentences? Jane Austen from *Pride and Prejudice* or J. K. Rowling from *Harry Potter*?

“The thought of the confined creature was so dreadful to him...”

“...was as if something turned over, and the point of view altered...”

Examples: Text as data (predictors)



You will work on these text data in your labs with your GSI, and make predictions on your own

Example: Disease risk prediction

Example 3. Disease risk prediction

We are given massive biobank data for patients with stroke history that contain information:

1. $Y \in \{0,1\}$: if the patient has developed Alzheimer's disease at the time of recruit
2. $X \in \mathbb{R}^p$: covariates information including individual lifestyles (insomnia, current smoking status, beef lover, etc.), baseline biomarker information (gender, age, education, and family AD history), and genetic information (SNPs)

Goal:

1. Predict which patients that are at high risk of developing AD (why important?)
2. Find any lifestyle factors can reduce the risk of AD

How about the other patients that are not diagnosed? Can we use their information to improve our prediction?

Semi-supervised Learning

Semi-supervised Learning

In medical record database, we often encounter the following scenario:

1. **Labeled data** $\{Y_i, X_i\}_{i=1}^n$ — $Y_i \in \{0,1\}$ represents if the patient is diagnosed with certain disease at the time of recruit, and $X_i \in \mathbb{R}^p$ represents patient all available information
2. **Unlabeled data** $\{X_j\}_{j=n+1}^{n+N}$ — $X_j \in \mathbb{R}^p$ represents unlabeled patient information

Question: Can we improve the prediction results given massive amount of unlabeled data? If so, when?

Learning objectives in the first half of the semester

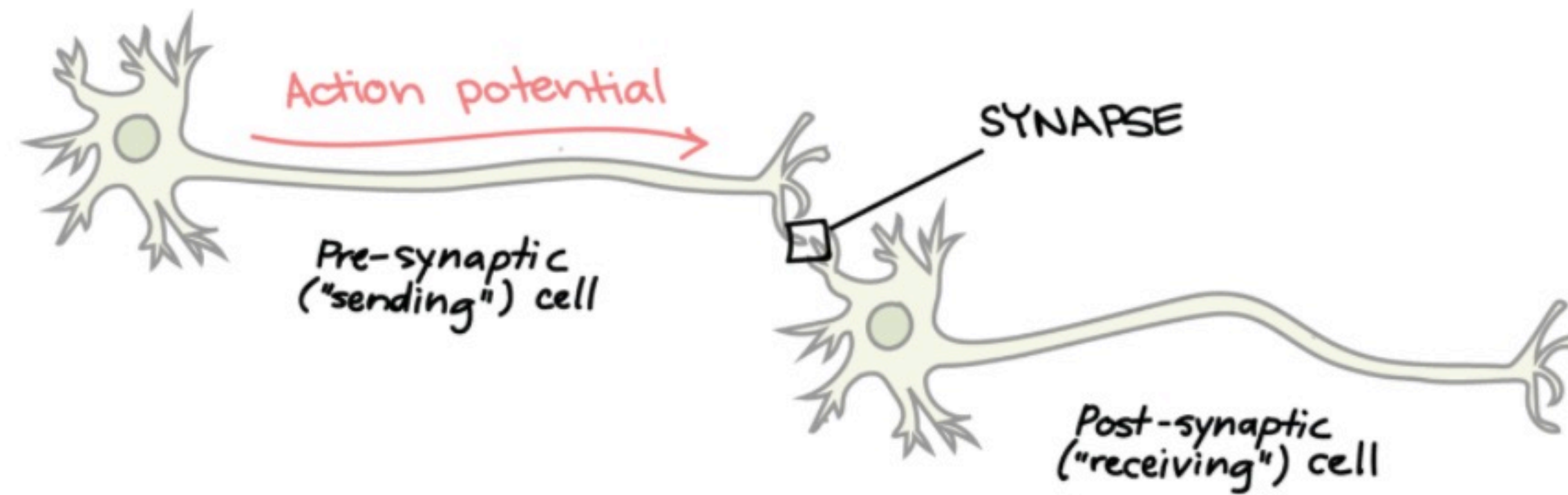
Supervised Learning (classical approaches)

1. GLM/SVM (09/08)
2. Kernel-based Methods (09/15)
3. Metric Learning (09/22)
4. Tree-based Methods (09/29, 10/06)

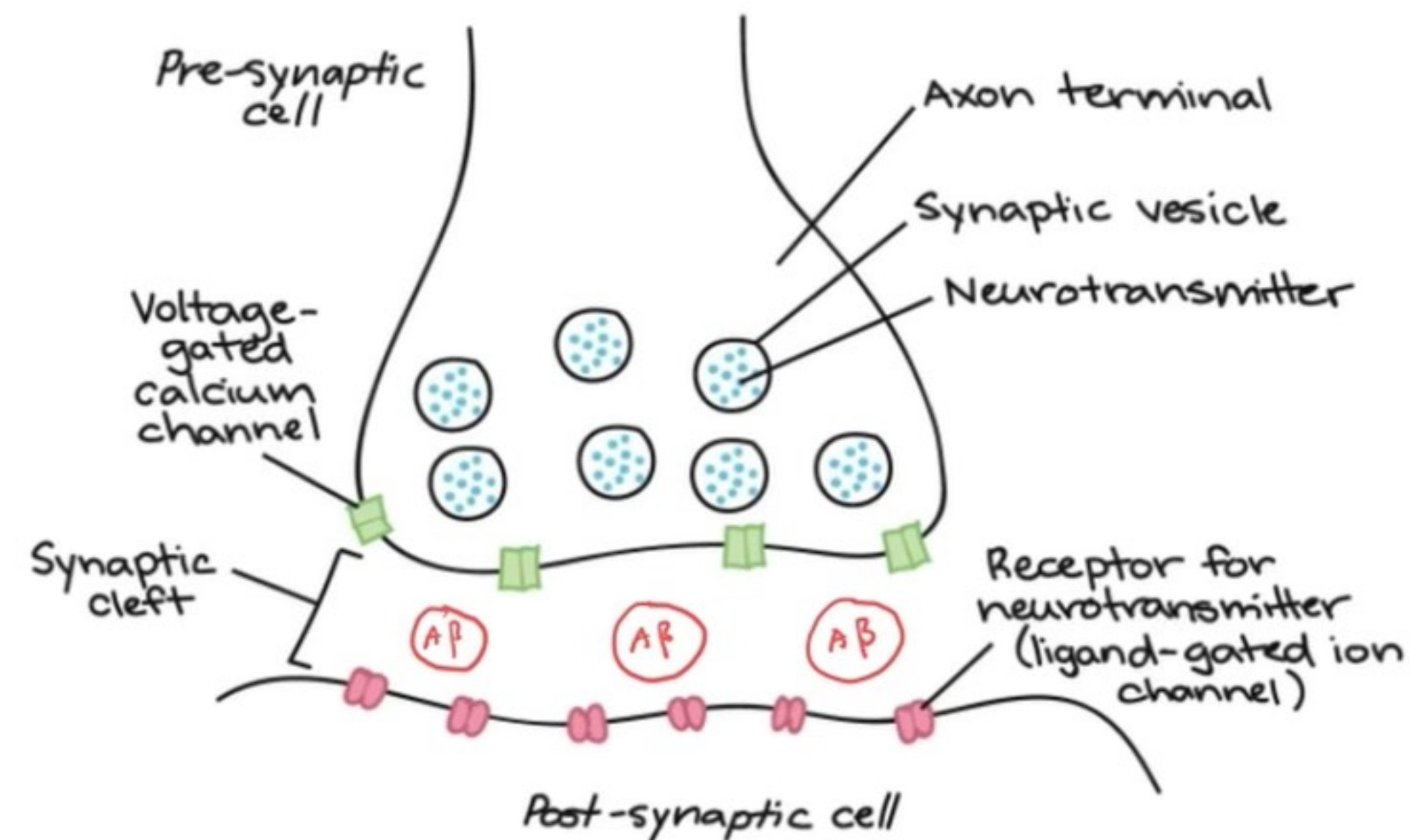
Semi-supervised Learning (10/13)

And then?

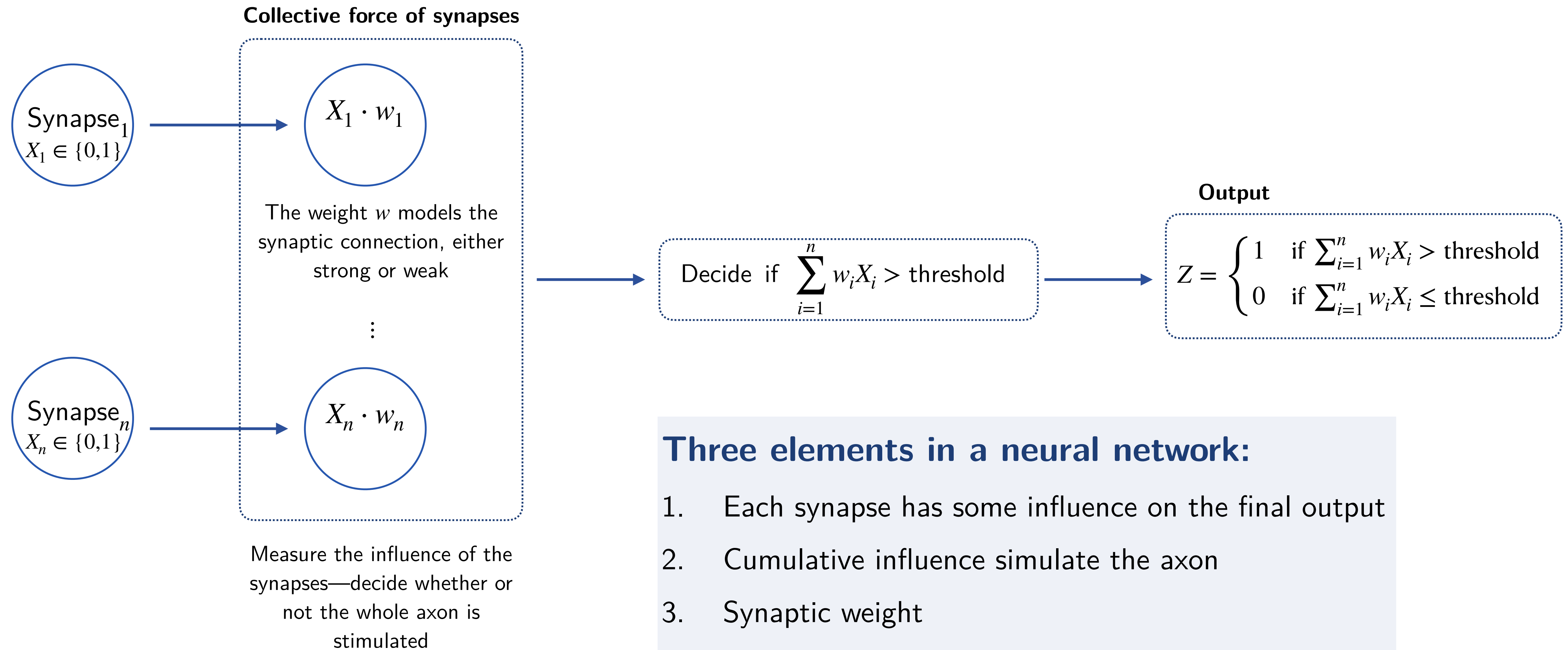
Alzheimer's disease: a little more detail..



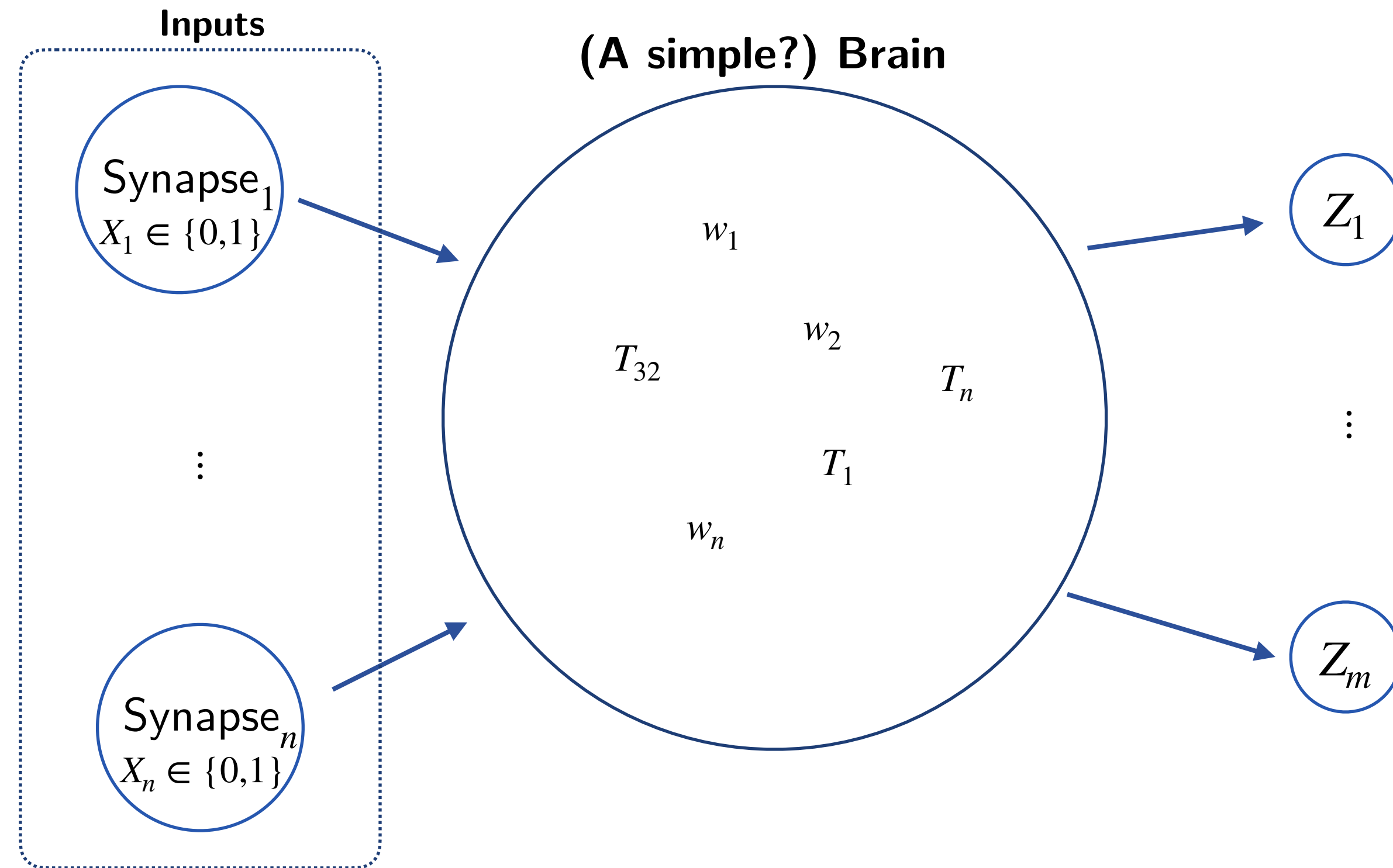
Alzheimer's disease is a neurodegenerative disease often characterized by dementia, accumulation of beta-amyloid ($A\beta$) plaques and tau proteins on neurons, and brain inflammation and atrophy.



Neural Networks (1)



Neural Networks (2)



Neural Network (NN)

A NN learns the functional form of

$$Z_j = f_j(X_1, \dots, X_n, w_1, \dots, w_n, T_1, \dots, T_n), \quad j = 1, \dots, m,$$

where we need to adjust the weights w_i and the thresholds T_i so that what we get out is the outcome Z_j .

Example: Disease risk prediction with NN

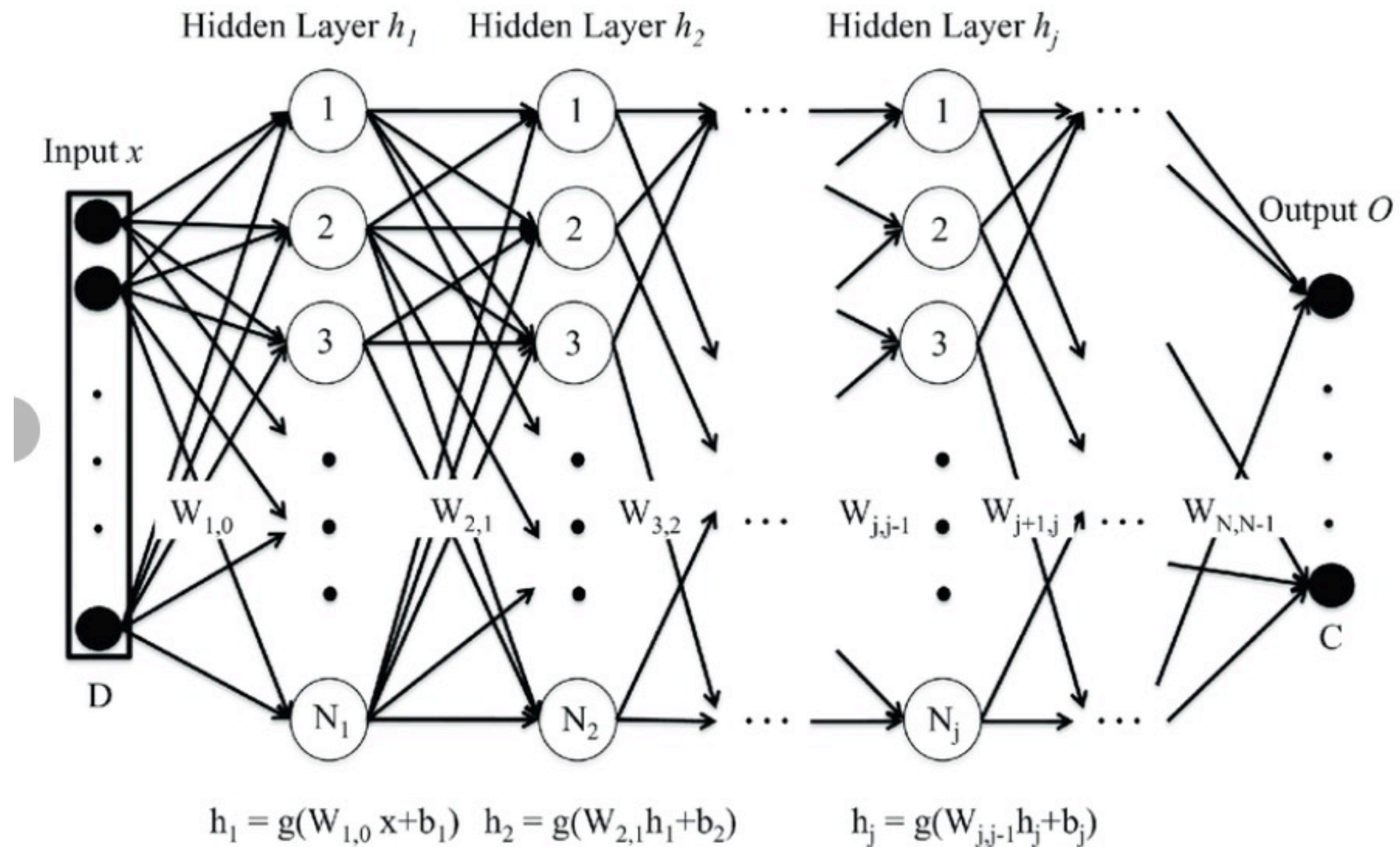
Example 3. Disease risk prediction (revisit)

We are given massive biobank data for patients with stroke history that contain information:

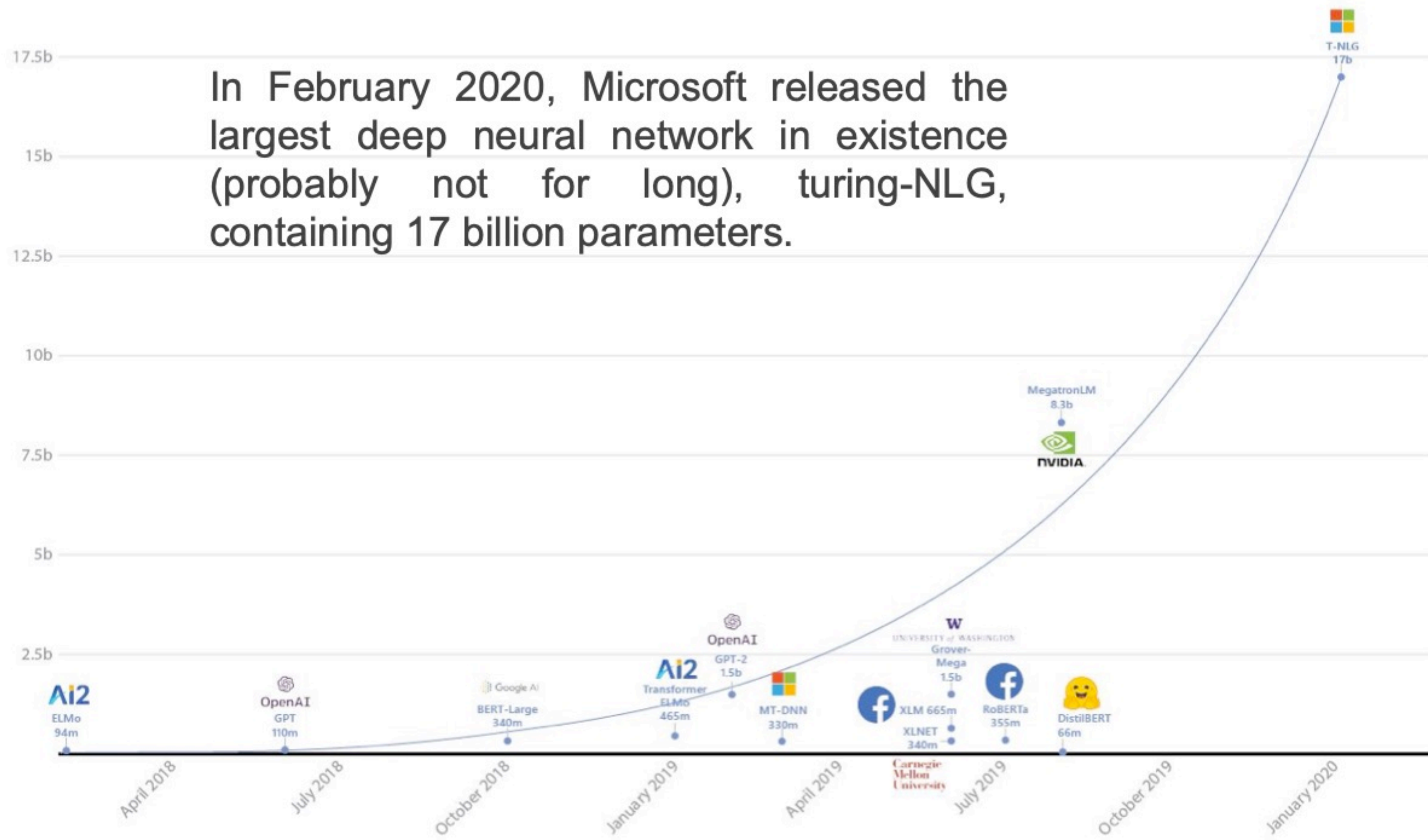
1. $Y \in \{0,1\}$: if the patient has developed Alzheimer's disease at the time of recruit
2. $X \in \mathbb{R}^p$: covariates information including individual lifestyles (insomnia, current smoking status, beef lover, etc.), baseline biomarker information (gender, age, education, and family AD history), and genetic information (SNPs)

To have binary input, we can transform the covariates into dummy variables. Can we still obtain the effect of certain lifestyle on lowering the disease risk?

Deep Neural Networks



Deep Neural Networks



Challenges in Research

IBM Watson for oncology

EDITORS' PICK | Feb 19, 2017, 03:49pm EST

MD Anderson Benches IBM Watson In Setback For Artificial Intelligence In Medicine



Matthew Herper Former Staff
Healthcare

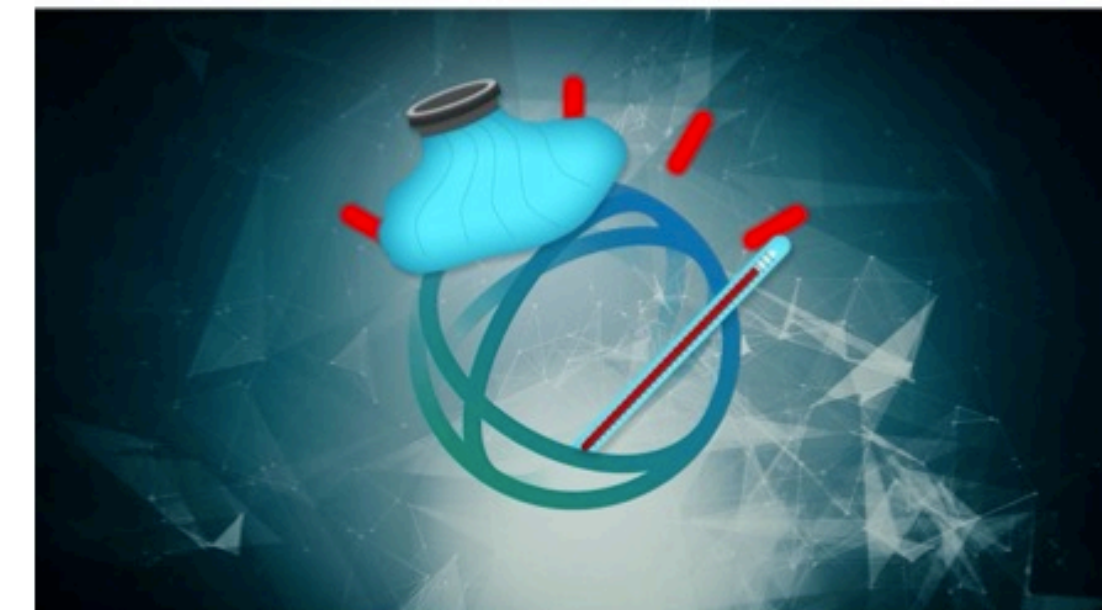
I cover science and medicine, and believe this is biology's century.

1. MD Anderson taps IBM Watson to power "Moon Shots" mission aimed at ending cancer, starting with Leukemia
2. Big data insights to help accelerate translation of cancer-fighting knowledge to cutting edge medical practices
3. Link to the news: <https://www.ibm.com/products/clinical-decision-support-oncology>
4. IBM's Watson supercomputer recommended 'unsafe and incorrect' cancer treatments, internal documents show

How to make machine learning methods more trustworthy?

IBM's Watson supercomputer recommended 'unsafe and incorrect' cancer treatments, internal documents show

By Casey Ross @caseymross and Ike Swetlitz • July 25, 2018



Alex Hogan/STAT

Internal IBM documents show that its Watson supercomputer often spit out erroneous cancer treatment advice and that company medical specialists and customers identified "multiple examples of unsafe and incorrect treatment recommendations" as IBM was promoting the product to hospitals and physicians around the world.

Learning objectives

1. Neural Networks (10/20)

Convolutional Networks, Recurrent Networks, and their applications in medical research

2. Deep Neural Networks (10/27)

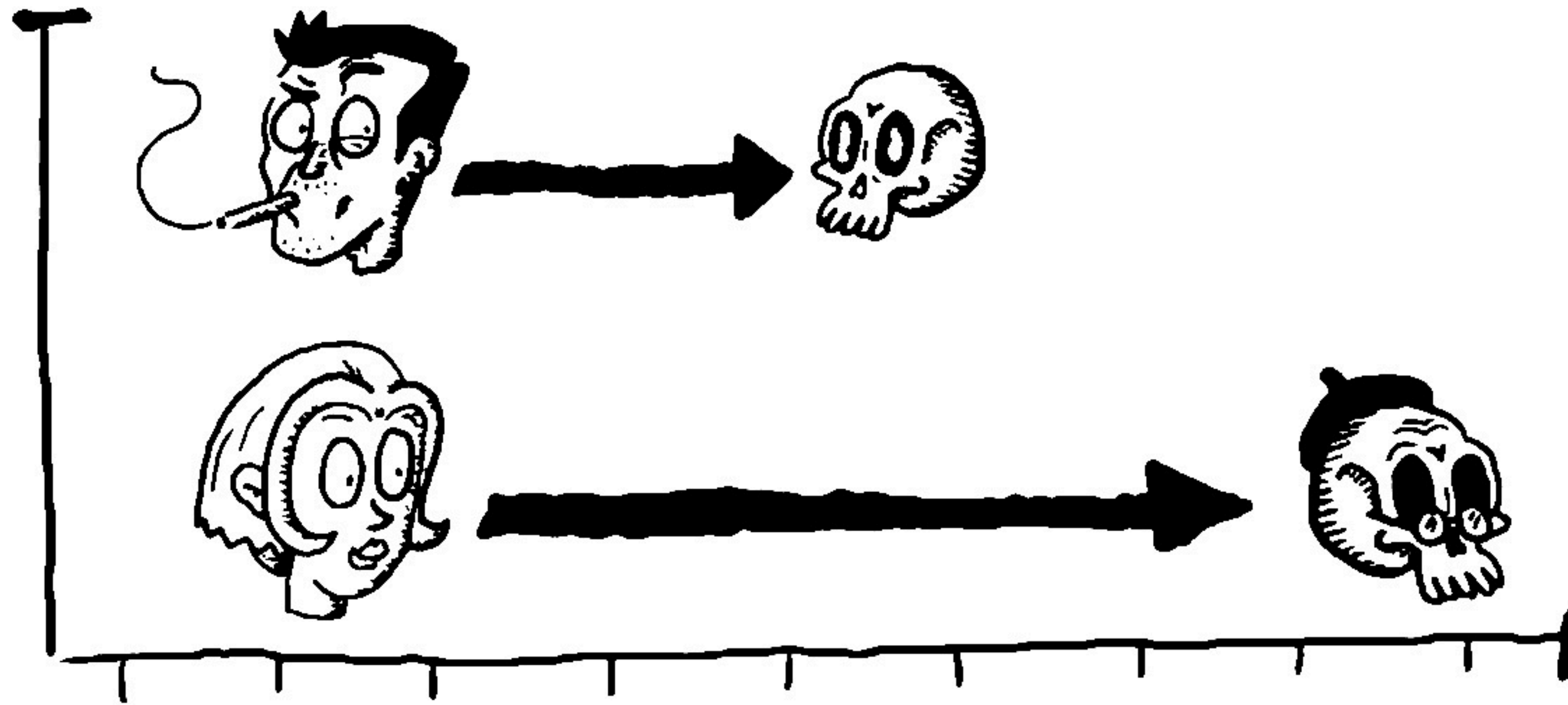
Observational data

1. Disease risk prediction — supervised learning (including NN and DNN), semi-supervised learning
2. Learning effective treatments for different patients

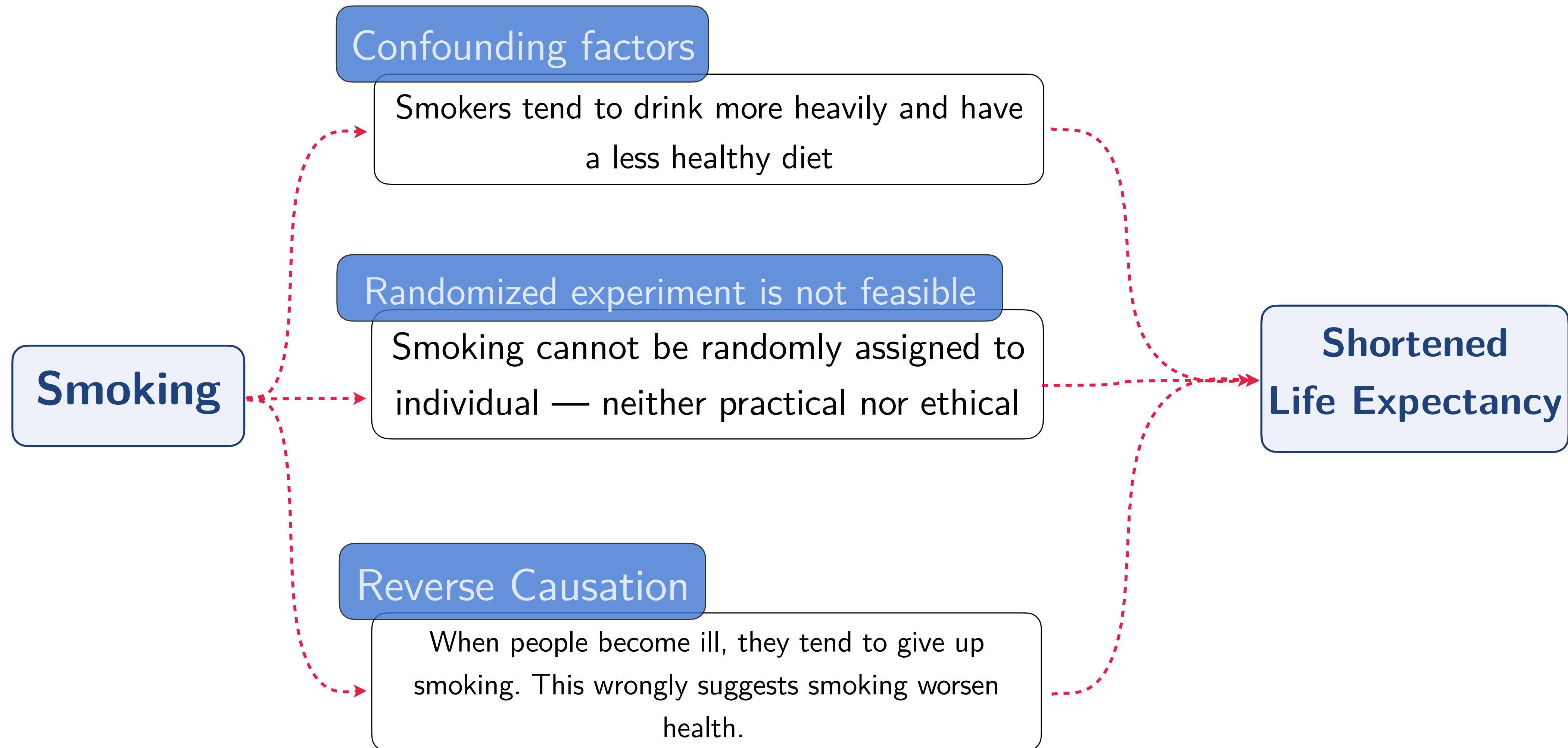
?

1. Will you trust the findings from observational data?
2. If not, what can we do?
3. How can we guide future clinical decisions?

Effect of Smoking on Life Expectancy?



What hinders our understanding?



Last learning objectives

Causal Inference and Clinical Trials

1. Nature's Experiments: Mendelian Randomization (11/10)
2. Bayesian Inference and Design of Experiments (11/17)
3. Adaptive Clinical Trial and Reinforcement Learning (12/01)

Observational data

1. Disease risk prediction — supervised learning (including NN and DNN), semi-supervised learning
2. Learning effective treatments for different patients

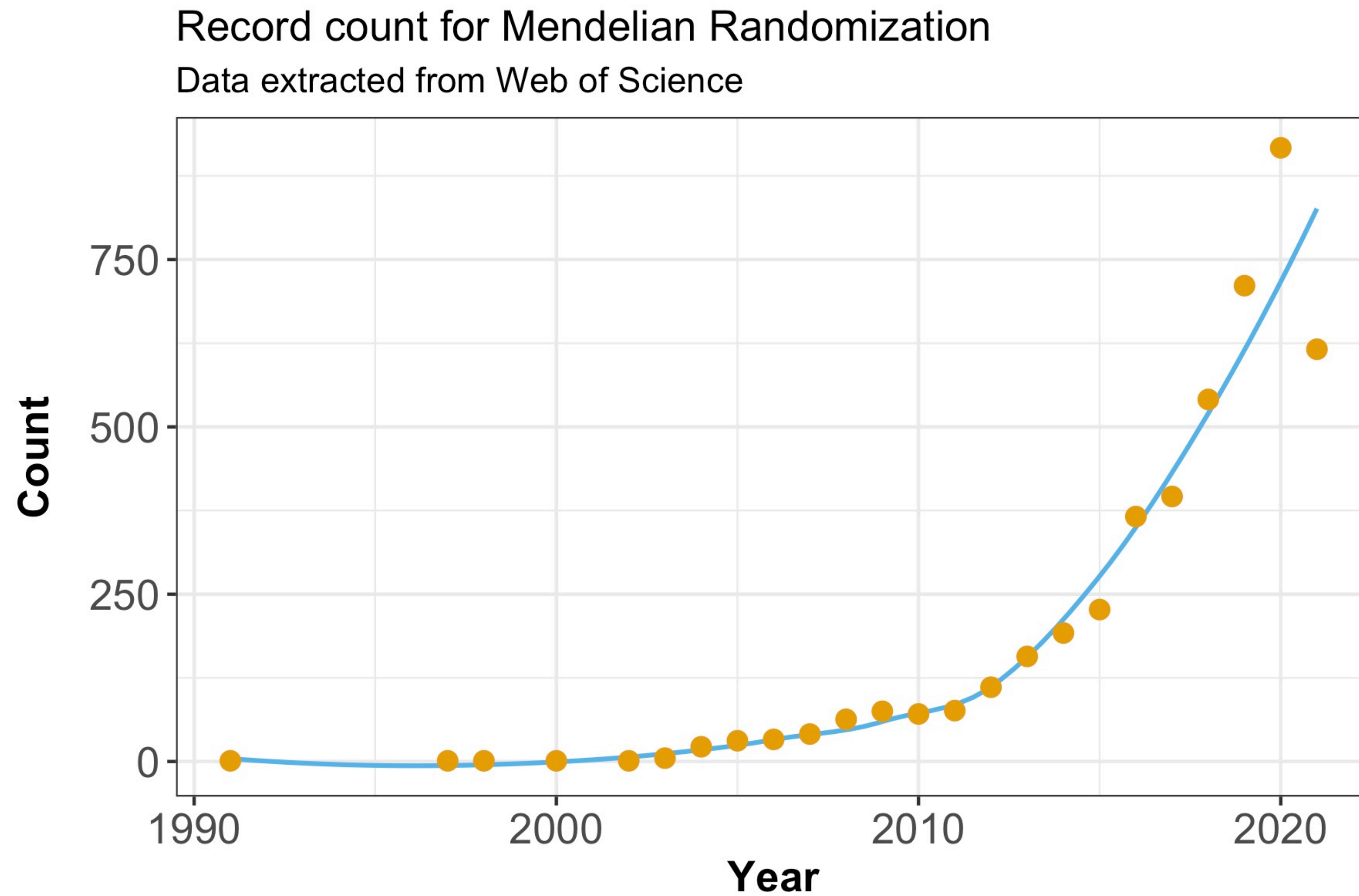


Causal Inference

1. Randomized control trials (RCT) can help us to verify our finds from observation study
2. When RCT is not available, we may rely on Nature's experiment (Mendelian Randomization)

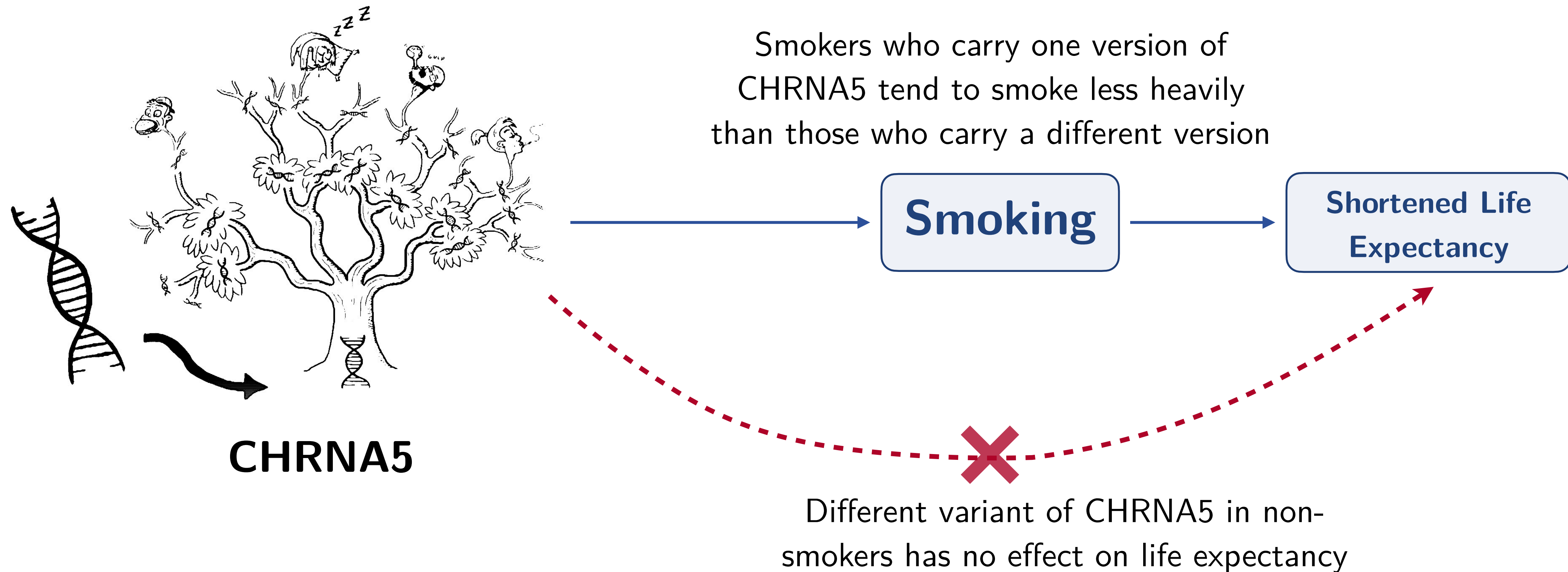
Reinforcement

Mendelian Randomization: Popularity



Mendelian Randomization

Randomly inherited genes are **not** associated with any confounding factors



Labs

Please fill in the lab
time change pool!

Supervised Learning (classical approaches)

1. GLM/SVM (09/08)
2. Kernel-based Methods (09/15)
3. Metric Learning (09/22)
4. Tree-based Methods (09/29, 10/06)

3 Labs

Semi-supervised Learning (10/13)

1. Neural Networks (10/20)
2. Deep Neural Networks (10/27)

1 Labs

Causal Inference and Clinical Trials

1. Nature's Experiments: Mendelian Randomization (11/10)
2. Bayesian Inference and Design of Experiments (11/17)
3. Adaptive Clinical Trial and Reinforcement Learning (12/01)

2 Labs