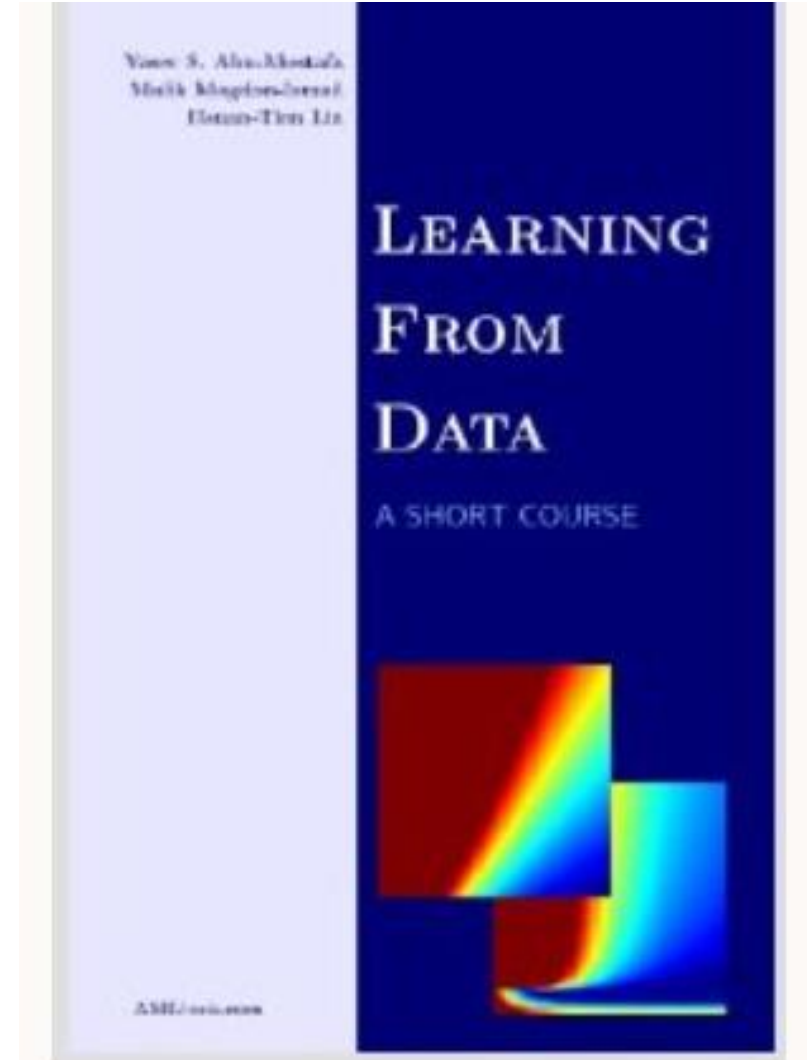# Machine Learning from Data

Lecture 1: Spring 2021

# Resources

- Textbook (Yaser S. Abu-Mostafa, Malik Magdon-Ismail)
- Website
- Homework Submission: Submitty

# Topics Covered in the Course

- What is Learning?
- Can we Do it? ✓
- How to Do it? ← **Models**
- How to do it well? **Regularization**
- General Principles of Learning
- Advanced techniques
- Other Learning Paradigms

# Today's lecture

- Motivation
- Learning Vs. Defining
- Formalize Learning
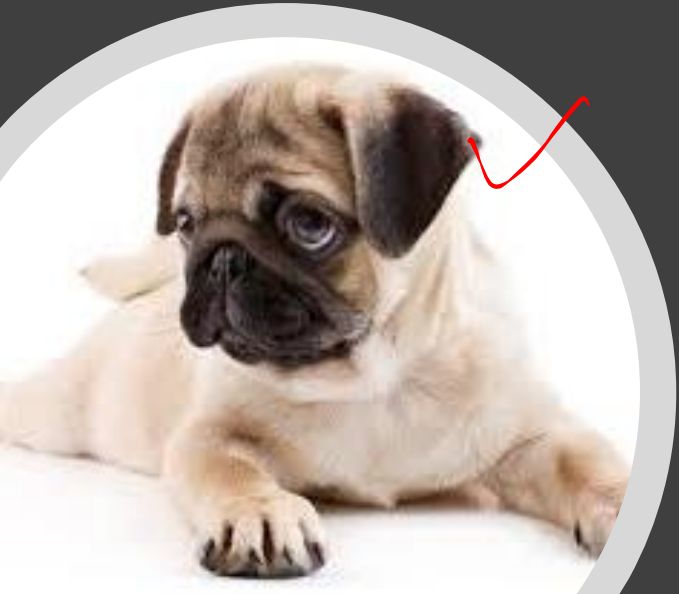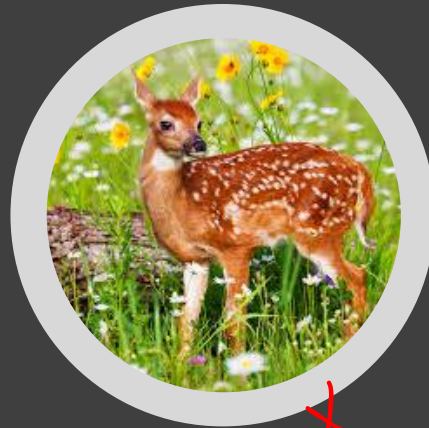- Set-up a Machine Learning Problem

# Machine Learning Everywhere

pic source: eduCBA

# What is ~~Machine~~ Learning in General

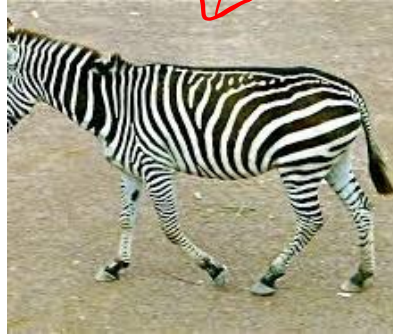- Ask a 5-year-old, is this a dog?
- Most likely, the answer is Yes

# Are these Dogs?

- It is easy for humans to identify.
- Has anyone ever defined dogs for us?
- We have learned from data.

# Learning: "Which ones are dogs?"

- Defining is hard.
- Recognizing is Easy.
- It is hard to give a mathematical definition of a Dog.
- A 5-year-old can tell the difference  (they learned from Data).
- Learning from Data is used when we do not have an analytic solution.
  - We have data to construct an analytic solution.

# The Netflix Problem

# Problem Setup

- Netflix Problem: Predict recommendations, get more subscriptions.

- Criteria used to rate movies – Unknown/Complex

- Create user and movie profiles.

- Calculate predicted rating.

- The learning algorithm 'reverse engineers' the factors based on past ratings (starting with random factors mostly).

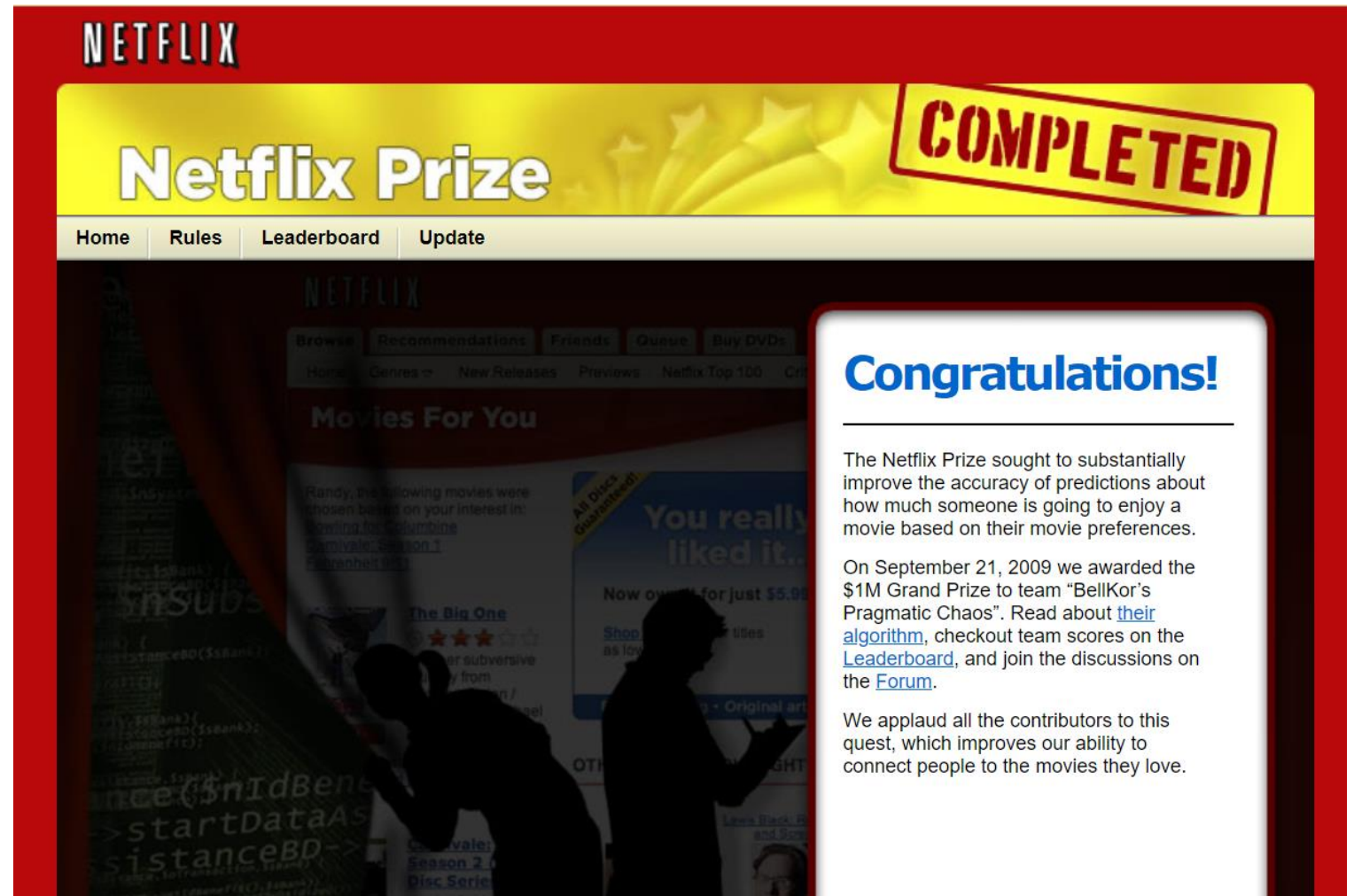- It tunes these factors to make them more aligned with real ratings of viewers.

$$2 \quad 5$$

$$V1 = [f1, f2 \cdots fn]$$

*Empirical solution*

$$M1 = [f1, f2 \cdots fn]$$

$$1 \quad 3$$

*historical*

*M2*



likes comedy? likes action? prefers blockbusters? likes Tom Cruise?

viewer ● ● ● ------- •

Match movie and viewer factors

add contributions from each factor → **predicted rating**

movie ● • ● ------ ●

comedy content action content blockbuster? Tom Cruise in it?

# Components of Learning

*Banks*

*Analytic*

- **The Credit Approval Problem:**

- Approve or not?

- No magic formula exists.

- Banks have data: customer information like salary and debt; whether they defaulted on their credit or not.

$$Cl = [ \text{---} \text{---} \text{---}$$

| | |
|---|---|
| age | 23 years ✓ |
| gender | male ✓ |
| annual salary | $30,000 ✓ |
| years in residence | 1 year ✓ |
| years in job | 1 year ✓ |
| current debt | $15,000 ✓ |
| . . . | . . . |

# Key Takeaway

- A pattern exists
- We do not know it
- We have data to learn it

# Formalize Components of Learning:

$input\ \mathbf{x} \in \mathbb{R}^d = \mathcal{X}.$

$output\ y \in \{-1, +1\} = \mathcal{Y}.$

*Real number*

$target\ function\ f : \mathcal{X} \mapsto \mathcal{Y}.$

**(The target $f$ is unknown.)**

$data\ set\ \mathcal{D} = (\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_N, y_N).$

$(y_n = f(\mathbf{x}_n).)$

*Vector $\rightarrow$ factors / attributes*

*N, Y*

*$c_1$* *$c_N$*

$y_n = f(x_n)$

$y_i = f(x_i)$

$x_1, x_2, x_3 \ldots$ *Credit level*

- **Input**: Salary, debt, years
  **Output**: Approve or not
  **Target function:** Relationship between X and Y

- Data on customers *※*

- X, Y and D will be given by the learning problem.

# The Learning Process

$H = \{$ Universe of all functions $\}$

one

- Start with a set of possible Hypothesis that are most likely to represent the target $f$.

  function $\rightarrow f$  $g$

- $H = \{h_1, h_2, \dots\}$ is the hypothesis set or the **model**.

  all possible

- Select a hypothesis $g$ from $H$. The way we did this selection (process) is the **learning algorithm**.

- Use this selected hypothesis to predict for new data (new customers). Our goal is to bring $g$ as close to $f$ as possible. The target $f$ is fixed but unknown.

  $y = g(x)$  $g \sim f$

- NOTE: We as ML practitioners will choose $H$ and the learning Algorithm.

# Summary of the Learning Set-Up



UNKNOWN TARGET FUNCTION
$$f : \mathcal{X} \mapsto \mathcal{Y}$$
*(ideal credit approval formula)*

$$y_n = f(\mathbf{x}_n)$$

TRAINING EXAMPLES
$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_N, y_N)$$
*(historical records of credit customers)*

LEARNING ALGORITHM
$$\mathcal{A}$$

FINAL HYPOTHESIS
$$g \approx f$$
*(learned credit approval formula)*

HYPOTHESIS SET
$$\mathcal{H}$$
*(set of candidate formulas)*

minimize error

Test

{ linear function, 2nd order }

$$a_1 u_1 + a_2 u_2 = y$$