# As you walk in
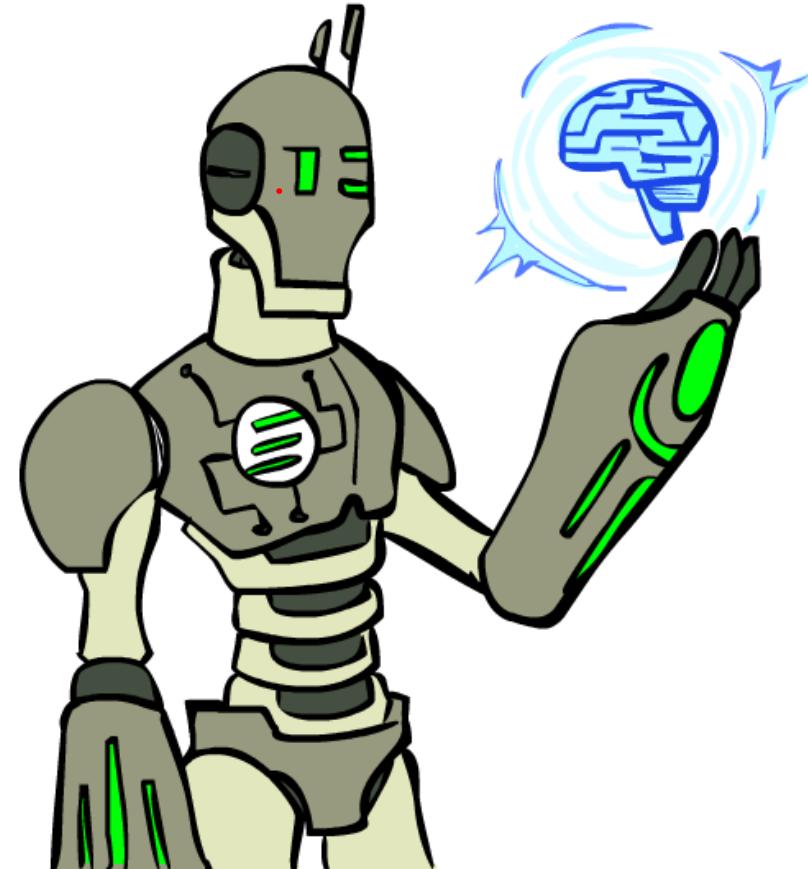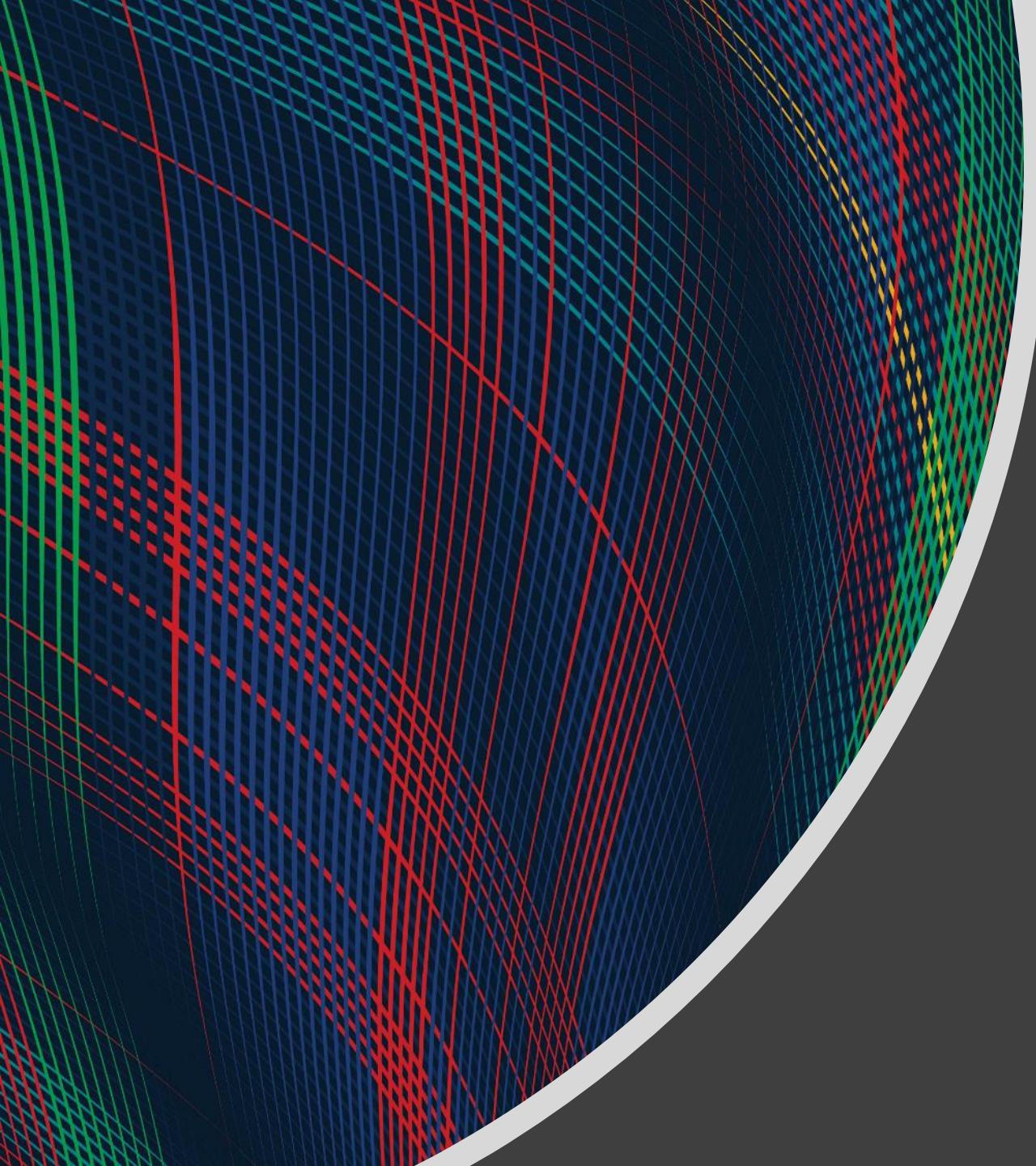
Welcome!

1) Sit at a table next to another student

2) Make name plate
   - Fold paper in half
   - Write preferred name
   - Below write you favorite fictional AI/robot

# 10-606
# Mathematical Foundations for Machine Learning

Instructor: Pat Virtue
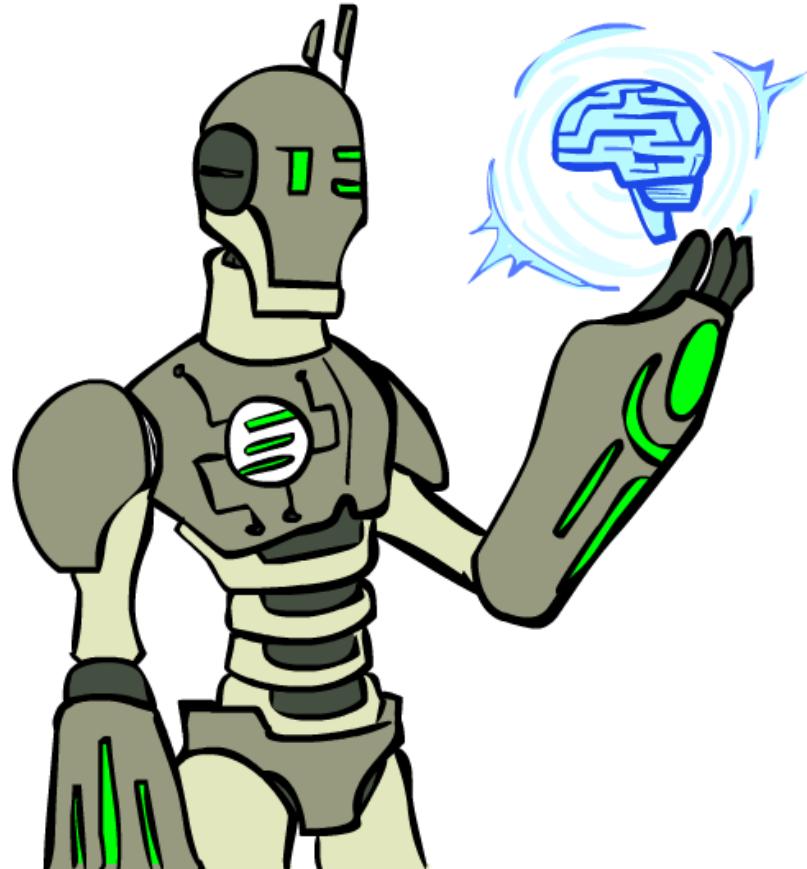
# Today

Course Info

Warm-up exercise

ML and Math Intro

Systems of equations

More Course Info

# Course Team

Instructor

Teaching Assistants



Pat
Virtue
pvirtue



Kellen
Gibson
kagibson



Ian
Char
ichar

# Course Team

Students!!

# Team Tips

Try not to act surprised

Here's a thing that happens a lot:

what's bash?

you don't know what bash is?!
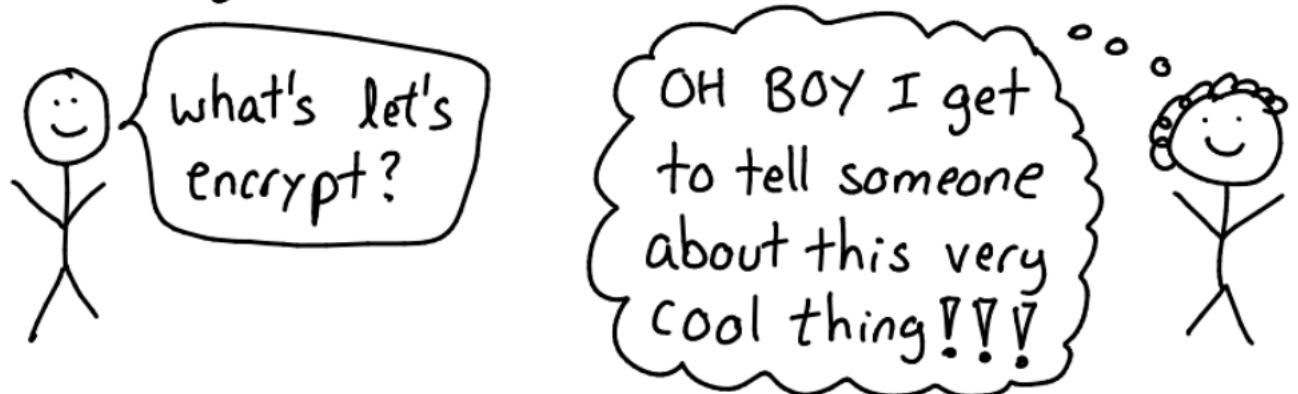
oooh I guess I was supposed to know that?

# Team Tips

Try not to act surprised

Here's a cool simple trick!

Don't act surprised when someone doesn't know something you thought they knew (even if you _are_ a little surprised!) It doesn't help.
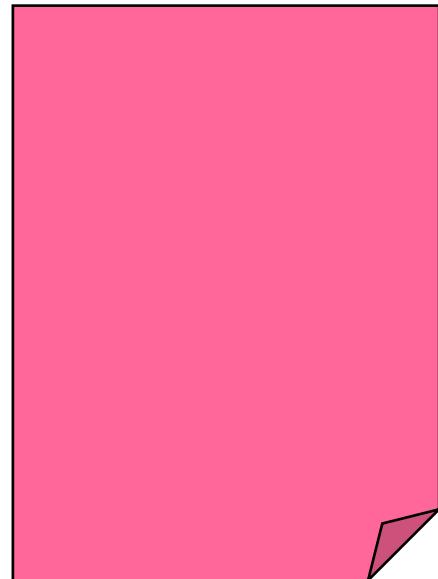
Then you get to have fun times like this:

what's let's encrypt?

OH BOY I get to tell someone about this very cool thing!!!

And it gets easier with practice!

# Team Tips

Jargon card

# Warm-up Exercise

Suppose we want to find the parameter $\hat{\theta}$ that minimizes the function $J(\theta)$.

For three different sets of assumptions about $J(\theta)$, write down an algorithm (steps, pseudocode, etc.) that will try to find a good estimate $\hat{\theta}$.

1. Assumptions: $J: \mathbb{R} \to \mathbb{R}$. We are able to call $J(\theta)$ as much as we like. That's all we know.
   Algorithm: ?

   def J(theta):

   Notation alert!

2. Assumptions: you choose
   Algorithm: ?

   convex, first deriv = 0
   $J'(\theta)$

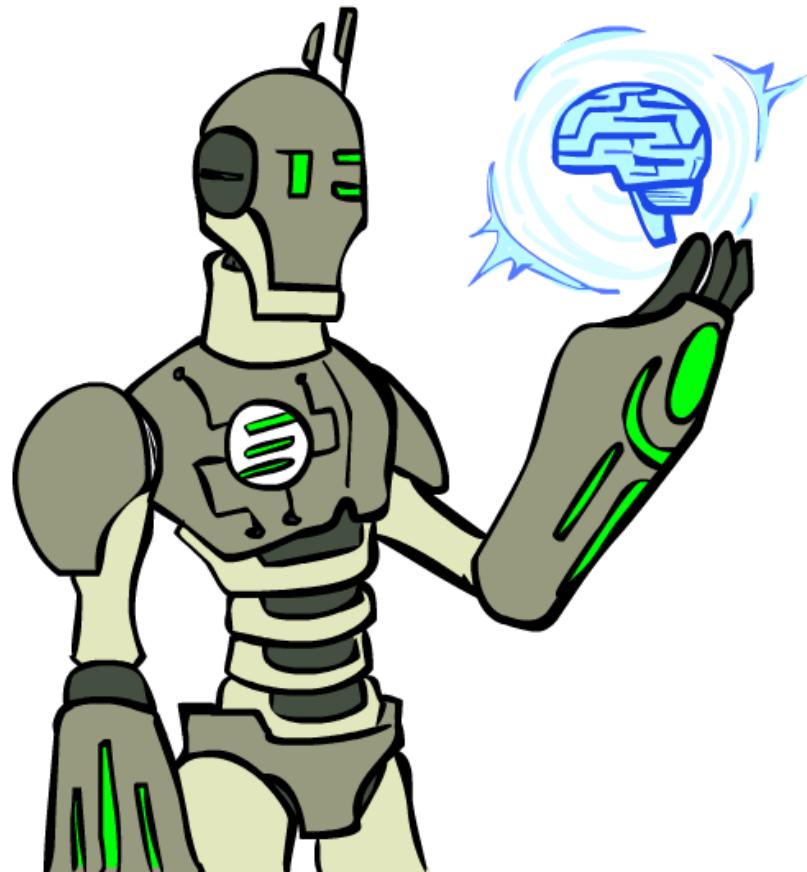3. Assumptions: you choose
   Algorithm: ?

# Today

Course Info

Warm-up exercise

ML and 606/607

Systems of equations
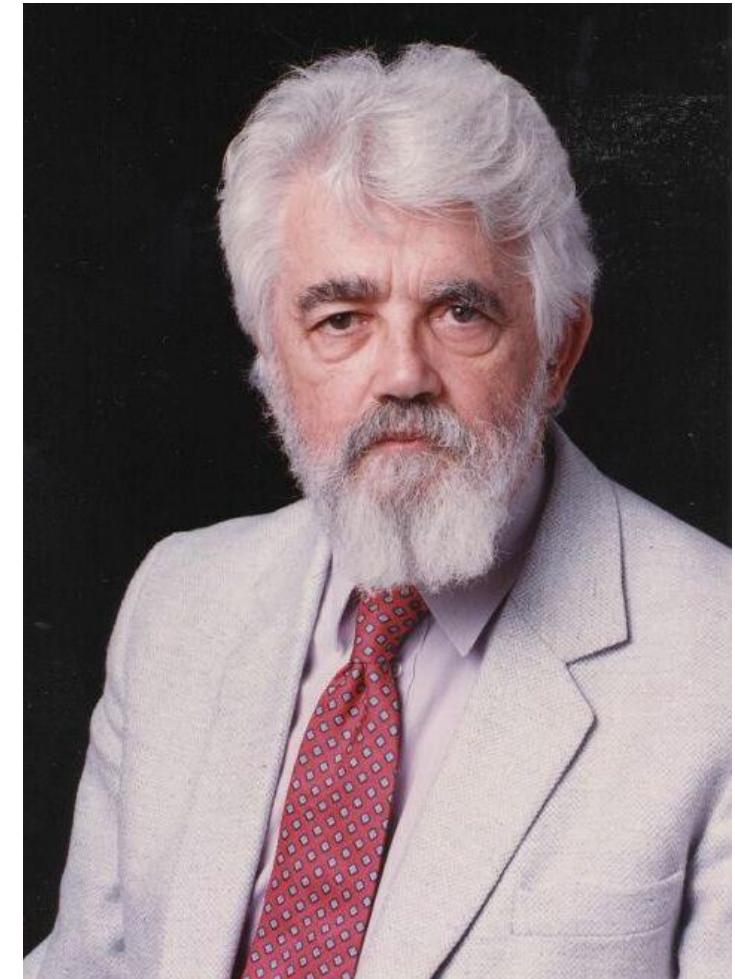
More Course Info

Images: ai.berkeley.edu

# AI Definition by John McCarthy

## What is artificial intelligence

- It is the science and engineering of making intelligent machines, especially intelligent computer programs

## What is intelligence

- Intelligence is the computational part of the ability to achieve goals in the world

http://www-formal.stanford.edu/jmc/whatisai/whatisai.html

# AI Stack for CMU AI

"AI must understand the human needs and it must make smart design decisions based on that understanding"



Autonomy

Human AI Interaction

Planning & Acting

Decision Support

Modeling

Machine Learning

Massive Data Management

Devices

Computing

Ethics

AI Stack

https://ai.cs.cmu.edu/about

# AI Stack for CMU AI
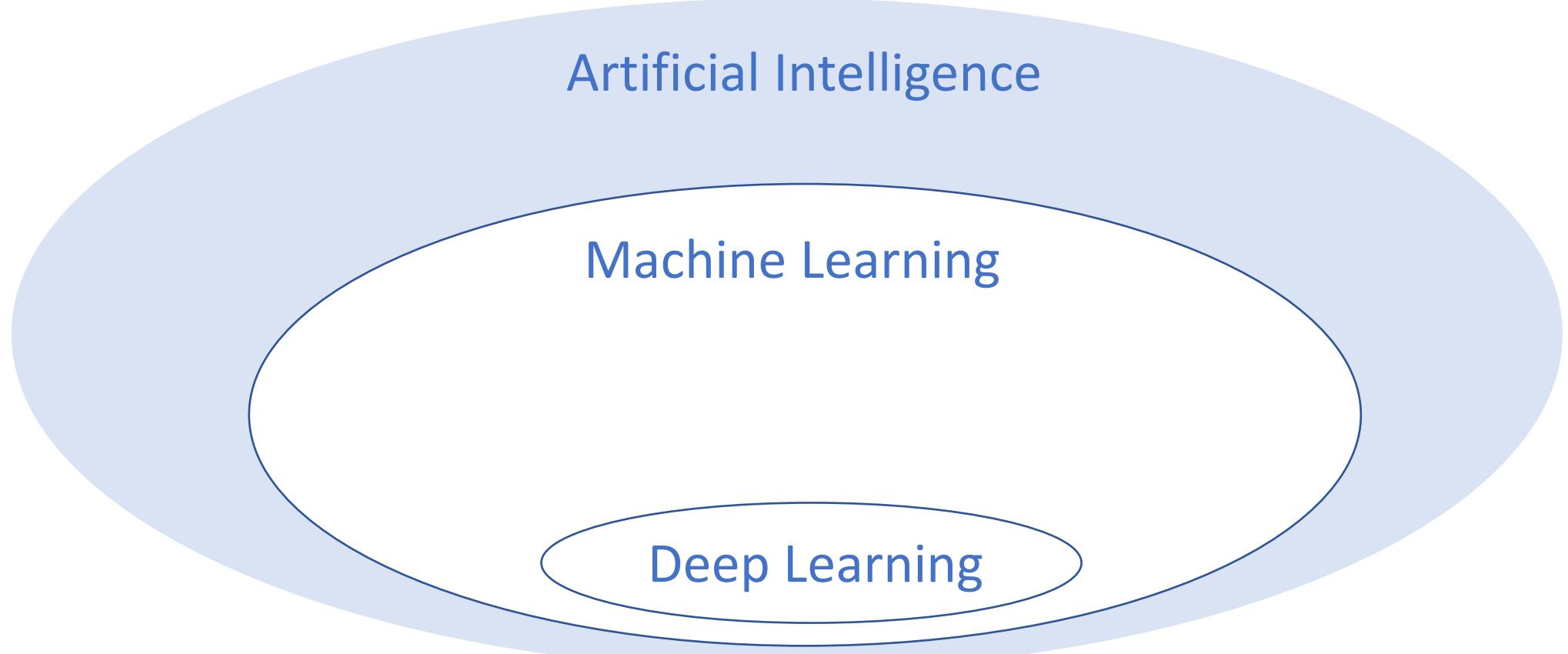
"Machine learning focuses on creating programs that learn from experience."

"It advances computing through exposure to new scenarios, testing and adaptation, while using pattern- and trend-detection to help the computer make better decisions in similar, subsequent situations."



| AI Stack |
|---|
| Autonomy / Human AI Interaction |
| Planning & Acting |
| Decision Support |
| Modeling |
| Machine Learning |
| Massive Data Management |
| Devices |
| Computing |

Ethics

https://ai.cs.cmu.edu/about

# Artificial Intelligence vs Machine Learning?

# A Brief History of AI

# A Brief History of AI



AI Excitement! 1950-1970

Knowledge Based Systems 1970-1990

Statistical Approaches 1990-

Deep Learning Era 2012-

machine learning

artificial intelligence

deep learning

formal logic

0.000400%
0.000350%
0.000300%
0.000250%
0.000200%
0.000150%
0.000100%
0.000050%
0.000000%

1940  1950  1960  1970  1980  1990  2000  2010

https://books.google.com/ngrams

# A Brief History of AI

## 1940-1950: Early days

- 1943: McCulloch & Pitts: Boolean circuit model of brain
- 1950: Turing's "Computing Machinery and Intelligence"

## 1950—70: Excitement: Look, Ma, no hands!

- 1950s: Early AI programs, including Samuel's checkers program, Newell & Simon's Logic Theorist, Gelernter's Geometry Engine
- 1956: Dartmouth meeting: "Artificial Intelligence" adopted

## 1970—90: Knowledge-based approaches

- 1969—79: Early development of knowledge-based systems
- 1980—88: Expert systems industry booms
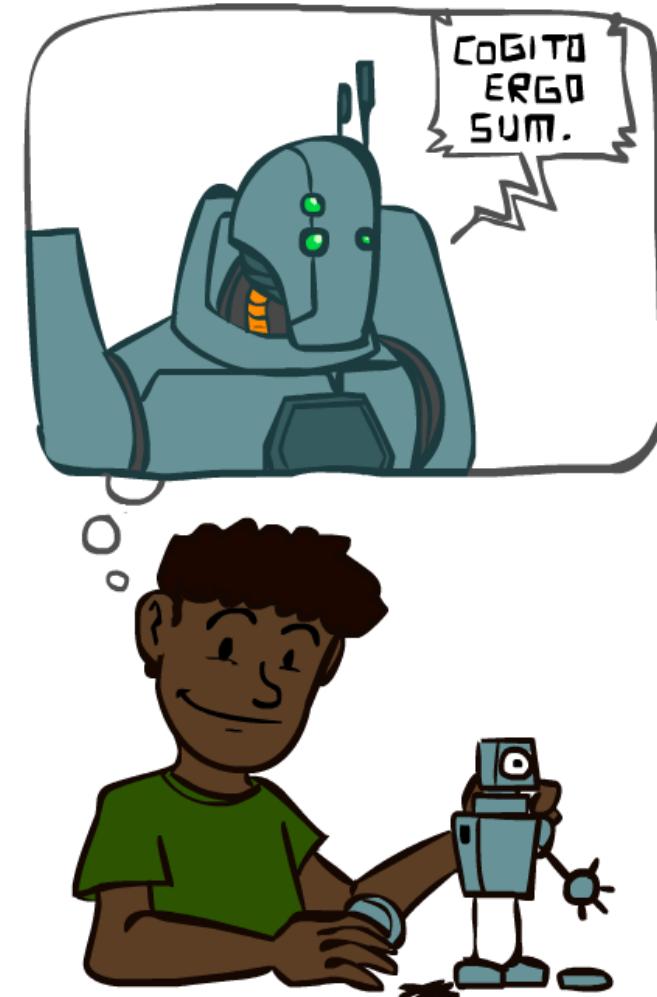- 1988—93: Expert systems industry busts: "AI Winter"

## 1990—: Statistical approaches

- Resurgence of probability, focus on uncertainty
- General increase in technical depth
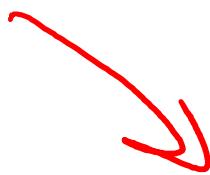- Agents and learning systems… "AI Spring"?

## 2012—: Deep learning

- 2012: ImageNet & AlexNet

Images: ai.berkeley.edu

# ML Applications?



Slide credit: CMU MLD Matt Gormley

# Speech Recognition

## 1. Learning to recognize spoken words

| THEN | NOW |
|------|-----|
| "…the SPHINX system (e.g. Lee 1989) learns speaker-specific strategies for recognizing the primitive sounds (phonemes) and words from the observed speech signal…neural network methods…hidden Markov models…" |  |

(Mitchell, 1997)

**Source:** https://www.stonetemple.com/great-knowledge-box-showdown/#VoiceStudyResults

19

# Robotics

## 2. Learning to drive an autonomous vehicle

| THEN | NOW |
|---|---|
| "…the ALVINN system (Pomerleau 1989) has used its learned strategies to drive unassisted at 70 miles per hour for 90 miles on public highways among other cars…" <br><br> (Mitchell, 1997) |  <br> https://www.geek.com/wp-content/uploads/2016/03/uber.jpg |

# Games / Reasoning

## 3. Learning to beat the masters at board games

| THEN | NOW |
|------|-----|

"…the world's top computer program for backgammon, TD-GAMMON (Tesauro, 1992, 1995), learned its strategy by playing over one million practice games against itself…"
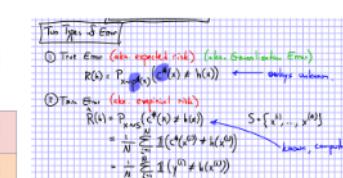
(Mitchell, 1997)

# Computer Vision

## 4. Learning to recognize images

| THEN | NOW |
|---|---|

"…The recognizer is a convolution network that can be spatially replicated. From the network output, a hidden Markov model produces word scores. The entire system is globally trained to minimize word-level errors.…"

Figure 2: Convolutional neural network character recognizer. This architecture is robust to local translations and distortions, with subsampling, shared weights, and local receptive fields.

number of subsampling layers and the sizes of the kernels are chosen, the sizes of all the layers, including the input, are determined unambiguously. The only architectural parameters that remain to be selected are the number of feature maps in each layer, and the information as to what feature map is connected to what other feature map. In our case, the subsampling rates were chosen as small as possible (2 × 2), and the kernels as small as possible in the first layer (3 × 3) to limit the total number of connections. Kernel sizes in the upper layers are chosen to be as small as

(LeCun et al., 1995)



Revolution of Depth

152 layers — 3.57

ILSVRC'15 ResNet — 3.57
ILSVRC'14 GoogleNet — 6.7 (22 layers)
ILSVRC'14 VGG — 7.3 (19 layers)
ILSVRC'13 — 11.7 (8 layers)
ILSVRC'12 AlexNet — 16.4 (8 layers)
ILSVRC'11 — 25.8 (shallow)
ILSVRC'10 — 28.2

ImageNet Classification top-5 error (%)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Learning Theory

- 5. In what cases and how well can we learn?

### Sample Complexity Results

**Definition 0.1.** The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).
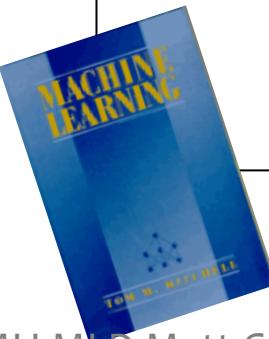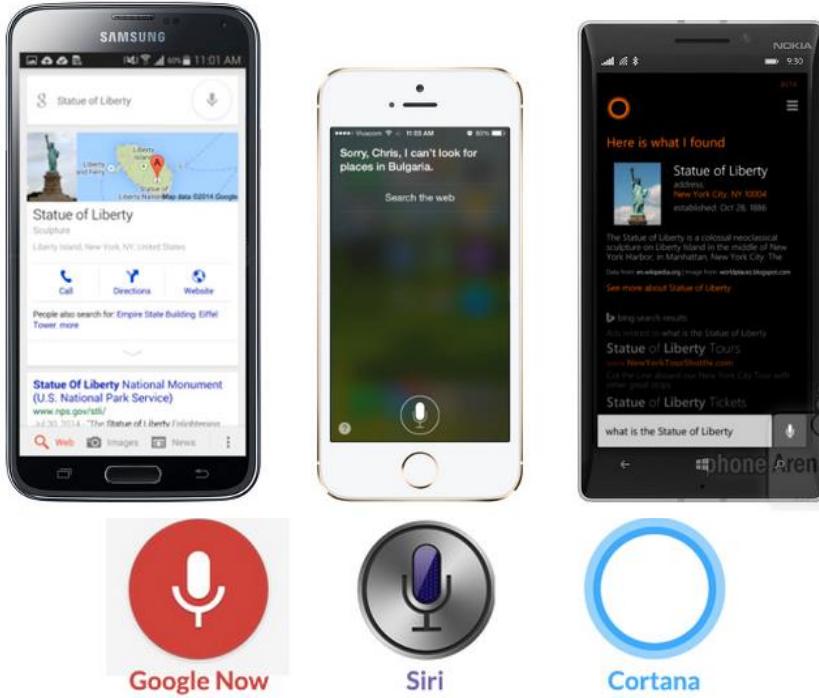
**Four Cases we care about...**

|  | Realizable | Agnostic |
|---|---|---|
| Finite $|\mathcal{H}|$ | $N \geq \frac{1}{\epsilon}\left[\log(|\mathcal{H}|) + \log(\frac{1}{\delta})\right]$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$. | $N \geq \frac{1}{2\epsilon^2}\left[\log(|\mathcal{H}|) + \log(\frac{2}{\delta})\right]$ labeled examples are sufficient so that with probability $(1-\delta)$ for all $h \in \mathcal{H}$ we have that $|R(h) - \hat{R}(h)| < \epsilon$. |
| Infinite $|\mathcal{H}|$ | $N = O(\frac{1}{\epsilon}\left[VC(\mathcal{H})\log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})\right])$ labeled examples are sufficient so that with probability $(1-\delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$. | $N = O(\frac{1}{\epsilon^2}\left[VC(\mathcal{H}) + \log(\frac{1}{\delta})\right])$ labeled examples are sufficient so that with probability $(1-\delta)$ for all $h \in \mathcal{H}$ we have that $|R(h) - \hat{R}(h)| \leq \epsilon$. |

Two Types of Error

① True Error (aka. expected risk) (aka. Generalization Error)
$$R(h) = P_{x \sim p^*(x)}\left(c^*(x) \neq h(x)\right) \quad \leftarrow \text{ always unknown}$$

② Train Error (aka. empirical risk)
$$\hat{R}(h) = P_{x \sim S}\left(c^*(x) \neq h(x)\right)$$
$$= \frac{1}{N}\sum_{i=1}^{N} \mathbb{I}\left(c^*(x^{(i)}) \neq h(x^{(i)})\right)$$
$$= \frac{1}{N}\sum_{i=1}^{N} \mathbb{I}\left(y^{(i)} \neq h(x^{(i)})\right)$$

$S = \{x^{(i)}, ..., x^{(N)}\}$

known, computable

PAC Learning

Q: Can we bound $R(h)$ in terms of $\hat{R}(h)$?
A: Yes!

PAC stands for Probably Approximately Correct

PAC learner yields hypothesis $h$, which is approximately correct $R(h) \approx 0$ with high probability $Pr(R(h) \approx 0) \approx 1$

Def: PAC Criterion
$$Pr(\forall h, |R(h) - \hat{R}(h)| \leq \epsilon) \geq 1 - \delta$$

1. How many examples do we need to learn?
2. How do we quantify our ability to generalize to unseen data?
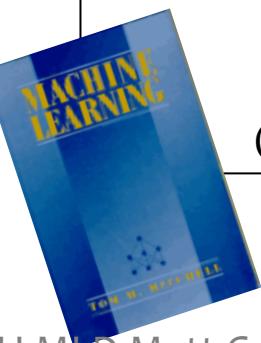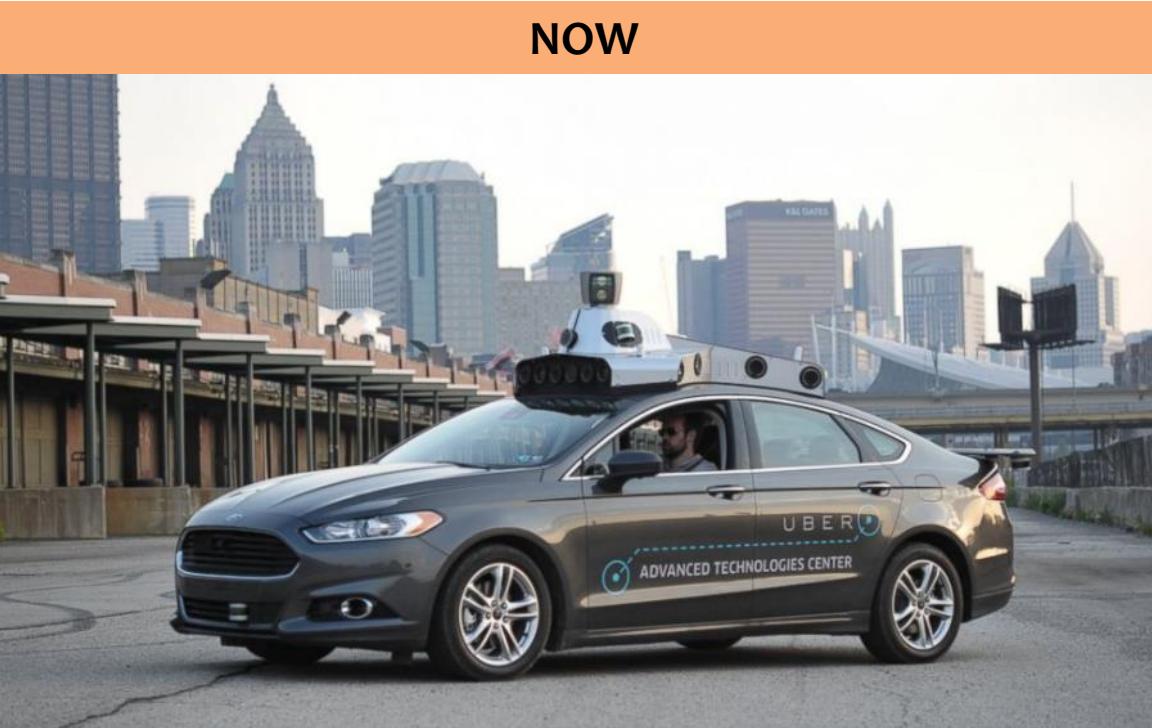3. Which algorithms are better suited to specific learning settings?

# What is ML?

607

Computer Science

607

Optimization

606

Calculus

Machine Learning

Domain of Interest

Statistics

Linear Algebra

Probability

Measure Theory

# Binary Logistic Regression

$\vec{x} \in R^M \quad \vec{\theta} \in R^M$

Gradient

$$J(\theta) = -\frac{1}{N} \sum_{i=1}^{N} \underbrace{\log p(y^{(i)} | \vec{x}^{(i)}, \vec{\theta})}_{J^{(i)}(\theta)}$$

$\nabla_\theta J(\theta)$

$R^M \longrightarrow R^M$

$$\nabla_\theta J^{(i)}(\theta) = \begin{bmatrix} \frac{\partial J^{(i)}}{\partial \theta_1} \\ \vdots \\ \frac{\partial J^{(i)}}{\partial \theta_m} \\ \vdots \end{bmatrix} \longrightarrow \frac{\partial J^{(i)}}{\partial \theta_m} = -\left(y^{(i)} - \sigma(\vec{\theta}^T \vec{x})\right) x_m^{(i)}$$

Example slide from 10-601

SGD $\longrightarrow \nabla_\theta J^{(i)} = -\left(y^{(i)} - \sigma(\vec{\theta}^T \vec{x})\right) \vec{x}$

# Naïve Bayes MLE

$$\mathcal{D} = \{y^{(n)}, \boldsymbol{x}^{(n)}\}_{n=1}^{N}$$
$$y^{(n)} \in \{0,1\}$$
$$\boldsymbol{x}^{(n)} \in \{0,1\}^{M}$$
$$\phi \in [0,1]$$
$$\Theta \in [0,1]^{Mx2}$$

$$L(\phi, \Theta) = p(\mathcal{D} \mid \phi, \Theta)$$

$$= \Pi_{n=1}^{N} p(\mathcal{D}^{(n)} \mid \phi, \Theta) \quad \text{i.i.d assumption}$$

$$= \Pi_{n=1}^{N} p(y^{(n)}, \boldsymbol{x}^{(n)} \mid \phi, \Theta)$$

$$= \Pi_{n=1}^{N} p(y^{(n)} \mid \phi) p(\boldsymbol{x}^{(n)} \mid y^{(n)}, \Theta) \quad \text{Generative model}$$

$$= \Pi_{n=1}^{N} p(y^{(n)} \mid \phi) p(x_1^{(n)}, x_2^{(n)}, \dots, x_M^{(n)} \mid y^{(n)}, \Theta)$$

$$= \Pi_{n=1}^{N} p(y^{(n)} \mid \phi) \Pi_{m=1}^{M} p(x_m^{(n)} \mid y^{(n)}, \theta_{m,y}) \quad \text{Naïve Bayes}$$

$$= \Pi_{n=1}^{N} \phi^{y^{(n)}} (1-\phi)^{1-y^{(n)}} \Pi_{m=1}^{M} \theta_{m,1}^{\mathbb{I}\left(y^{(n)}=1 \wedge x_m^{(n)}=1\right)} (1-\theta_{m,1})^{\mathbb{I}\left(y^{(n)}=1 \wedge x_m^{(n)}=0\right)}$$
$$\theta_{m,0}^{\mathbb{I}\left(y^{(n)}=0 \wedge x_m^{(n)}=1\right)} (1-\theta_{m,0})^{\mathbb{I}\left(y^{(n)}=0 \wedge x_m^{(n)}=0\right)}$$

$$= \phi^{N_{y=1}} (1-\phi)^{N_{y=0}} \Pi_{m=1}^{M} \theta_{m,1}^{N_{y=1,x_m=1}} (1-\theta_{m,1})^{N_{y=1,x_m=0}} \theta_{m,0}^{N_{y=0,x_m=1}} (1-\theta_{m,0})^{N_{y=0,x_m=0}}$$

Example slide from 10-601

# Network Optimization: Layer Implementation

$$J(\mathbf{w}) = z_3$$
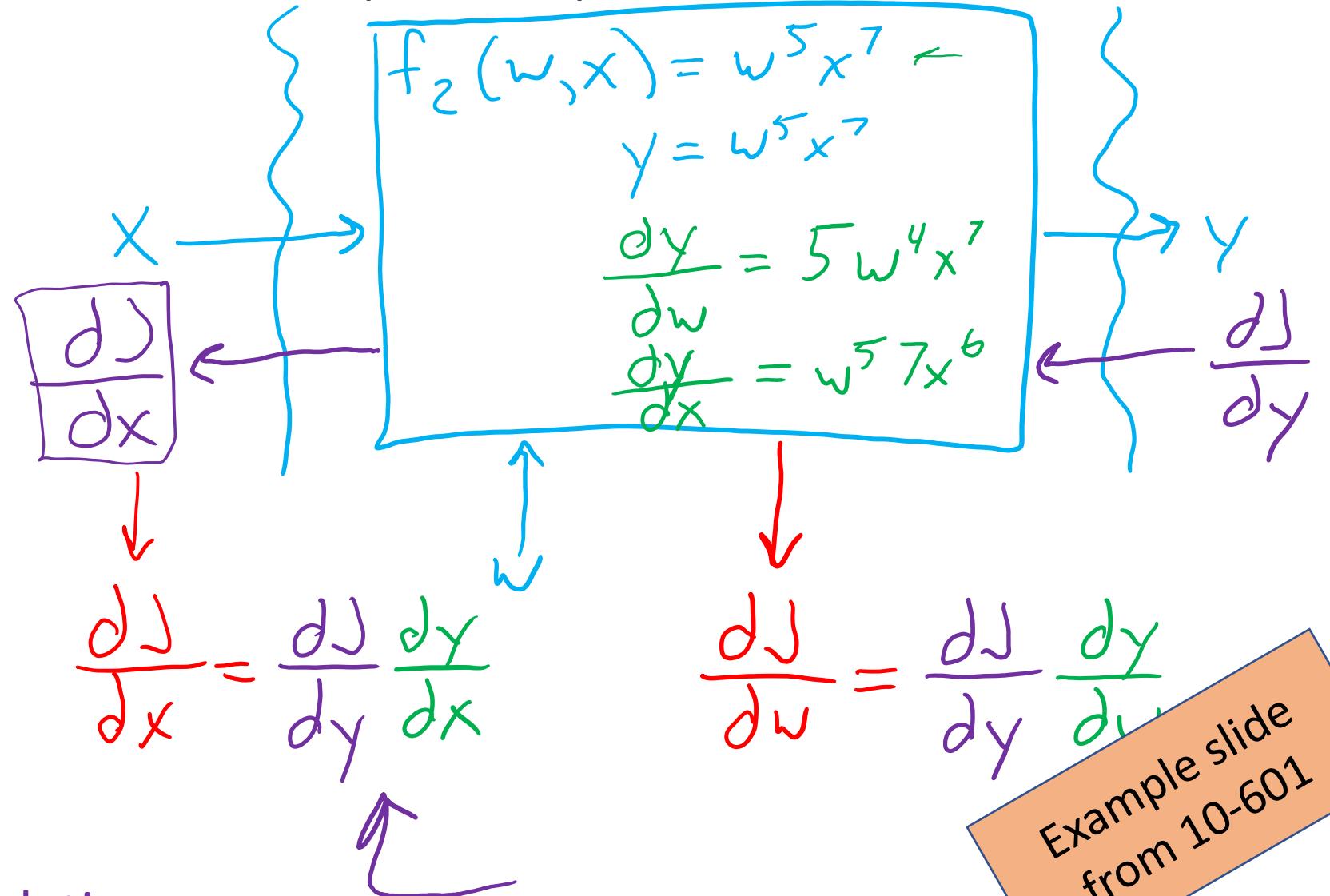
$$z_3 = f_3(w_3, z_2)$$

$$z_2 = f_2(w_2, z_1)$$

$$z_1 = f_1(w_1, x)$$

$$\frac{\partial J}{\partial w_3} = \frac{\partial J}{\partial z_3}\frac{\partial z_3}{\partial w_3}$$

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial z_3}\frac{\partial z_3}{\partial z_2}\frac{\partial z_2}{\partial w_2}$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial z_3}\frac{\partial z_3}{\partial z_2}\frac{\partial z_2}{\partial z_1}\frac{\partial z_1}{\partial w_1}$$

Lots of repeated calculations

$$f_2(w, x) = w^5 x^7$$

$$y = w^5 x^7$$

$$\frac{\partial y}{\partial w} = 5 w^4 x^7$$

$$\frac{\partial y}{\partial x} = w^5 7 x^6$$

$$x \longrightarrow \boxed{\;} \longrightarrow y$$

$$\frac{\partial J}{\partial x} \longleftarrow \boxed{\;} \longleftarrow \frac{\partial J}{\partial y}$$

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial y}\frac{\partial y}{\partial x}$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial y}\frac{\partial y}{\partial w}$$

Example slide from 10-601

# PCA: the First Principal Component

To find the first principal component, we wish to solve the following constrained optimization problem (variance minimization).

$$\mathbf{v}_1 = \underset{\mathbf{v}:||\mathbf{v}||^2=1}{\operatorname{argmax}} \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} \tag{1}$$

So we turn to the method of Lagrange multipliers. The Lagrangian is:

$$\mathcal{L}(\mathbf{v}, \lambda) = \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} - \lambda(\mathbf{v}^T \mathbf{v} - 1) \tag{2}$$

Taking the derivative of the Lagrangian and setting to zero gives:

$$\frac{d}{d\mathbf{v}} \left( \mathbf{v}^T \boldsymbol{\Sigma} \mathbf{v} - \lambda(\mathbf{v}^T \mathbf{v} - 1) \right) = 0 \tag{3}$$

$$\boldsymbol{\Sigma} \mathbf{v} - \lambda \mathbf{v} = 0 \tag{4}$$

$$\boldsymbol{\Sigma} \mathbf{v} = \lambda \mathbf{v} \tag{5}$$

Recall: For a square matrix $\mathbf{A}$, the vector $\mathbf{v}$ is an **eigenvector** iff there exists **eigenvalue** $\lambda$ such that:

$$\mathbf{A}\mathbf{v} = \lambda \mathbf{v} \tag{6}$$

Example slide from 10-601

30

# 10-606 and 10-607

- Mini Courses
  - 10-606 ←
  - 10-607

- Intro ML Courses
  - 10-315
  - 10-301/601
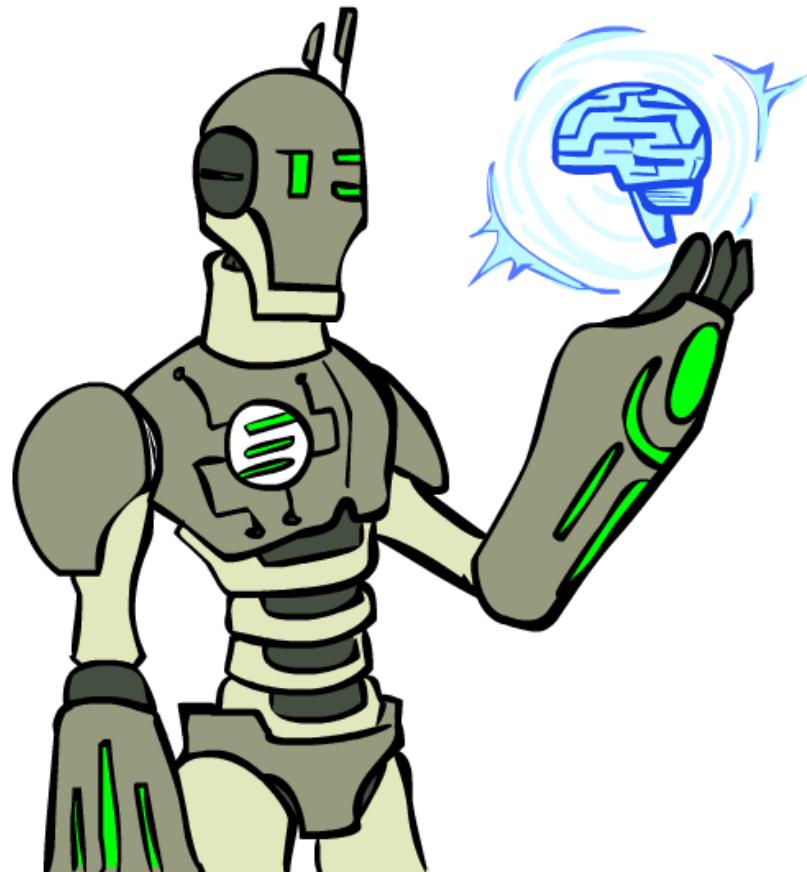  - 10-701
  - 10-715

- Prerequisites ←

# Today

Course Info

Warm-up exercise

ML and 606/607
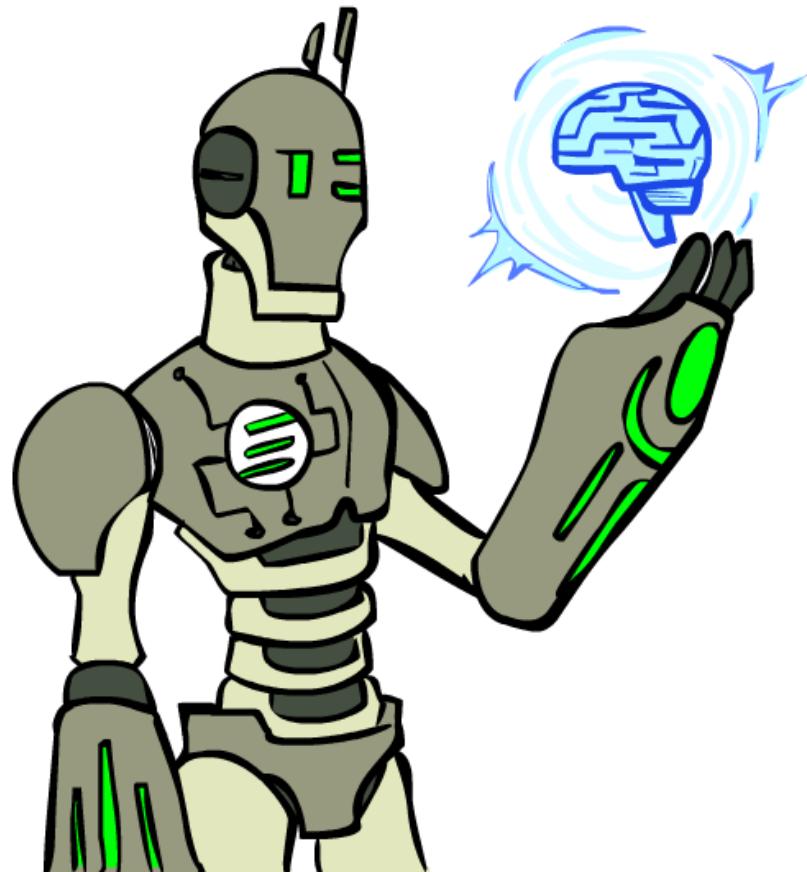
Systems of equations

More Course Info

Images: ai.berkeley.edu

# Today

Course Info

Warm-up exercise

ML and 606/607

Systems of equations

More Course Info

# Course Information

Website: https://www.cs.cmu.edu/~10606

Canvas: canvas.cmu.edu

Gradescope: gradescope.com

Communication:

piazza.com

E-mail (if piazza doesn't work):

pvirtue@andrew.cmu.edu

# Course Information

## Lectures

- Lectures are recorded
  - Shared with our course and ML course staff only
- Participation point earned by answering Piazza polls in lecture
- Quizzes will in lecture, announced two days ahead of time
- Slides will be posted

## Recitations

- Recommended attendance
- No plans to record at this point
- No participation points in recitation
- Recitation materials are in-scope for quizzes and exams

# Course Information

## Office Hours

- OH calendar on course website

- OH-by-appointment requests are certainly welcome

## Mental Health