# Deep Learning
## Chris Piech
## CS109, Stanford University

# Innovations in deep learning



AlphaGO (2016)

Deep learning (neural networks) is the core idea driving the current revolution in AI.

Notes:
- Checkers is the last **solved** game (from game theory, where perfect player outcomes can be fully predicted from any gameboard).
  https://en.wikipedia.org/wiki/Solved_game
- The first machine learning algorithm defeated a world champion in Chess in 1996.
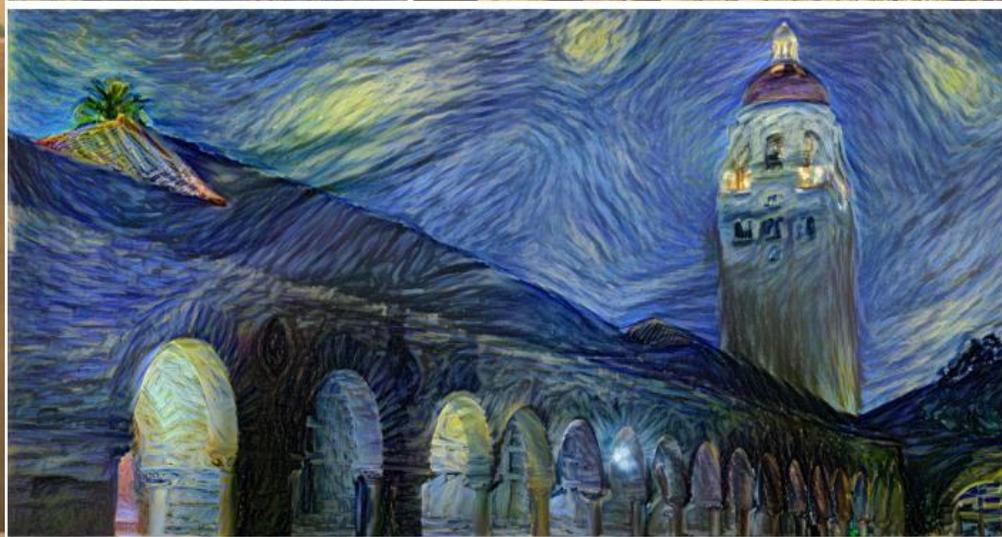  https://en.wikipedia.org/wiki/Deep_Blue_(chess_computer)
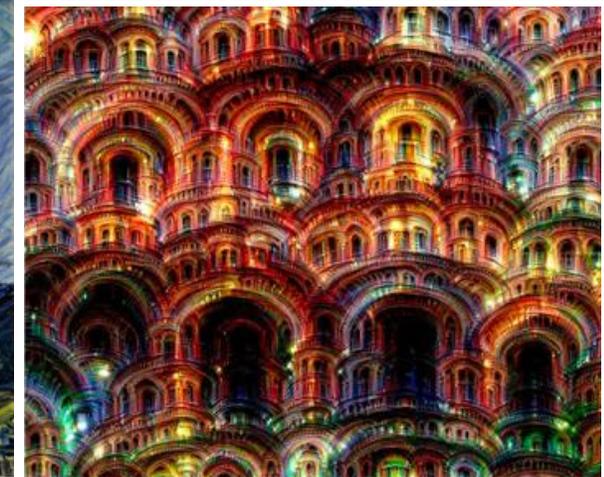
Self Driving Cars

# Computers making art



The Next Rembrandt
https://medium.com/@DutchDigital/the-next-rembrandt-bringing-the-old-master-back-to-life-35dfb1653597

A Neural Algorithm of Artistic Style
https://arxiv.org/abs/1508.06576
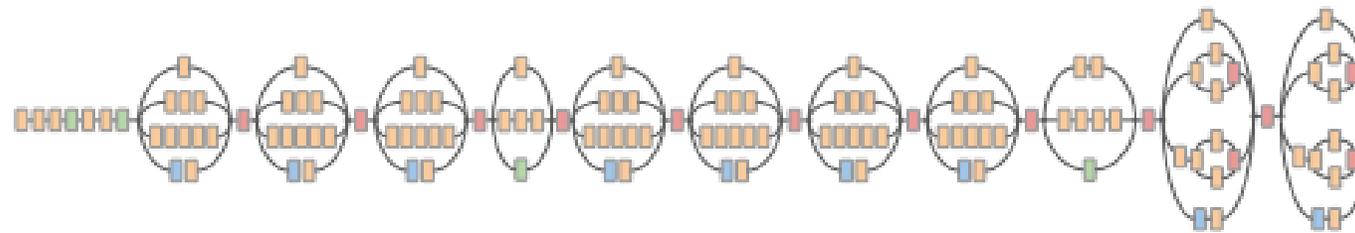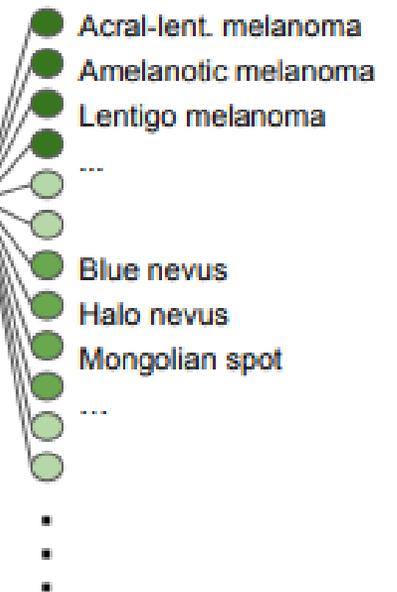https://github.com/jcjohnson/neural-style

Google Deep Dream
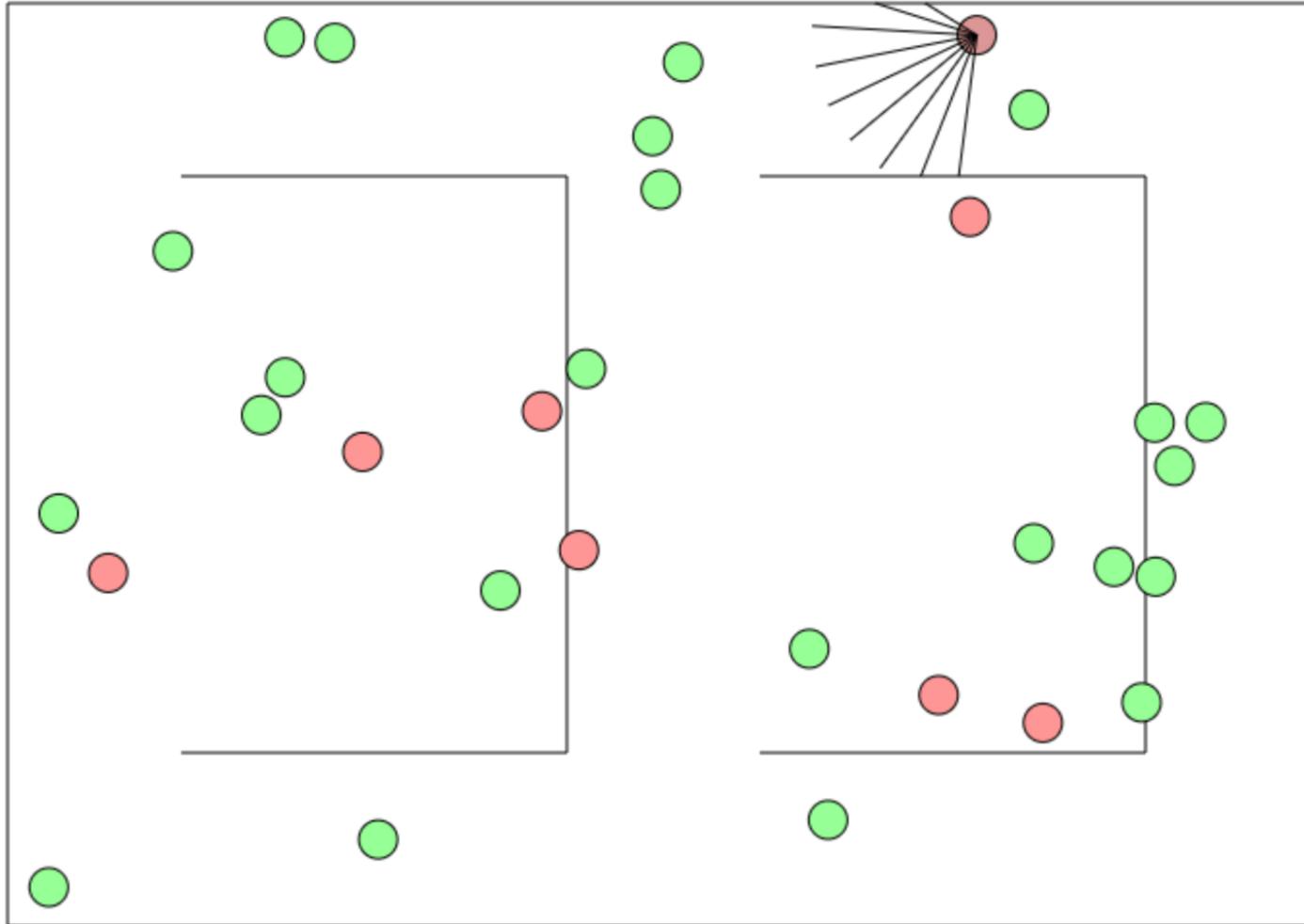https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html

# Detecting skin cancer



Esteva, Andre, et al. "Dermatologist-level classification of skin cancer with deep neural networks." *Nature* 542.7639 (2017): 115-118.

# Our Little Buddy from Last Class



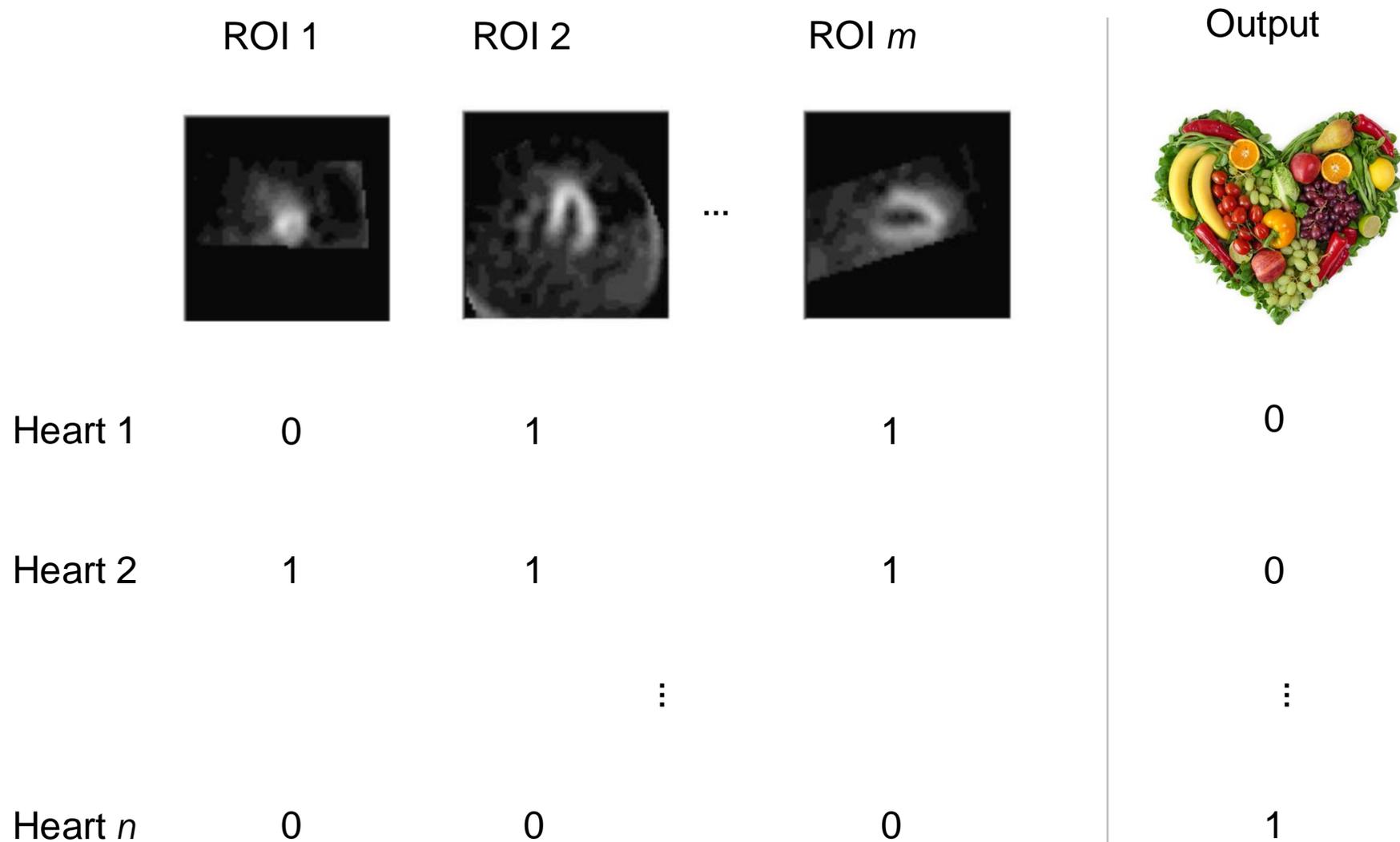http://cs.stanford.edu/people/karpathy/convnetjs/demo/rldemo.html

# Logistics

# After Break Schedule

| | Monday | Tuesday | Wednesday | Thursday | Friday |
|---|---|---|---|---|---|
| W10 | Deep Learning | PEP | Diffusion | | PSet 6 is Due |
| W11 | | Final Exam | | | |

# Review

# Classification Task

| | ROI 1 | ROI 2 | | ROI $m$ | Output |
|---|---|---|---|---|---|
| |  |  | ... |  |  |
| Heart 1 | 0 | 1 | | 1 | 0 |
| Heart 2 | 1 | 1 | | 1 | 0 |
| | | ⋮ | | | ⋮ |
| Heart $n$ | 0 | 0 | | 0 | 1 |

# Machine Learning

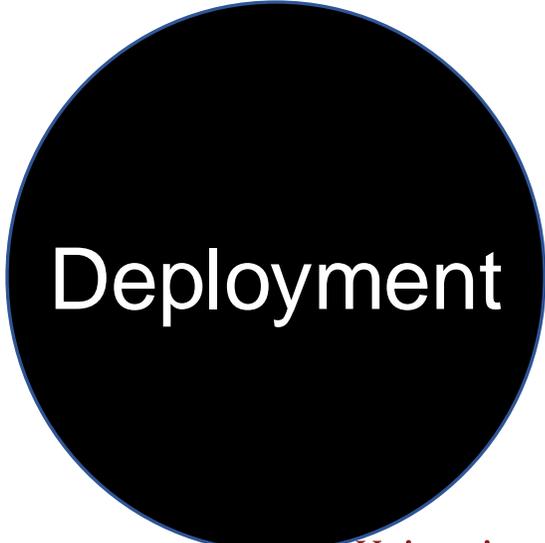$\mathbf{x}$

(inputs)

$\theta$

(model)

$\hat{y}$

(prediction)

# The Training / Testing Paradigm

**Dataset**

Deployment

# The Training / Testing Paradigm

**Training**

**Testing**

If your model passes testing...

Learn your parameters

Make sure that they work

Deployment

# Logistic Regression

$\boldsymbol{x}$:   0   1   1

$x_0$

$\theta_0$

$x_1$

$\theta_1$

$x_2$

$\theta_2$

$z$

$\sigma(z)$

$P(Y = 1|x)$

$x_3$

$\theta_3$

$$P(Y = 1|X = \mathbf{x}) = \sigma(\theta^T \mathbf{x})$$
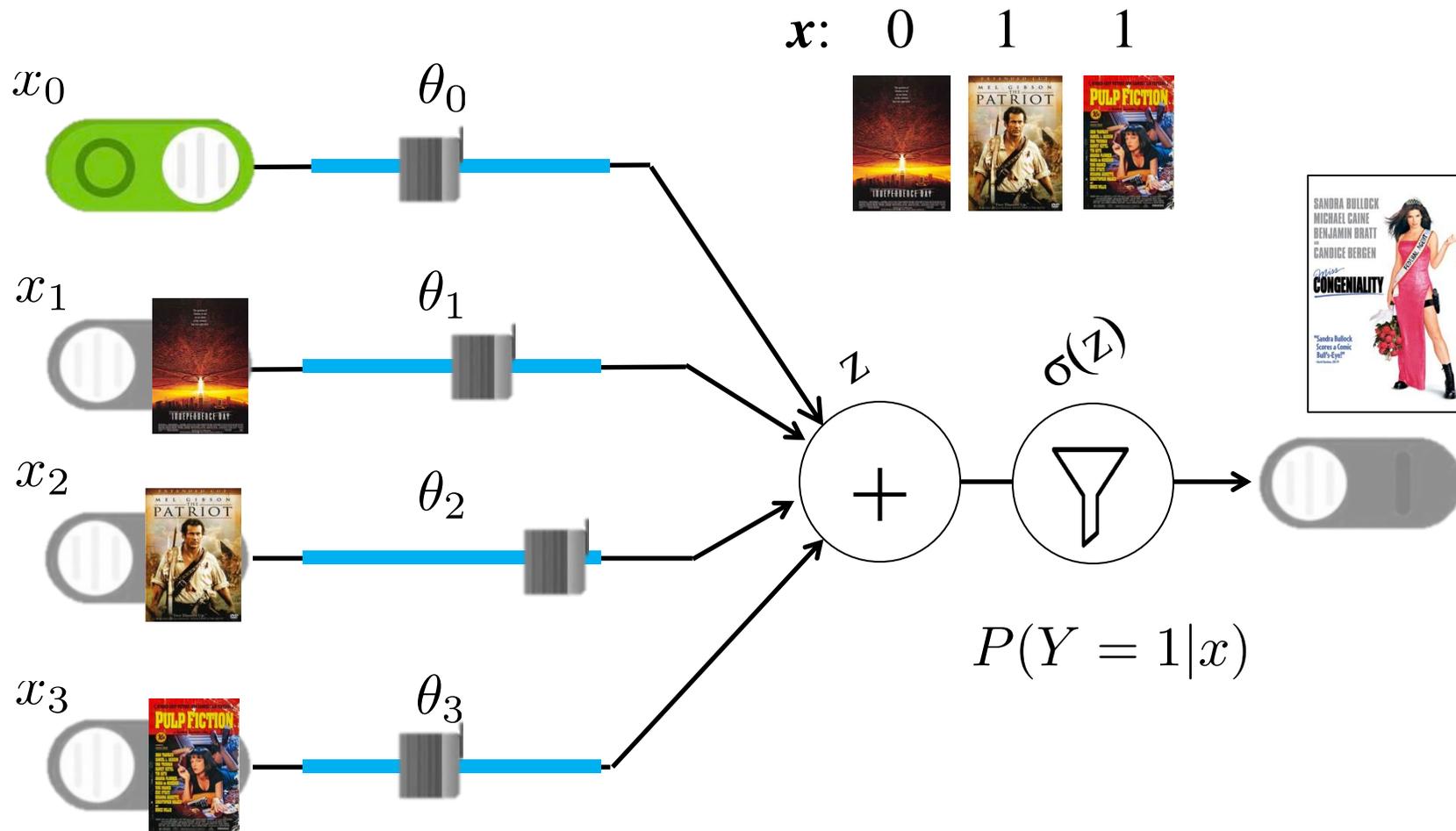
# Logistic Regression



$$P(Y = 1 | X = \mathbf{x}) = \sigma(\theta^T \mathbf{x})$$

# Logistic Regression

$x$:   0   1   1

$x_0$   $\theta_0$

$x_1$   $\theta_1$

$x_2$   $\theta_2$

$x_3$   $\theta_3$
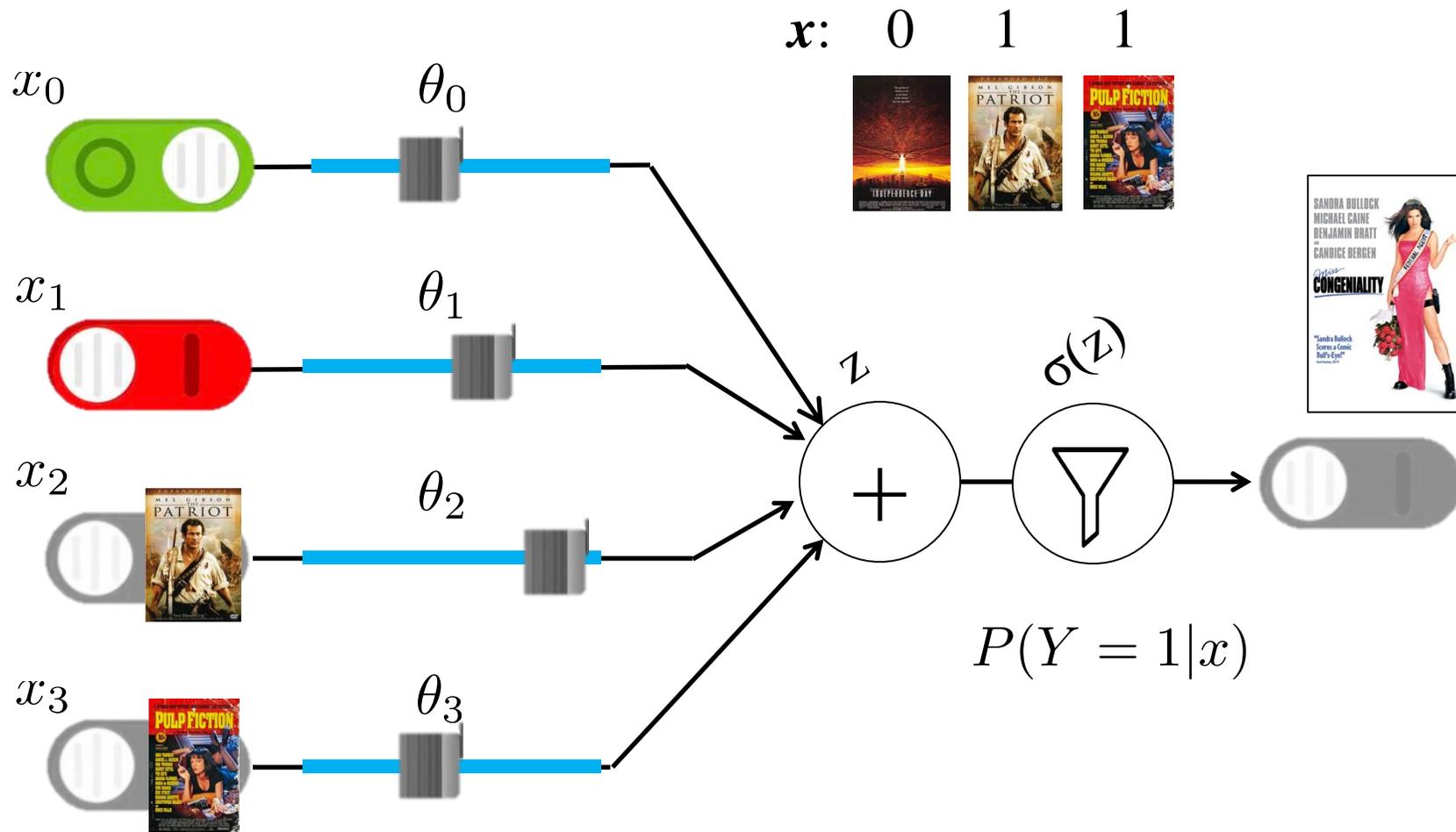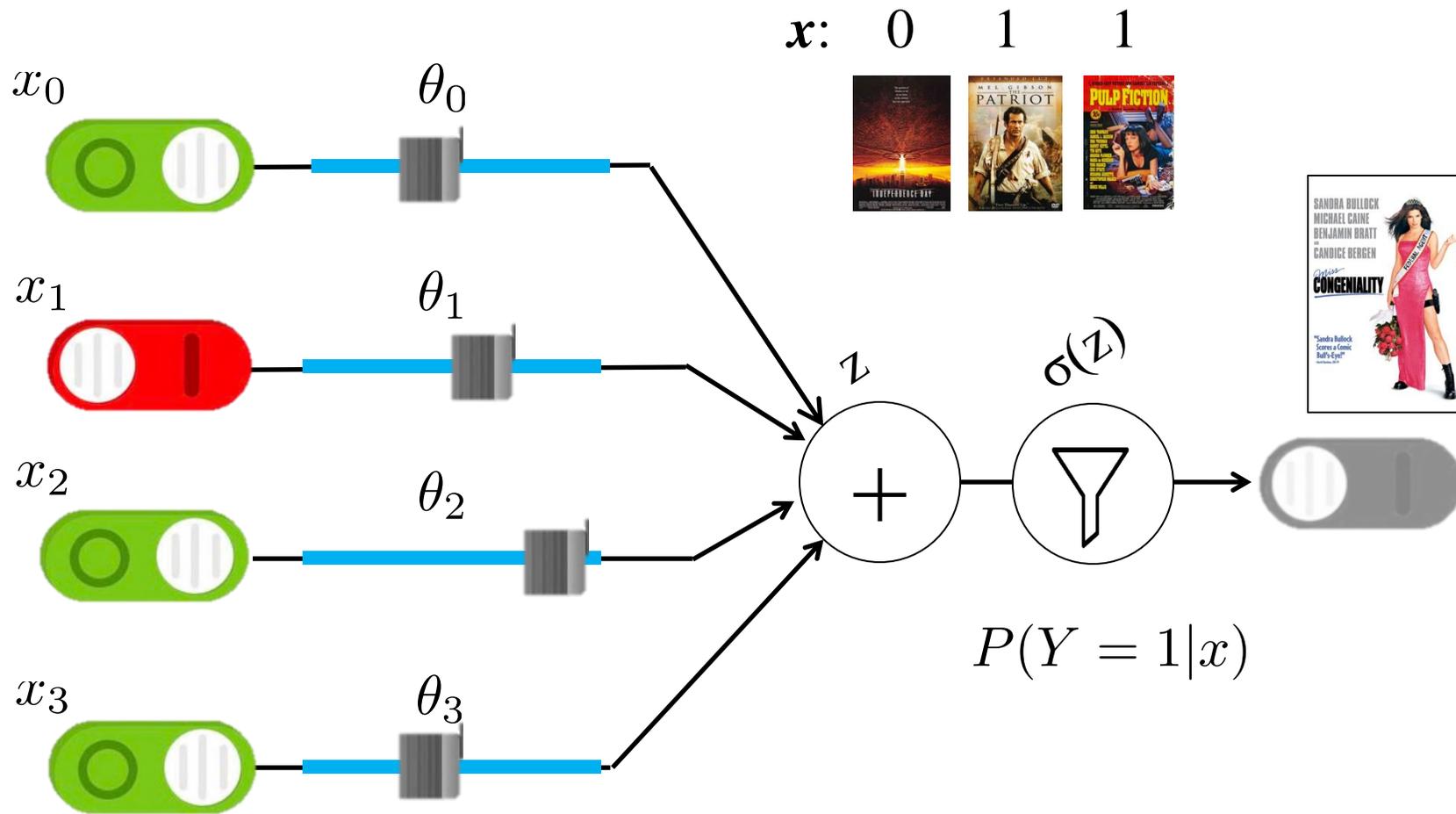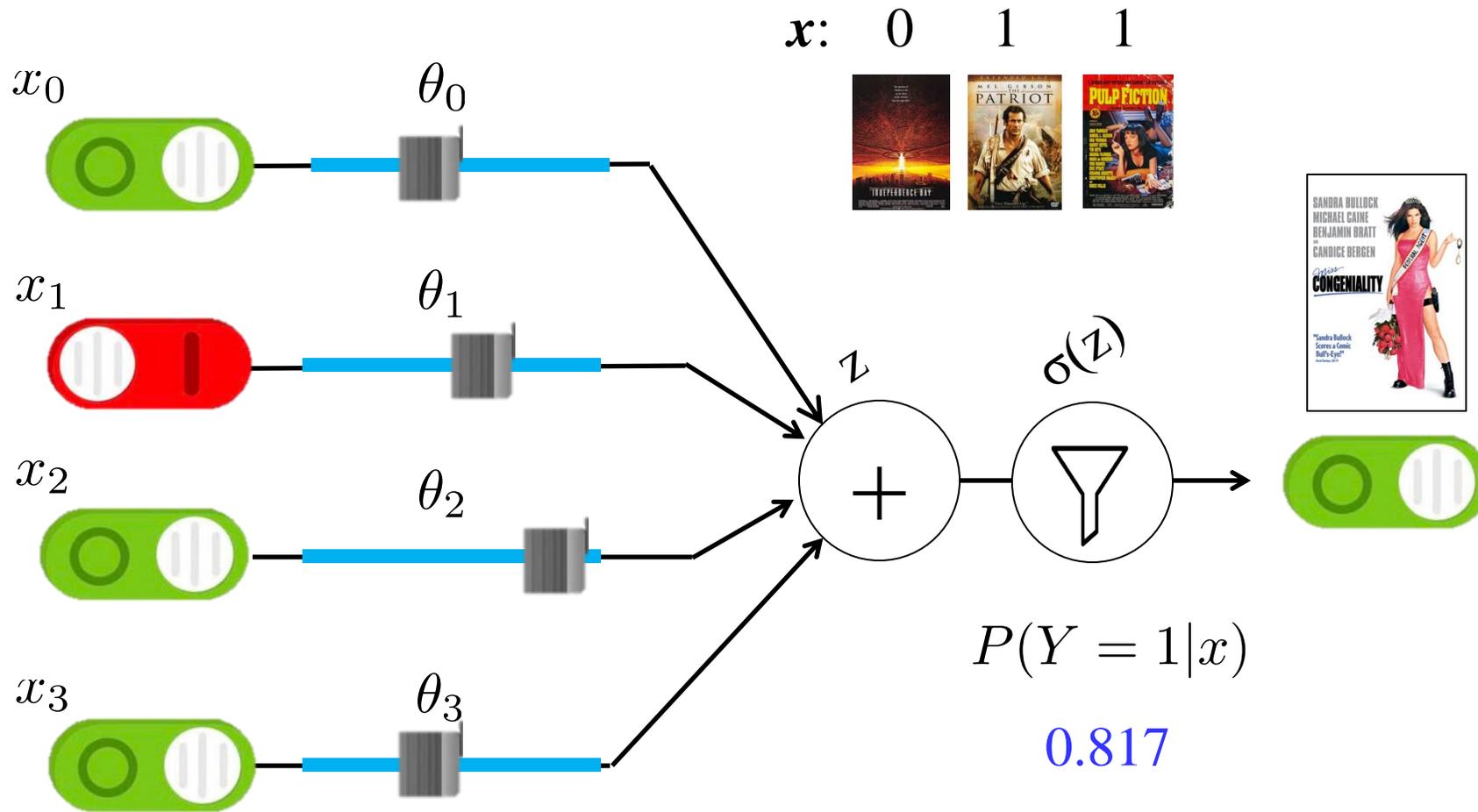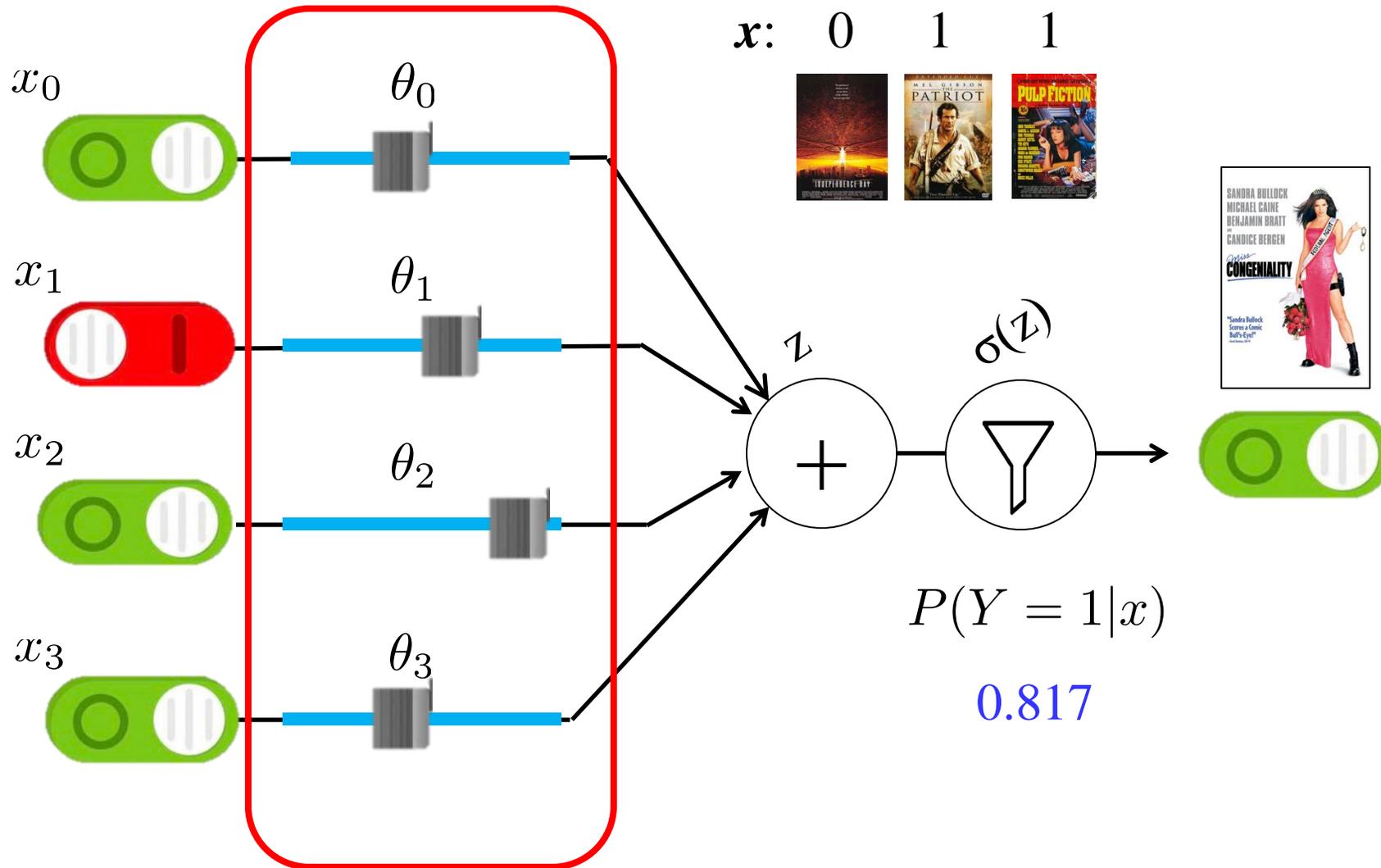
$z$   $\sigma(z)$

$+$

$P(Y = 1|x)$

$$P(Y = 1|X = \mathbf{x}) = \sigma(\theta^T \mathbf{x})$$

# Logistic Regression



$$P(Y = 1|X = \mathbf{x}) = \sigma(\theta^T \mathbf{x})$$

# Logistic Regression

$x$:   0   1   1

$x_0$     $\theta_0$

$x_1$     $\theta_1$

$x_2$     $\theta_2$

$x_3$     $\theta_3$

$z$

$\sigma(z)$

$P(Y = 1|x)$

0.817

$$P(Y = 1|X = \mathbf{x}) = \sigma(\theta^T \mathbf{x})$$

# Logistic Regression



$$P(Y = 1 | X = \mathbf{x}) = \sigma(\theta^T \mathbf{x})$$

# Math for Logistic Regression

**1**   Make logistic regression assumption

$$P(Y = 1 | X = \mathbf{x}) = \sigma(\theta^T \mathbf{x})$$

$$P(Y = 0 | X = \mathbf{x}) = 1 - \sigma(\theta^T \mathbf{x})$$

*Often call this* $\hat{y}$

**2**   Calculate the log likelihood for all data

$$LL(\theta) = \sum_{i=1}^{n} y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log[1 - \sigma(\theta^T \mathbf{x}^{(i)})]$$

**3**   Get derivative of log likelihood with respect to thetas

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^{n} \left[ y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)}) \right] x_j^{(i)}$$

# Logistic Regression Training

Initialize: $\theta_j = 0$ for all $0 \leq j \leq m$

Repeat many times:

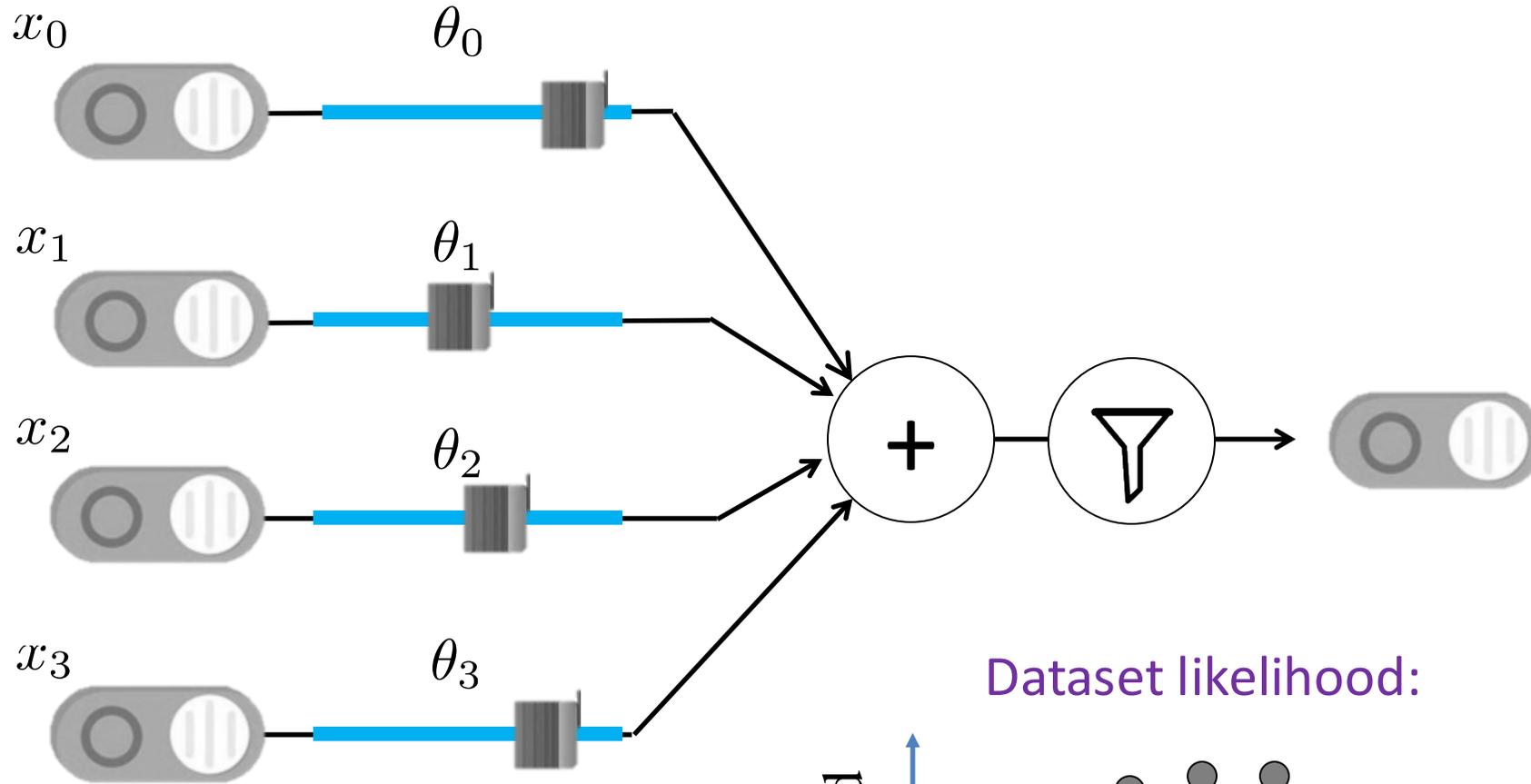gradient[j] = 0 for all $0 \leq j \leq m$
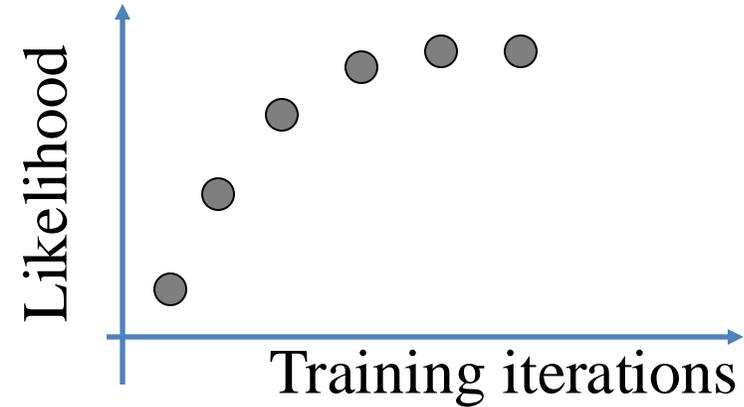
For each parameter j

For each training example (x, y):

$$\text{gradient[j]} \mathrel{+}= x_j \left( y - \frac{1}{1 + e^{-\theta^{\mathrm{T}}\mathbf{x}}} \right)$$

$\theta_j$ += η * gradient[j] for all $0 \leq j \leq m$

# Training



$x_0$     $\theta_0$

$x_1$     $\theta_1$
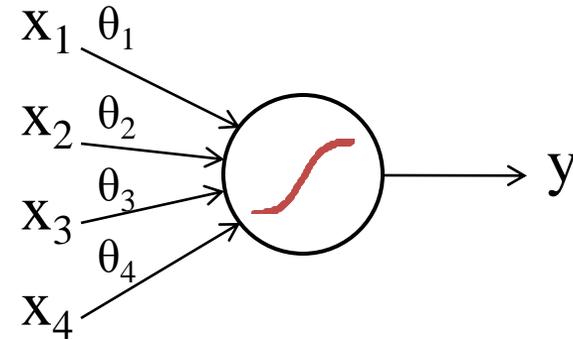
$x_2$     $\theta_2$

$x_3$     $\theta_3$

$+$
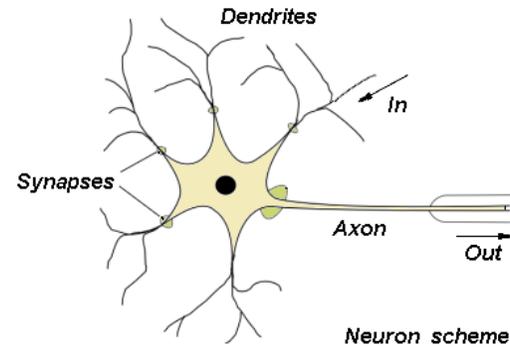
Dataset likelihood:

Likelihood

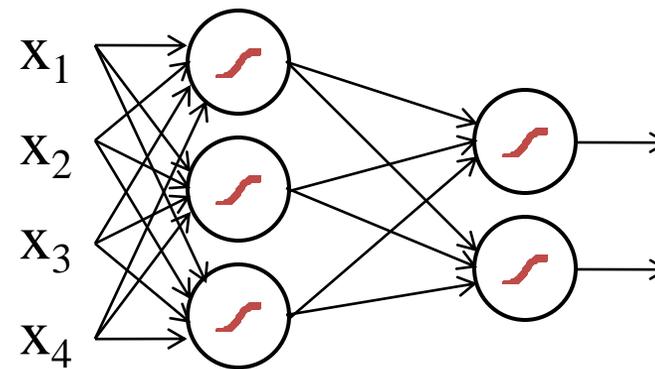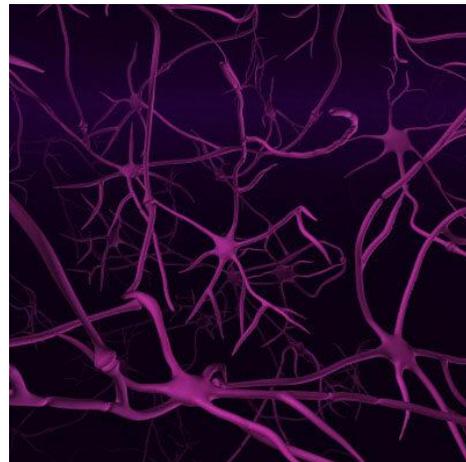Training iterations

# Artificial Neurons

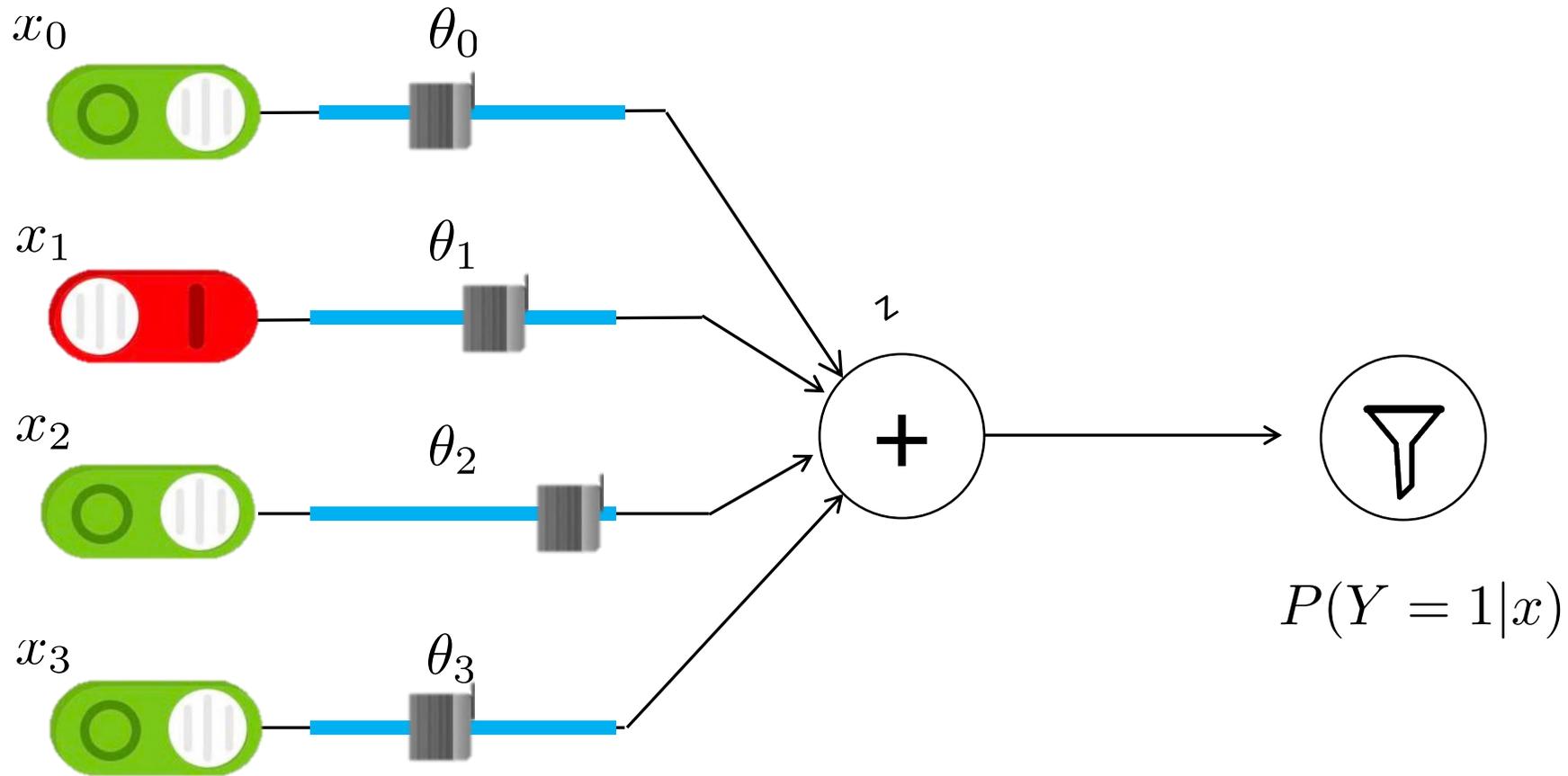# Biological Basis for Neural Networks

A neuron



Your brain



**Actually, it's probably someone else's brain**

# Beyond Classification

# Binary Prediction



$$P(Y = 1|\mathbf{X} = \mathbf{x}) = \sigma\left(\sum_i \theta_i x_i\right)$$

# Predict a Categorical



$P(Y = \text{"a"} | X = x)$

$P(Y = \text{"b"} | X = x)$

$P(Y = \text{"c"} | X = x)$

# Predict a Real Value (aka Regression)



y is distributed as a normal with this mean

$$\underset{\theta}{\mathrm{argmax}}\ LL(\theta) = \underset{\theta}{\mathrm{argmax}}\ -\sum_{i=1}^{n}\left(y^{(i)} - \theta^T x^{(i)}\right)^2$$

Hey it's the sum of squared errors!

niversity

# Core idea behind the revolution in AI

(aka Neural Networks)

**Deep learning** is (at its core) many logistic regression pieces stacked on top of each other.

# Digit Recognition Example

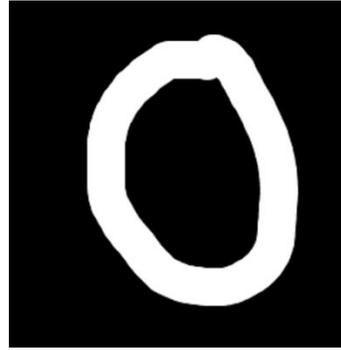Let's make feature vectors from pictures of numbers



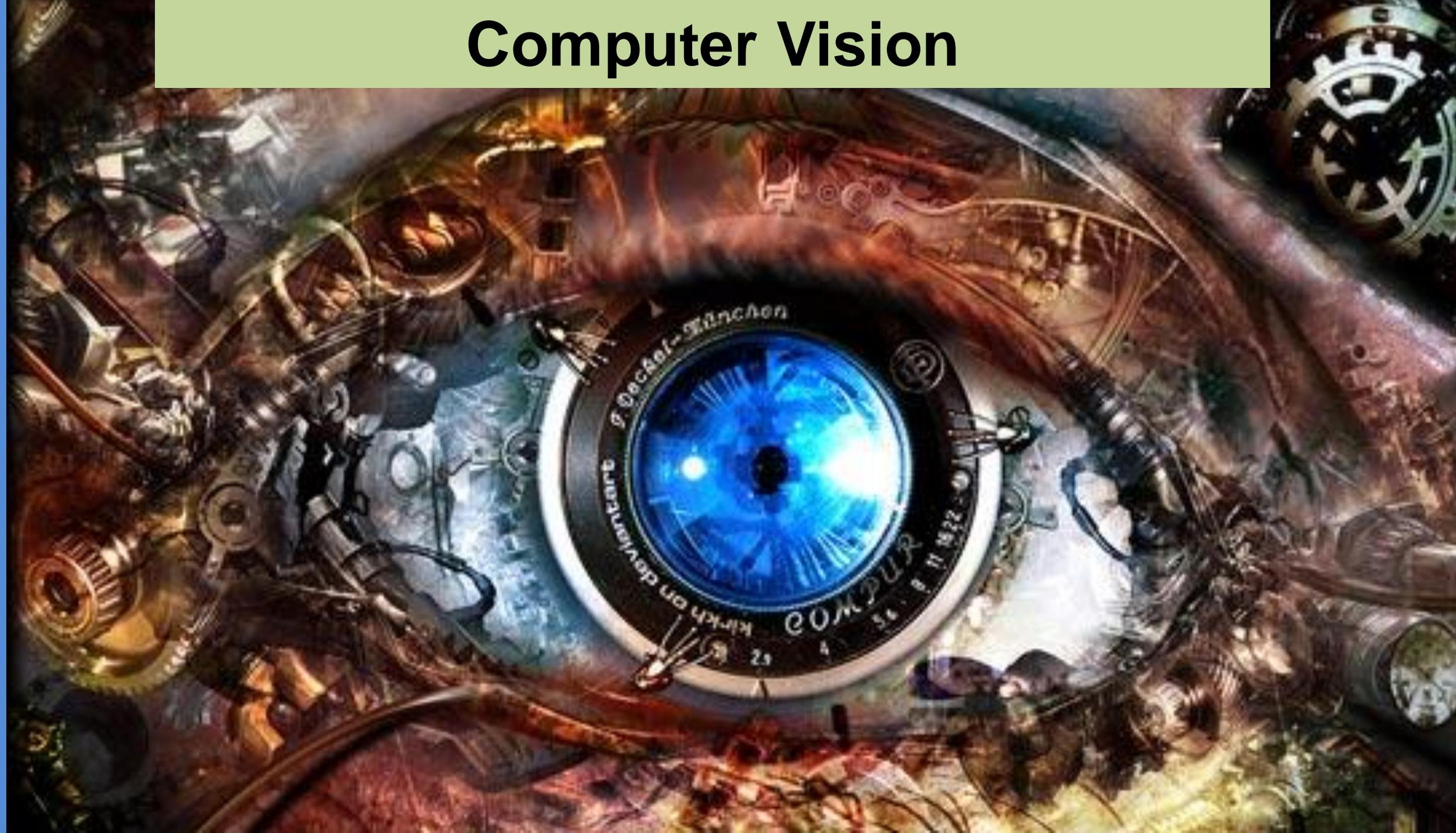$$\mathbf{x}^{(i)} = [0, 0, 0, 0, \ldots, 1, 0, 0, 1, \ldots 0, 0, 1, 0]$$
$$y^{(i)} = 0$$



$$\mathbf{x}^{(i)} = [0, 0, 1, 1, \ldots, 0, 1, 1, 0, \ldots 0, 1, 0, 0]$$
$$y^{(i)} = 1$$

# Computer Vision

# Vision in your Brain

Hundreds of millions of neurons [1]

Visual neurons make up up 30% of your cortex [1]

[1] http://discovermagazine.com/1993/jun/thevisionthingma227

# Logistic Regression



This means it predicts a 0

# Logistic Regression



Indicates logistic regression connection

This means it predicts a 0

# Logistic Regression



This means it predicts a 1

# Not So Good



This means it predicts a 1

# We Can Put Neurons Together



This means it predicts a 0

# We Can Put Neurons Together



There is a parameter for every connection

This means it predicts a 0

Look at a single "hidden" neuron

# We Can Put Neurons Together



There is a parameter for every connection

This means it predicts a 0

Look at another "hidden" neuron

# We Can Put Neurons Together



This means it predicts a 0

# We Can Put Neurons Together



There is a parameter for every connection

This means it predicts a 0

Look at another neuron

# We Can Put Neurons Together



This means it predicts a 0

# We Can Put Neurons Together



*lots

# We Can Put Neurons Together

# We Can Put Neurons Together



*lots

# We Can Put Neurons Together



*lots

*logistic regression

# Deep learning

<u>def</u> **Deep learning** is
maximum likelihood estimation
with neural networks.

<u>def</u> A **neural network** is
(at its core) many logistic
regression pieces stacked on
top of each other.

$[1, 0, \ldots, 1]$
$\boldsymbol{x}$, input

LOL

Lots of Logistic
(regressions)

$\hat{y}$, output
$P(Y|\boldsymbol{X} = \boldsymbol{x})$

> 0.5?

Yes.
Predict 1

# Demonstration



https://adamharley.com/nn_vis/cnn/2d.html

🔑 Deep learning gets its *intelligence* from its thetas (aka its parameters)

EXPLORER

∨ 25
> MNIST
🐍 train.py

🐍 train.py ⟩ ⬡ main

```python
import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader
from torchvision import datasets
from torchvision.transforms import ToTensor
import matplotlib.pyplot as plt

def main():
    # get the data
    train_loader, test_loader = download_data()
    print(f"Training examples: {len(train_loader.dataset)}")
    print(f"Test examples: {len(test_loader.dataset)}")

    # a very simple and fast nn
    model = nn.Sequential(
        nn.Flatten(),
        nn.Linear(28*28, 1024),
        nn.Sigmoid(),
        nn.Linear(1024, 1024),
        nn.Sigmoid(),
        nn.Linear(1024, 10)
    )

    # define the loss and optimizer ONCE outside run_train
    loss_function = nn.CrossEntropyLoss()
    optimizer = optim.Adam(model.parameters(), lr=0.001)
```

OUTPUT   PORTS   COMMENTS   PROBLEMS   DEBUG CONSOLE   **TERMINAL**

```
Untrained accuracy: 9.58%
Training ...
Epoch 1, Test Accuracy: 94.16%
Epoch 2, Test Accuracy: 95.87%
Epoch 3, Test Accuracy: 96.80%
Epoch 4, Test Accuracy: 97.53%
Epoch 5, Test Accuracy: 97.64%
Epoch 6, Test Accuracy: 97.42%
Epoch 7, Test Accuracy: 97.93%
Epoch 8, Test Accuracy: 97.95%
```

Ln 17, Col 22    Spaces: 4    UTF-8    LF    {} Python    3.12.3

Stanford University   51

# How do we train?

# MLE of Thetas!

# First: Learning Goals…

# 1. Understand Chain Rule as ♡ of Deep Learning

# 2. Demystify:
# Deep Learning is MLE

# 3. Become experts of logistic regression

# Math worth knowing:

# New Notation

Layer **x**        Layer **h**        Layer **ŷ**

# New Notation

Layer $\mathbf{x}$        Layer $\mathbf{h}$        Layer $\hat{\mathbf{y}}$



$\theta^{(h)}_{i,j}$

$\mathbf{x}_i$

$\mathbf{h}_j$

$$\mathbf{h}_j = \sigma \left( \sum_{i=0}^{m_x} \mathbf{x}_i \theta^{(h)}_{i,j} \right)$$

# New Notation

Layer $\mathbf{x}$

Layer $\mathbf{h}$

Layer $\hat{\mathbf{y}}$



$$\hat{y} = \sigma\left(\sum_{j=0}^{m_h} \mathbf{h}_j \theta_j^{(\hat{y})}\right)$$

# Forward Pass

Layer $\mathbf{x}$　　　　Layer $\mathbf{h}$　　　Layer $\hat{\mathbf{y}}$

# Forward Pass

Layer **x**          Layer **h**          Layer **ŷ**

# Forward Pass

Layer $\mathbf{x}$    Layer $\mathbf{h}$    Layer $\hat{\mathbf{y}}$



$$\mathbf{h}_j = \sigma\left(\sum_{i=0}^{m_x} \mathbf{x}_i \theta_{i,j}^{(h)}\right)$$

# Forward Pass

Layer $\mathbf{x}$    Layer $\mathbf{h}$    Layer $\hat{\mathbf{y}}$

$$LL(\theta) = y \log \hat{y}$$
$$+ (1-y) \log[1 - \hat{y}]$$

$$\hat{y} = \sigma \left( \sum_{j=0}^{m_h} \mathbf{h}_j \theta_j^{(\hat{y})} \right)$$

$$\mathbf{h}_j = \sigma \left( \sum_{i=0}^{m_x} \mathbf{x}_i \theta_{i,j}^{(h)} \right)$$

# All Together

Neural Network

$\mathbf{x}$ $\mathbf{h}$ $\hat{\mathbf{y}}$

$\theta^{(h)}$ $\theta^{(\hat{y})}$

$|\mathbf{x}| = 40$

$|\mathbf{h}| = 20$

How many parameters in $\theta^{(\hat{y})}$ ?

a) 2      b) 20      c) 40      d) 800

# Smoke Check 2

Neural Network

$\mathbf{x}$ $\quad$ $\mathbf{h}$ $\quad$ $\hat{\mathbf{y}}$

$\theta^{(h)}$ $\quad$ $\theta^{(\hat{y})}$

$|\mathbf{x}| = 40$

$|\mathbf{h}| = 20$

How many parameters in $\theta^{(h)}$ ?

a) 2 $\qquad$ b) 20 $\qquad$ c) 40 $\qquad$ d) 800

# Smoke Check 3

Neural Network

$\mathbf{x}$ $\qquad$ $\mathbf{h}$ $\qquad$ $\hat{\mathbf{y}}$

$\theta^{(h)}$ $\qquad$ $\theta^{(\hat{y})}$

$|\mathbf{x}| = 40$

$|\mathbf{h}| = 20$

How many parameters in total?

a) 800      b) 20      c) 820      d) 16000

# Today: Do Something Brave

# Forward Pass

Layer **x**          Layer **h**          Layer **ŷ**

800 parameters need setting

20 parameters need setting

# Only Have to Do Three Things

**1** Make deep learning assumption

**2** Calculate the log probability for all data

**3** Get partial derivative of log likelihood with respect to each theta

( 3 )  Get partial derivative of log likelihood with respect to each theta

# Why?

# Why We Calculate Partial Derivatives

A deep learning model gets its **intelligence** by having **useful thetas**.

We can find **useful thetas**, by searching for ones that **maximize likelihood** of our training data

We can **maximize likelihood** using **optimization techniques** (such as gradient ascent).

In order to use **optimization techniques**, we need to calculate the **partial derivative** of likelihood with respect to thetas.

Basically MLE is hard because it has so many details

Okay gang, let's see what deep learning really is.

Thanks to Keith Eicher

Convex Optimization??

# Only Have to Do Three Things

**1** Make deep learning assumption

$$P(Y = 1|X = \mathbf{x}) = \hat{y}$$
$$P(Y = 0|X = \mathbf{x}) = 1 - \hat{y}$$

**2** Calculate the log probability for all data

# Same Assumption, Same LL

$$P(Y = 1 | X = \mathbf{x}) = \hat{y}$$

$$\hat{y} = \sigma\left(\sum_{j=0}^{m_h} \mathbf{h}_j \theta_j^{(\hat{y})}\right) \qquad \mathbf{h}_j = \sigma\left(\sum_{i=0}^{m_x} \mathbf{x}_i \theta_{i,j}^{(h)}\right)$$

*For one datum*

$$P(Y = y | \mathbf{X} = \mathbf{x}) = (\hat{y})^y (1 - \hat{y})^{1-y}$$

*Feel the Bern!*
$$Y \sim \mathrm{Bern}(\hat{y})$$

*For IID data*

$$L(\theta) = \prod_{i=1}^{n} P(Y = y^{(i)} | X = \mathbf{x}^{(i)})$$

$$= \prod_{i=1}^{n} (\hat{y}^{(i)})^{y^{(i)}} \cdot \left[1 - (\hat{y}^{(i)})\right]^{(1 - y^{(i)})}$$

*Take the log*

$$LL(\theta) = \sum_{i=1}^{n} y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log[1 - \hat{y}^{(i)}]$$

# Only Have to Do Three Things

(1) Make deep learning assumption

$$\hat{y} = \sigma \left( \sum_{j=0}^{m_h} \mathbf{h}_j \theta_j^{(\hat{y})} \right)$$

$$P(Y = 1 | X = \mathbf{x}) = \hat{y}$$

$$P(Y = 0 | X = \mathbf{x}) = 1 - \hat{y}$$

(2) Calculate the log probability for all data

$$LL(\theta) = \sum_{i=0}^{n} y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log[1 - \hat{y}^{(i)}]$$

(3) Get partial derivative of log likelihood with respect to each theta

# Derivative Goals

Loss with respect to output layer params

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}}$$

Loss with respect to hidden layer params

$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}}$$



Neural Network

$\mathbf{x}$  $\mathbf{h}$  $\hat{\mathbf{y}}$

$\theta^{(h)}$  $\theta^{(\hat{y})}$

# Bad Approach

$$LL(\theta) = y \log \hat{y} + (1 - y) \log[1 - \hat{y}]$$

---

$$\hat{y} = \sigma \left( \sum_{i=0}^{m_h} \mathbf{h}_i \theta_i^{(\hat{\mathbf{y}})} \right)$$

Math bug

$$= \sigma \left( \sum_{i=0}^{m_h} \left[ \sigma \left( \sum_{i=0}^{m_x} \mathbf{x}_i \theta_{i,j}^{(\mathbf{h})} \right) \right] \theta_i^{(\hat{\mathbf{y}})} \right)$$

Neural Network

$\mathbf{x}$ $\mathbf{h}$ $\hat{\mathbf{y}}$

$\theta^{(h)}$ $\theta^{(\hat{y})}$

# Derivatives Without Tears

# Big Idea #1: Chain Rule

Woah Mr Blanton, you were right.
Chain rule is useful!

$$\frac{\partial f(z)}{\partial x} = \frac{\partial f(z)}{\partial z} \cdot \frac{\partial z}{\partial x}$$

First use:

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}} = \frac{\partial LL}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \theta_i^{(\hat{y})}}$$

# Big Idea #2: Sigmoid Derivative

True fact about sigmoid functions

$$\frac{\partial}{\partial z}\sigma(z) = \sigma(z)[1 - \sigma(z)]$$

# Big Idea #3: Derivative of Sum

$$LL(\theta) = \sum_{i=0}^{n} y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log[1 - \hat{y}^{(i)}]$$

We only need to calculate the gradient
for one training example!

$$\frac{\partial}{\partial x} \sum f(x) = \sum \frac{\partial}{\partial x} f(x)$$

We will pretend we only have one
example

$$LL(\theta) = y \log \hat{y} + (1 - y) \log[1 - \hat{y}]$$

We can sum up the gradients of each
example to get the correct answer

# Recall

# Sigmoid has a Beautiful Slope

$$\frac{\partial}{\partial \theta_j} \sigma(\theta^T x)?$$

$$\frac{\partial}{\partial z} \sigma(z) = \sigma(z)[1 - \sigma(z)]$$

where $z = \theta^T x$

$$\frac{\partial}{\partial \theta_j} \sigma(\theta^T x) = \frac{\partial}{\partial z} \sigma(z) \cdot \frac{\partial z}{\partial \theta_j}$$

Chain rule!

$$\frac{\partial}{\partial \theta_j} \sigma(\theta^T x) = \sigma(\theta^T x)[1 - \sigma(\theta^T x)]x_j$$

Plug and chug

Sigmoid, you should be a ski hill
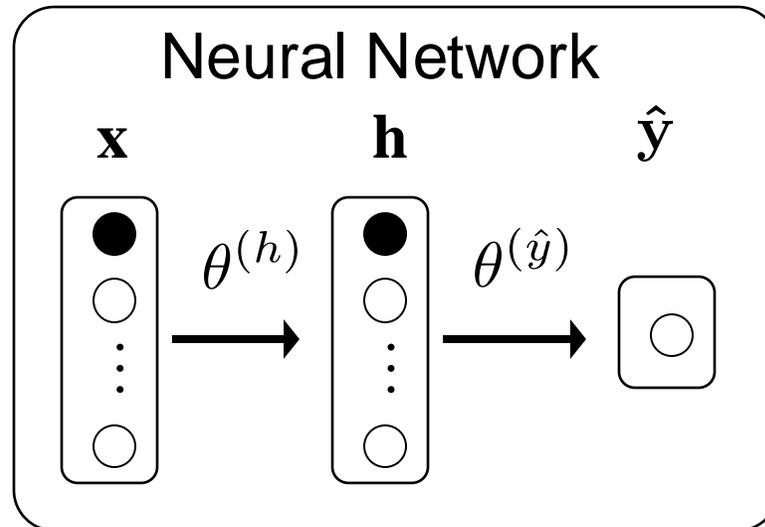
This is ~~Sparta~~!!!!!

↑
Stanford

# Derivative Goals

Loss with respect to output layer params

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}}$$
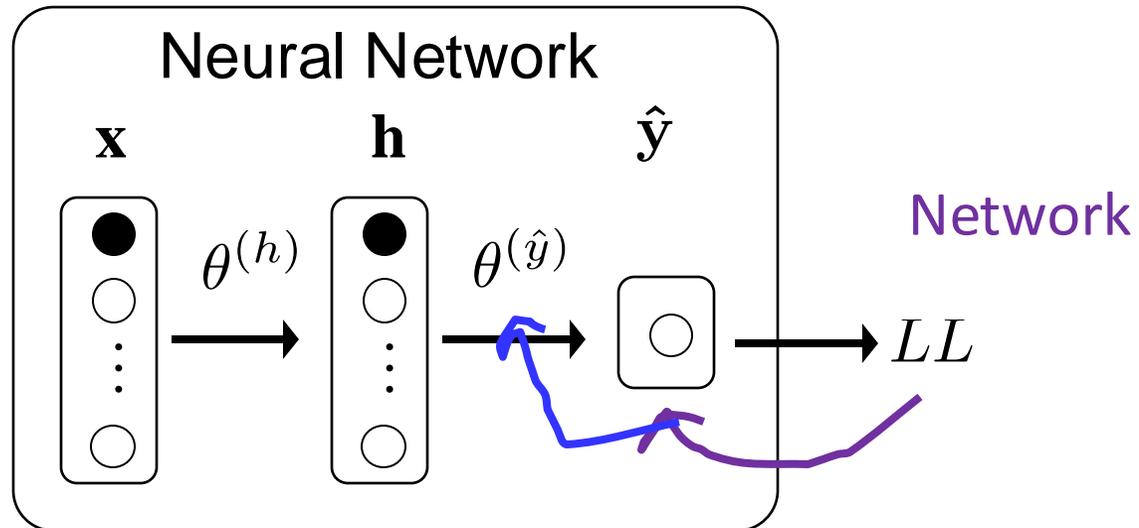
Loss with respect to hidden layer params

$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}}$$



Neural Network

$\mathbf{x}$ $\qquad$ $\mathbf{h}$ $\qquad$ $\hat{\mathbf{y}}$

$\theta^{(h)}$ $\qquad$ $\theta^{(\hat{y})}$

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}}$$

Goal



Neural Network

$\mathbf{x}$  $\mathbf{h}$  $\hat{\mathbf{y}}$

$\theta^{(h)}$  $\theta^{(\hat{y})}$  $LL$

Network

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}} = \frac{\partial LL}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \theta_i^{(\hat{y})}}$$

Decomposition

# Decomposition
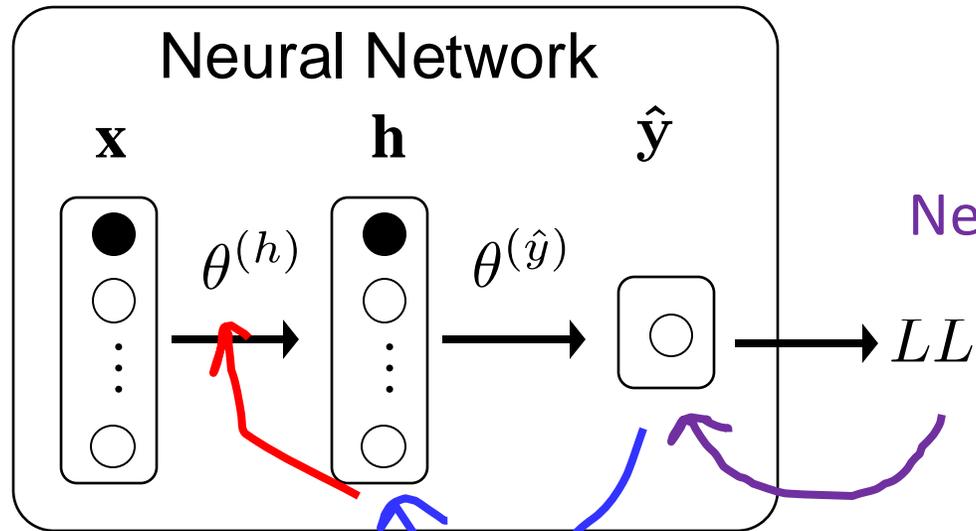
# Gradient of output layer params

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}} = \frac{\partial LL}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \theta_i^{(\hat{y})}}$$

---

# Gradient of output layer params
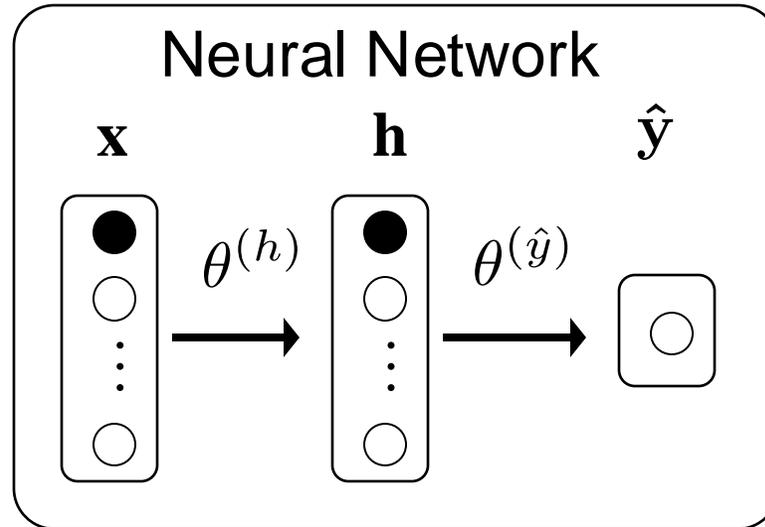
$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}} = \boxed{\frac{\partial LL}{\partial \hat{y}}} \cdot \frac{\partial \hat{y}}{\partial \theta_i^{(\hat{y})}}$$

---

$$LL(\theta) = y \log \hat{y} + (1 - y) \log[1 - \hat{y}]$$

$$\frac{\partial LL(\theta)}{\partial \hat{y}} = \frac{y}{\hat{y}} + \frac{(1 - y)}{(1 - \hat{y})} \cdot \frac{\partial (1 - \hat{y})}{\partial \hat{y}}$$

$$\frac{\partial LL(\theta)}{\partial \hat{y}} = \frac{y}{\hat{y}} - \frac{(1 - y)}{(1 - \hat{y})}$$

# Gradient of output layer params

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}} = \boxed{\frac{\partial LL}{\partial \hat{y}}} \cdot \boxed{\frac{\partial \hat{y}}{\partial \theta_i^{(\hat{y})}}}$$

$$\hat{y} = \sigma\left(\sum_{j=0}^{m_h} \mathbf{h}_j \theta_j^{(\hat{y})}\right) = \sigma(z) \qquad \text{where} \qquad z = \sum_{j=0}^{m_h} \mathbf{h}_j \theta_j^{(\hat{y})}$$

$$\frac{\partial \hat{y}}{\partial \theta_i^{(\hat{y})}} = \hat{y}[1 - \hat{y}] \cdot \frac{\partial}{\partial \theta_i^{(\hat{y})}} \sum_{j=0}^{m_h} \mathbf{h}_j \theta_j^{(\hat{y})}$$

$$= \hat{y}[1 - \hat{y}] \cdot h_i$$
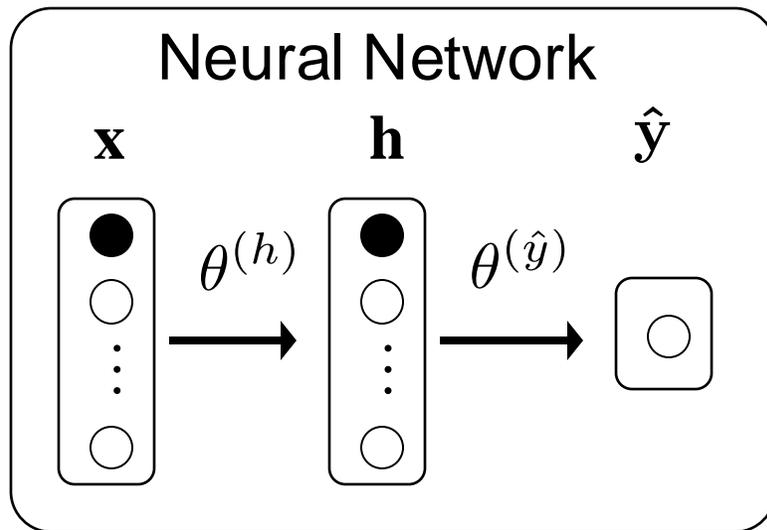
*What! That's not scary!*

# Make it Simple

$$\frac{\partial LL(\theta)}{\partial \theta_i^{(\hat{y})}} = \text{<image>} \cdot \text{<image>}$$

$$\text{<image>} = \frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})}$$

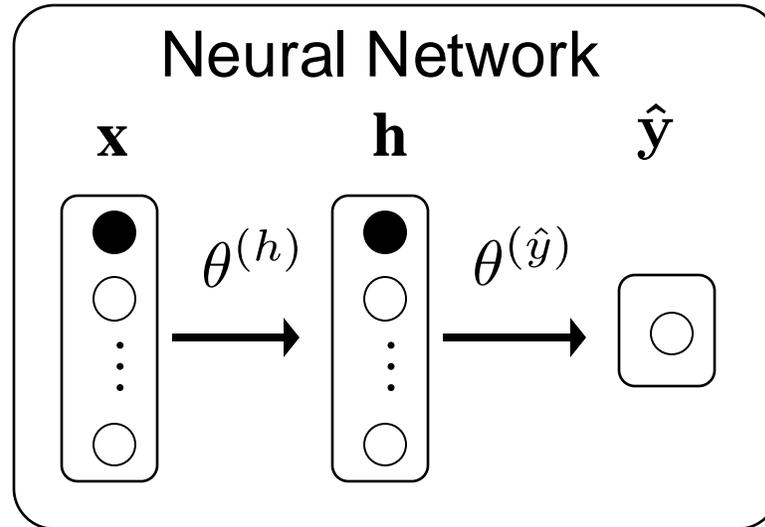$$\text{<image>} = \hat{y}[1-\hat{y}] \cdot h_i$$

# Boom!

$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}}$$

# Gradient of hidden layer params

$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}} = \frac{\partial LL}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \mathbf{h}_j} \cdot \frac{\partial \mathbf{h}_j}{\partial \theta_{i,j}^{(h)}}$$

# Gradient of hidden layer params

$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}} = \boxed{\frac{\partial LL}{\partial \hat{y}}} \cdot \boxed{\frac{\partial \hat{y}}{\partial \mathbf{h}_j}} \cdot \frac{\partial \mathbf{h}_j}{\partial \theta_{i,j}^{(h)}}$$

$$\hat{y} = \sigma\left(\sum_{i=0}^{m_h} \mathbf{h}_i \theta_i^{(\hat{y})}\right)$$

$$\frac{\partial \hat{y}}{\partial \mathbf{h}_j} = \hat{y}[1 - \hat{y}]\theta_j^{(\hat{y})}$$

Wait is it over?

# Gradient of hidden layer params

$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}} = \boxed{\frac{\partial LL}{\partial \hat{y}}} \cdot \boxed{\frac{\partial \hat{y}}{\partial \mathbf{h}_j}} \cdot \boxed{\frac{\partial \mathbf{h}_j}{\partial \theta_{i,j}^{(h)}}}$$

$$\mathbf{h}_j = \sigma \left( \sum_{k=0}^{m_x} \mathbf{x}_k \theta_{k,j} \right)$$

$$\frac{\partial \mathbf{h}_j}{\partial \theta_{i,j}^{(h)}} = \mathbf{h}_j [1 - \mathbf{h}_j] \mathbf{x}_i$$

That one too?

# Make it Simple

$$\frac{\partial LL(\theta)}{\partial \theta_{i,j}^{(h)}} = $$



$$= \frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})}$$

$$= \hat{y}[1 - \hat{y}]\theta_j^{(\hat{y})}$$

$$= \mathbf{h}_j[1 - \mathbf{h}_j]\mathbf{x}_j$$

Congrats. You now know Backpropagation

# Moment of silence

# What Would You Do Here?



Bigger Neural Network

$\mathbf{x}$ $\quad$ $\mathbf{h}^{(1)}$ $\quad$ $\mathbf{h}^{(2)}$ $\quad$ $\hat{\mathbf{y}}$

$\theta^{(1)}$ $\quad$ $\theta^{(2)}$ $\quad$ $\theta^{(\hat{y})}$
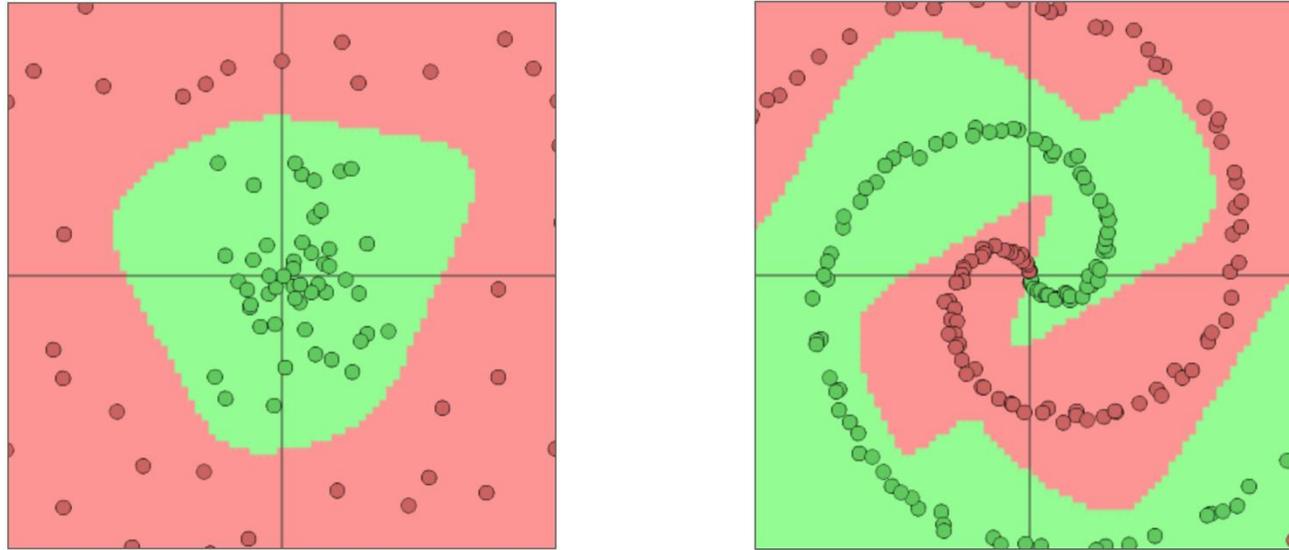
Chain rule:
Game changer for
artificial intelligence

# Neural Networks Can Learn Complex Functions

- Some data sets/functions are not separable



- These are classifiers learned by neural networks

http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html

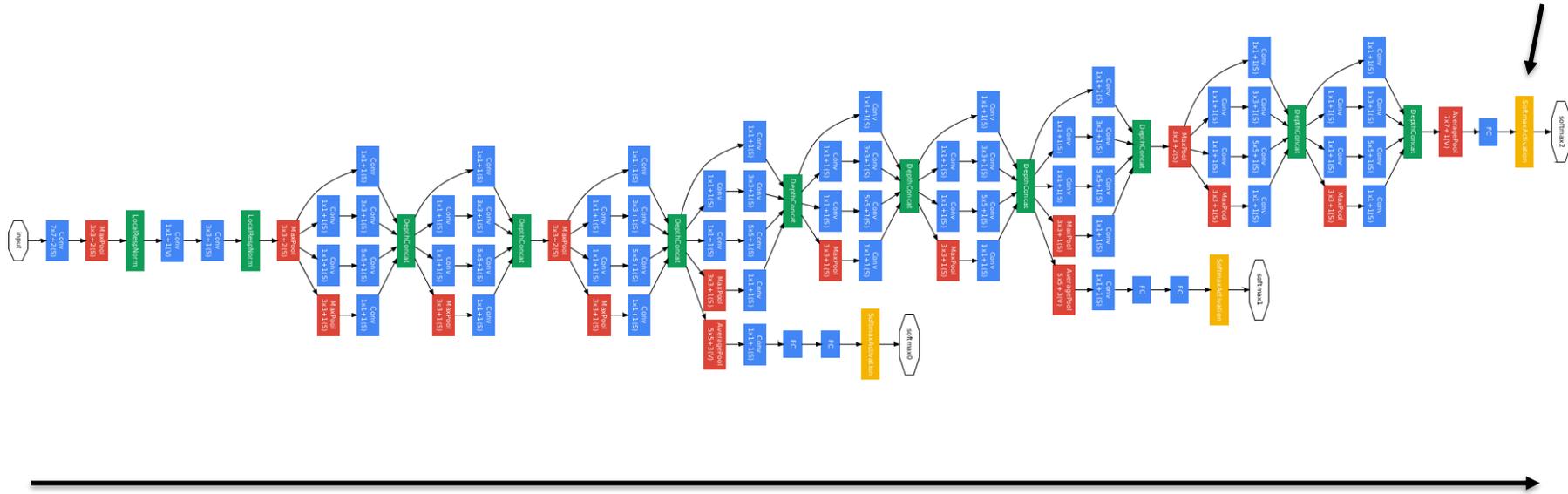# Works for any number of layers

# GoogLeNet Brain



1 Trillion Artificial Neurons

# GoogLeNet Brain



Multiple,
Multi class output

22 layers deep

Piech

# The Cat Neuron



Top stimuli from the test set

Optimal stimulus
by numerical optimization

Le, et al., *Building high-level features using large-scale unsupervised*

Hire the smartest people in the world

Invent cat detector

# Best Neuron Stimuli



Neuron 1

Neuron 2

Neuron 3

Neuron 4

Neuron 5

Le, et al., *Building high-level features using large-scale unsupervised*

# Best Neuron Stimuli



Neuron 6

Neuron 7

Neuron 8

Neuron 9

Le, et al., *Building high-level features using large-scale unsupervised*

# Best Neuron Stimuli

Neuron 10

Neuron 11

Neuron 12

Neuron 13

Le, et al., *Building high-level features using large-scale unsupervised*

# ImageNet Classification

22,000 categories

14,000,000 images

Hand-engineered features (SIFT, HOG, LBP),
Spatial pyramid,  SparseCoding/Compression

Le, et al., *Building high-level features using large-scale unsupervised*

# 22,000 is a lot!

...
smoothhound, smoothhound shark, Mustelus mustelus
American smooth dogfish, Mustelus canis
Florida smoothhound, Mustelus norrisi
whitetip shark, reef whitetip shark, Triaenodon obseus
Atlantic spiny dogfish, Squalus acanthias
Pacific spiny dogfish, Squalus suckleyi
hammerhead, hammerhead shark
smooth hammerhead, Sphyrna zygaena
smalleye hammerhead, Sphyrna tudes
shovelhead, bonnethead, bonnet shark, Sphyrna tiburo
angel shark, angelfish, Squatina squatina, monkfish
electric ray, crampfish, numbfish, torpedo
smalltooth sawfish, Pristis pectinatus
guitarfish
roughtail stingray, Dasyatis centroura
butterfly ray
eagle ray
spotted eagle ray, spotted ray, Aetobatus narinari
cownose ray, cow-nosed ray, Rhinoptera bonasus
manta, manta ray, devilfish
Atlantic manta, Manta birostris
devil ray, Mobula hypostoma
grey skate, gray skate, Raja batis
little skate, Raja erinacea
...

Stingray

Mantaray

# 0.005%

Random guess

# 1.5%

Pre Neural Networks

# ?

GoogLeNet

Le, et al., *Building high-level features using large-scale unsupervised*

0.005%     1.5%     43.9%

Random guess    Pre Neural Networks    GoogLeNet

Szegedy et al, Going Deeper With Convolutions, CVPR 2015
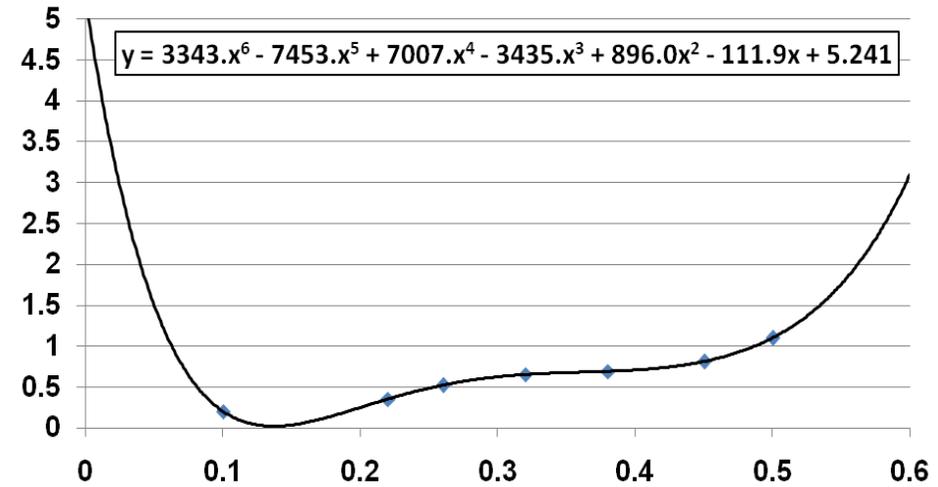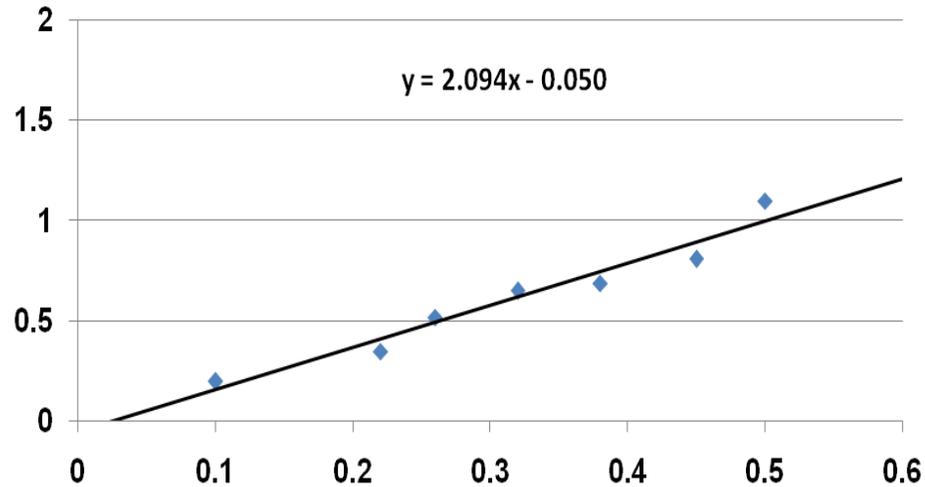
0.005%     1.5%     95.1%

Random guess     Pre Neural Networks     SE-ResNet

How many parameters
is too many?

# Good ML = Generalization

- Goal of machine learning: build models that **_generalize_** well to predicting new data

    - "Overfitting": fitting the training data too well, so we lose generality of model

$$y = 2.094x - 0.050$$

$$y = 3343.x^6 - 7453.x^5 + 7007.x^4 - 3435.x^3 + 896.0x^2 - 111.9x + 5.241$$

- Polynomial on the right fits training data perfectly!
- Which would you rather use to predict a new data point?

# Prevent Overfitting?

**Dropout** when your model is training, randomly turn off your neurons with probability 0.5. It will make your network more robust.



(a) Standard Neural Net

(b) After applying dropout.

# Shared Weights?

**Convolution** it turns out if you want to force some of your weights to be shared for different neurons, the math isn't that much harder. This is used a lot for vision (CNN).
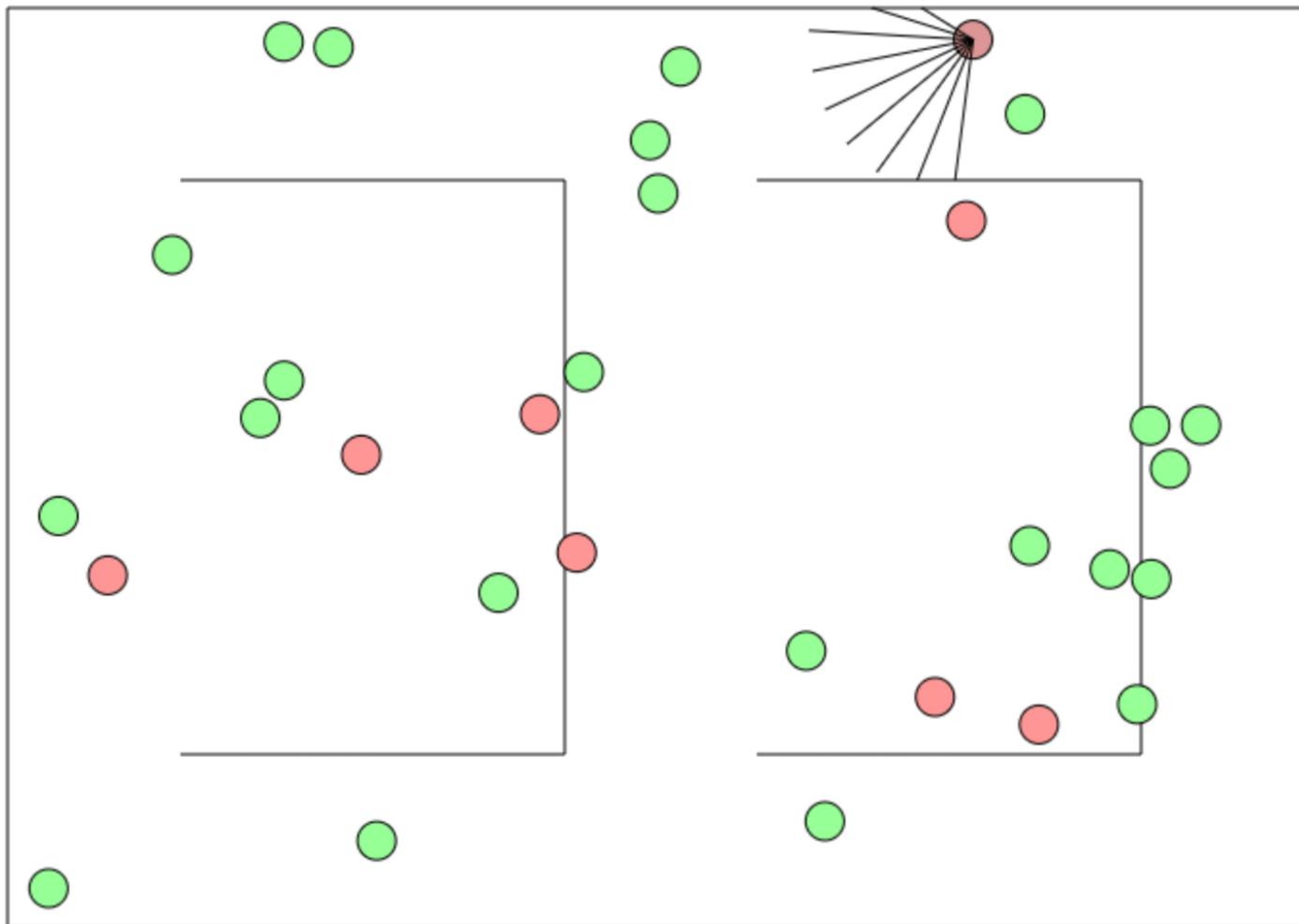
Not everything is classification

# Making Decisions?

**Deep Reinforcement Learning**
Instead of having the output of a model be a probability you can make it an expectation.

# Deep Reinforcement Learning



http://cs.stanford.edu/people/karpathy/convnetjs/demo/rldemo.html

# Deep Reinforcement Learning

R is a reward and $A_i$ is a legal action

Input is a representation of current state (S)

RL Neural Network

$\mathbf{x}$     $\mathbf{h}$     $\hat{\mathbf{y}}$

$\theta^{(h)}$     $\theta^{(\hat{y})}$

$E[R \mid A_1, S]$

$E[R \mid A_2, S]$

$E[R \mid A_3, S]$

Interpret outputs as expected reward for a given action

# Deep Mind Atari Games



| | |
|---|---|
| Video Pinball | 2539% |
| Boxing | 1707% |
| Breakout | 1327% |
| Star Gunner | 598% |
| Robotank | 508% |
| Atlantis | 449% |
| Crazy Climber | 419% |
| Gopher | 400% |
| Demon Attack | 294% |
| Name This Game | 278% |
| Krull | 277% |
| Assault | 246% |
| Road Runner | 232% |
| Kangaroo | 224% |
| James Bond | 145% |
| Tennis | 143% |
| Pong | 132% |
| Space Invaders | 121% |
| Beam Rider | 119% |
| Tutankham | 112% |
| Kung-Fu Master | 102% |
| Freeway | 102% |
| Time Pilot | 100% |
| Enduro | 97% |
| Fishing Derby | 93% |
| Up and Down | 92% |
| Ice Hockey | 79% |
| Q*bert | 78% |
| H.E.R.O. | 76% |

DQN

Score compared to best human

# Can A.I. Grade Your Next Test?

Neural networks could give online education a boost by providing automated feedback to students.

# Invented the Proto-Transformer

Rubric Level Accuracy on Few-Shot Grading a Novel Question

# Gave Feedback to 3,500 Real Students

Do you agree? AI feedback **97.9%**. Human feedback 96.7%



Algorithm uses attention to highlight where in the code the error comes from

AI generated feedback

Students evaluate the feedback

Syntax error (missing ") here would prevent auto graders from being useful.

**Rubric**
Mike

**Baselines**
Supervised
* Autograders

**Humans**
Stanford TAs
Code in Place TAs

**Historical Exams**

**Deployment**

**Validation**

Demographic Accuracy
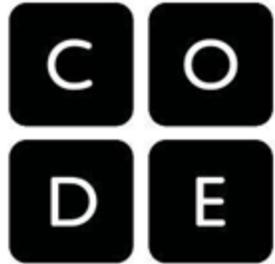
Teacher Rating

Constructive Feedback Rate

# But what about interactive, creative assignments?



1M ungraded code.org assignments.

The AI is shown a brand new student game. Does it work?

Simultaneously learn to grade and play to grade.

Majority class: 50%
Code-as-text:  67%
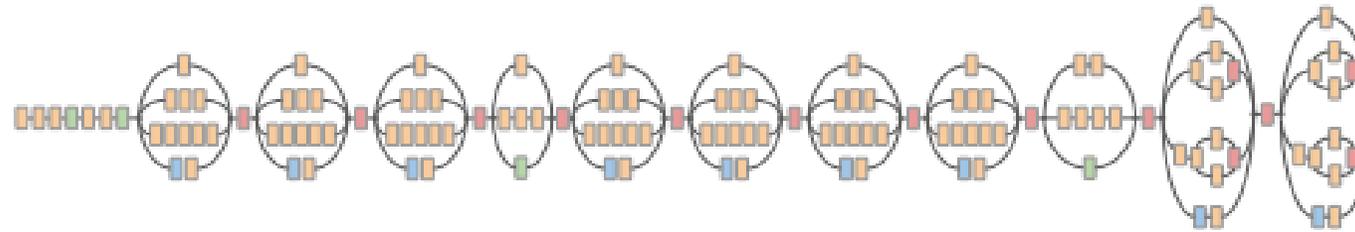Play-to-grade: **94%**
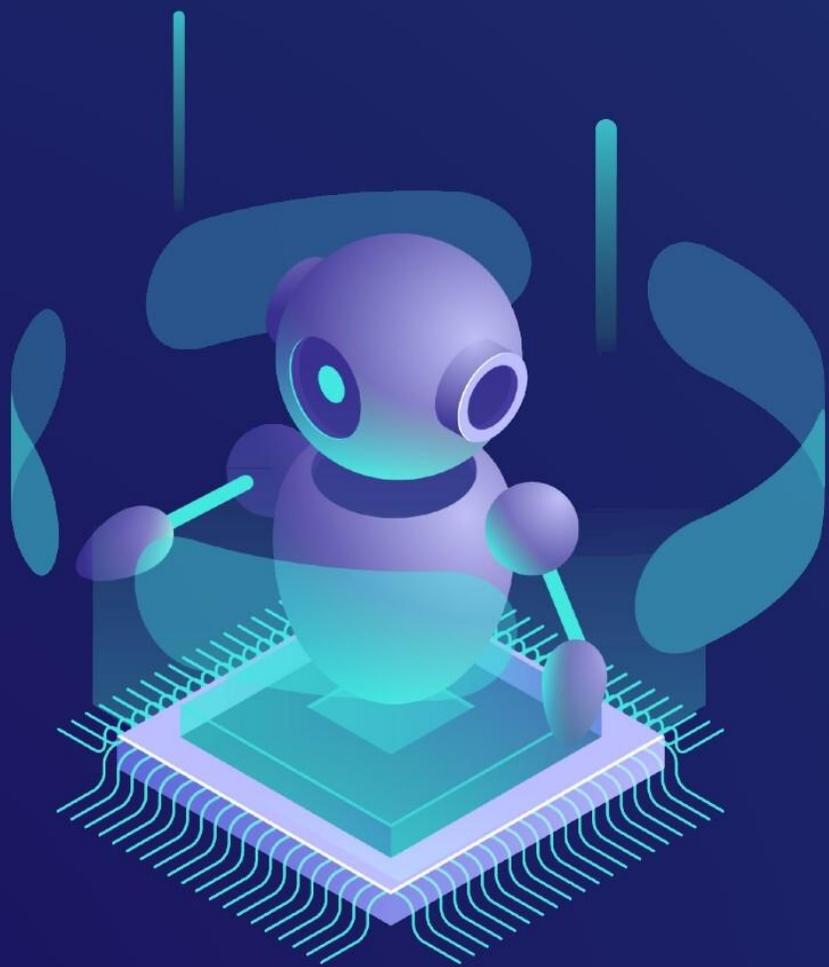
Piech

# Detecting skin cancer



Skin Lesion Image

Deep Convolutional Neural Network (Inception-v3)

Training Classes (757)

- Acral-lent. melanoma
- Amelanotic melanoma
- Lentigo melanoma
- ...
- Blue nevus
- Halo nevus
- Mongolian spot
- ...

Esteva, Andre, et al. "Dermatologist-level classification of skin cancer with deep neural networks." *Nature* 542.7639 (2017): 115-118.

Piech