

# 27: Intro to Deep Learning

---

Lisa Yan

November 13, 2020

# Quick slide reference

---

3	0.1% of Deep Learning	LIVE
41	Beyond the Basics	extra

# Deep Learning

# Innovations in deep learning



AlphaGO (2016)

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

Errata (misspoke):

- 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](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\\_\(chess\\_computer\)](https://en.wikipedia.org/wiki/Deep_Blue_(chess_computer))

# Computers making art



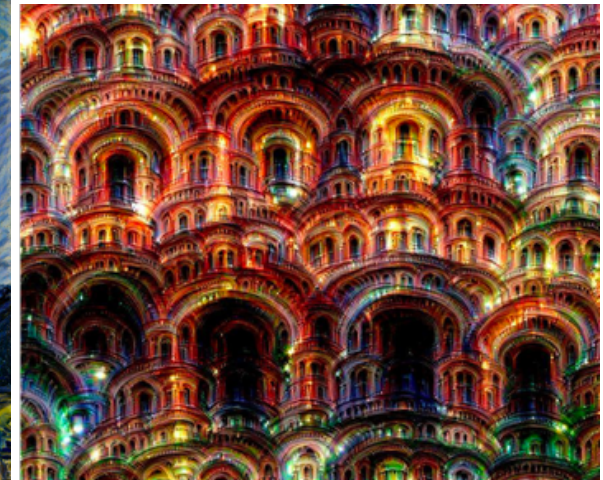
The Next Rembrandt

<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>

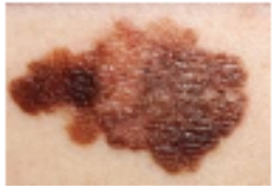


Google Deep Dream

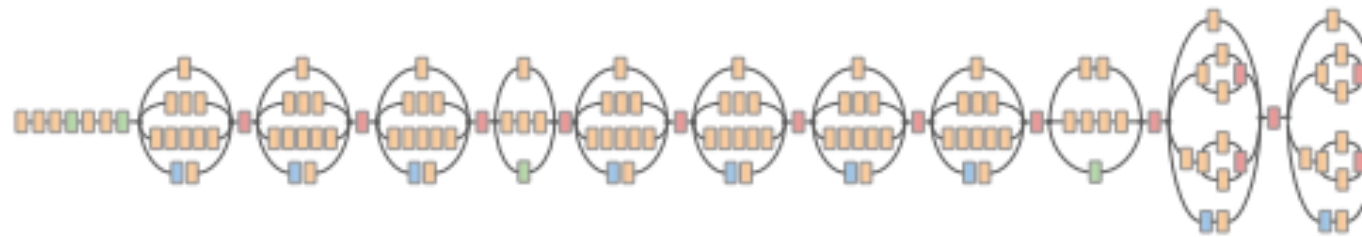
<https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>

# Detecting skin cancer

Skin Lesion Image



Deep Convolutional Neural Network (Inception-v3)



Training Classes (757)

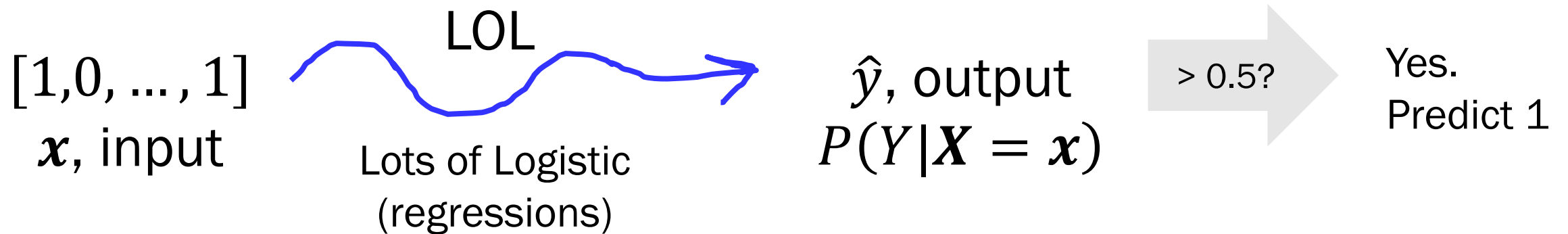
- Acral-lent. melanoma
- Amelanotic melanoma
- Lentigo melanoma
- ...
- Blue nevus
- Halo nevus
- Mongolian spot
- ...
- 
- 
- 

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

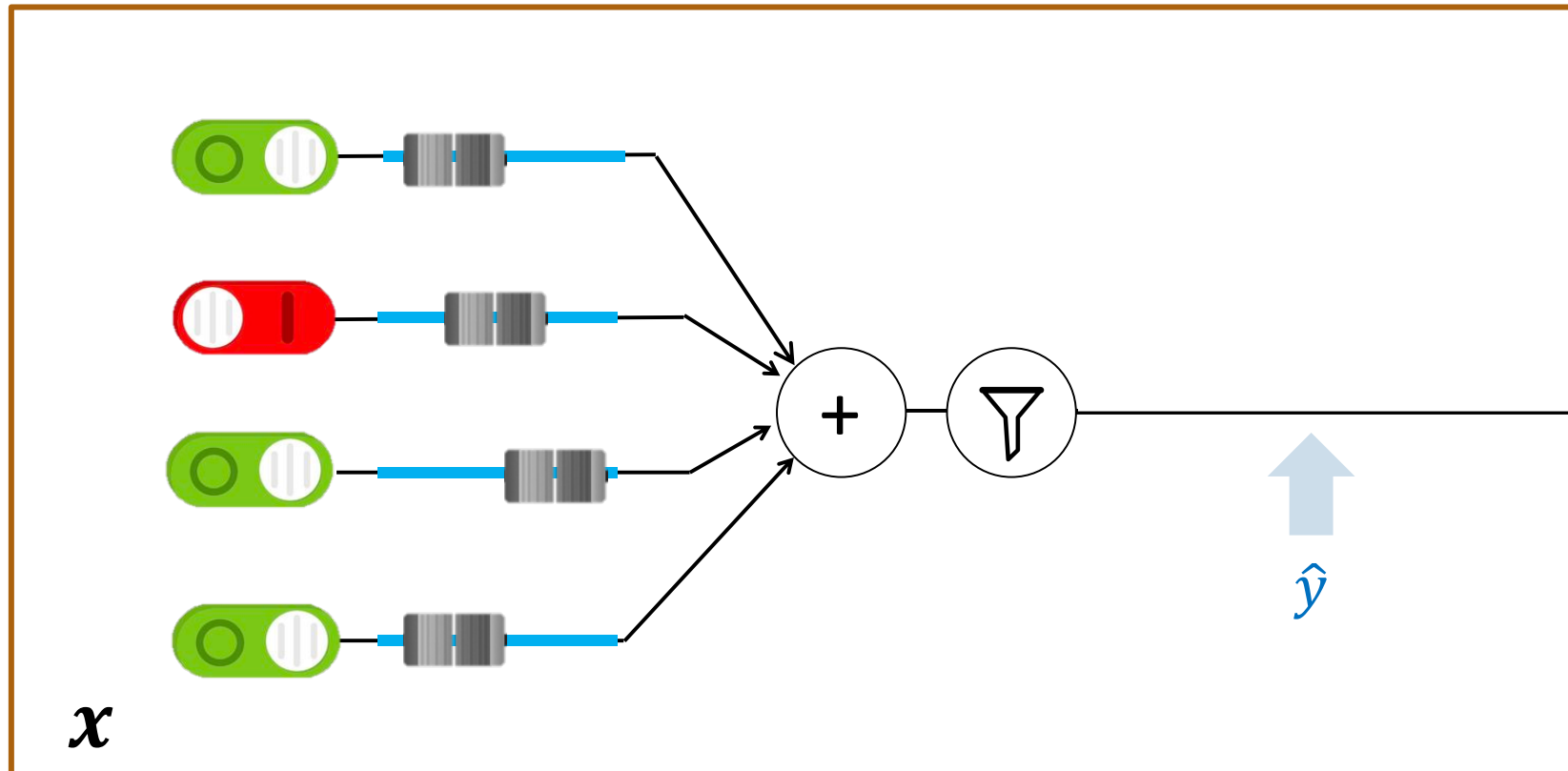
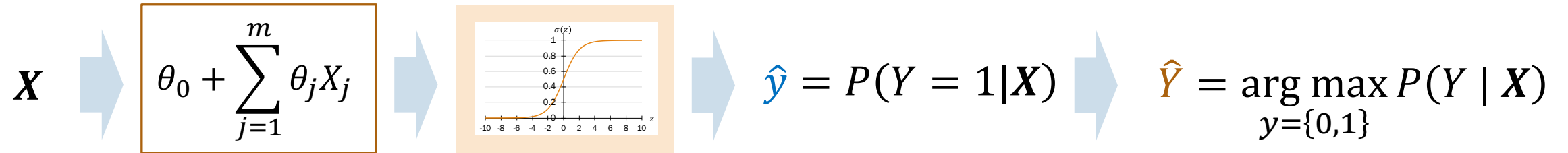
# Deep learning

def **Deep learning** is  
maximum likelihood estimation  
with neural networks.

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



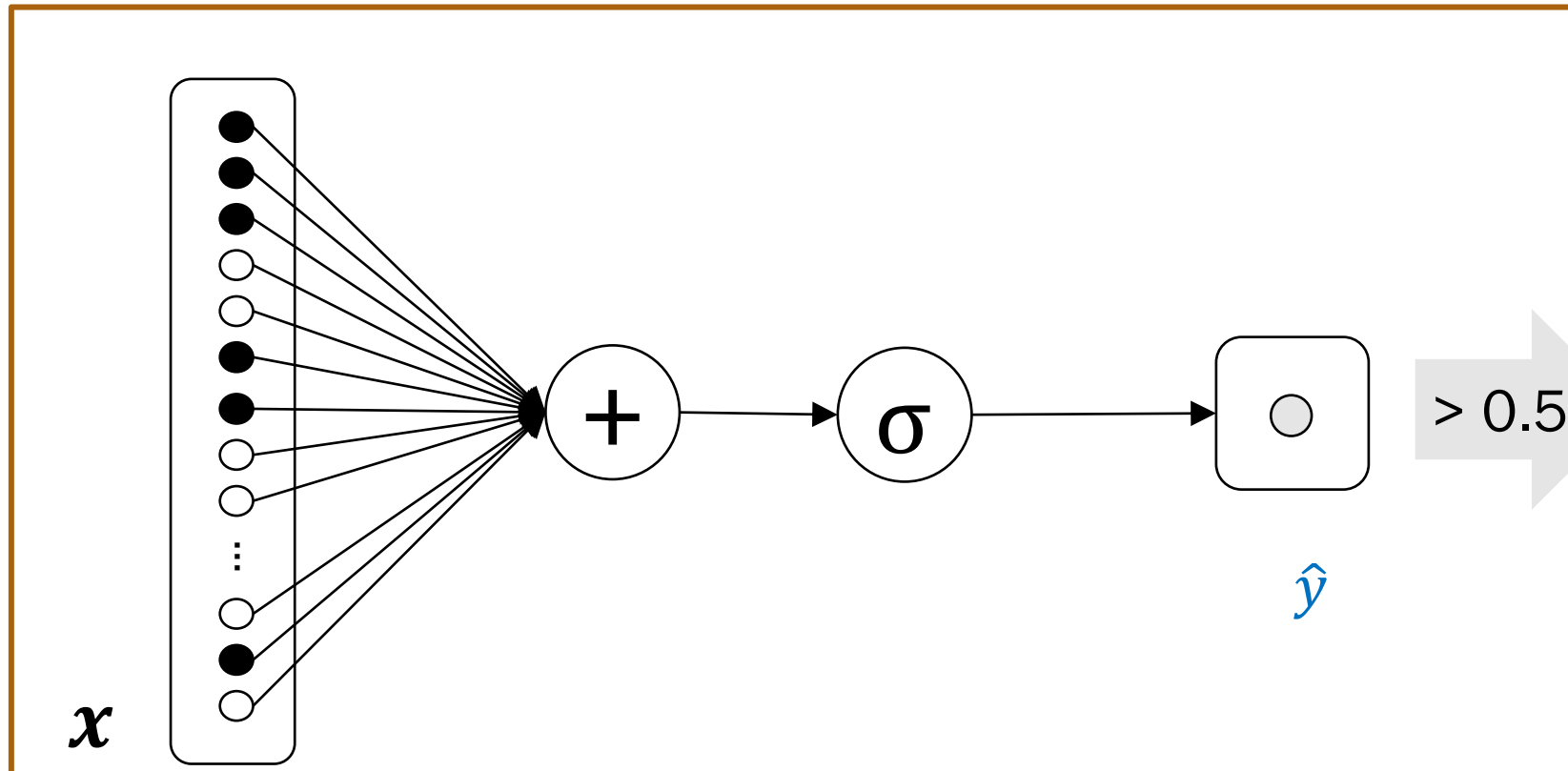
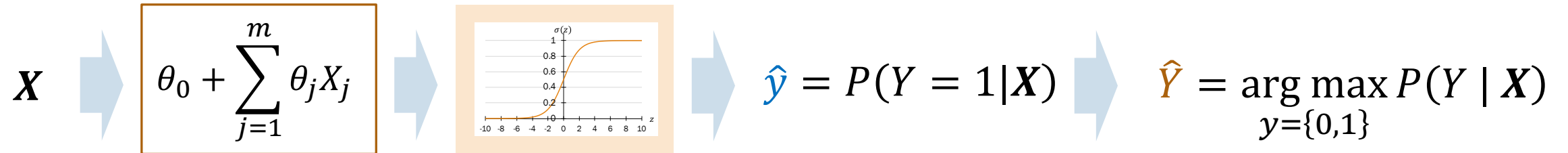
# Logistic Regression Model



Let's focus on the model up to  $\hat{y}$ .

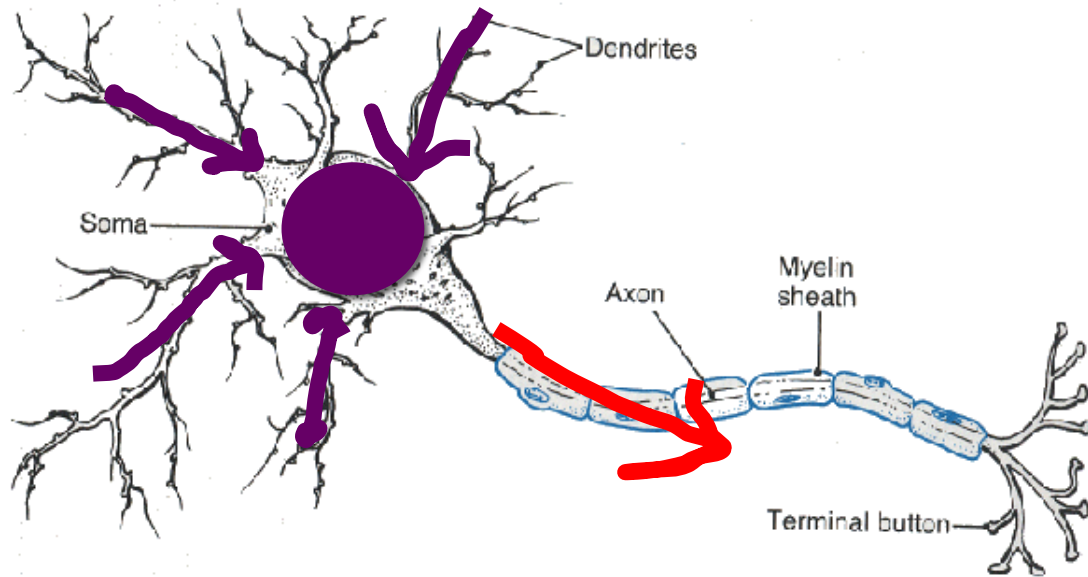


# Logistic Regression Model

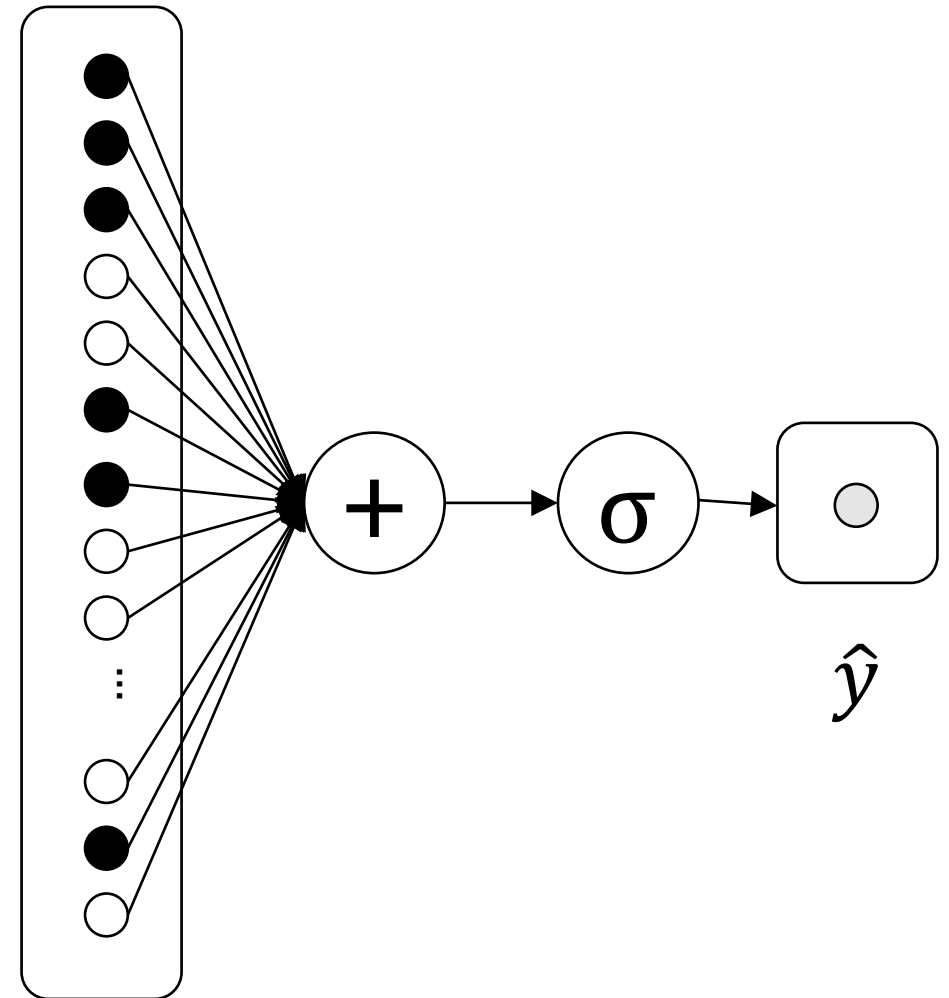


Let's focus on the model up to  $\hat{y}$ .

# One neuron = One logistic regression

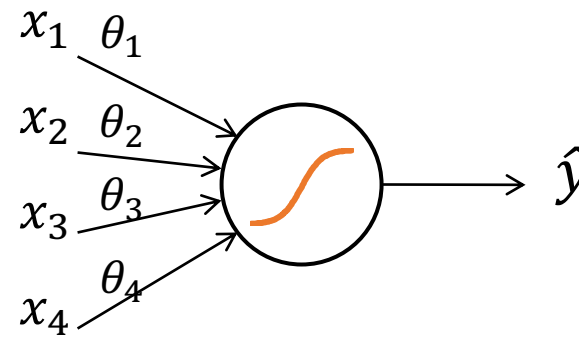
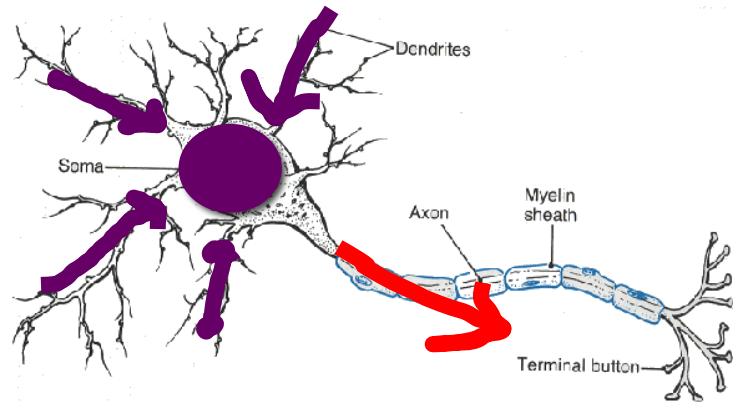


=



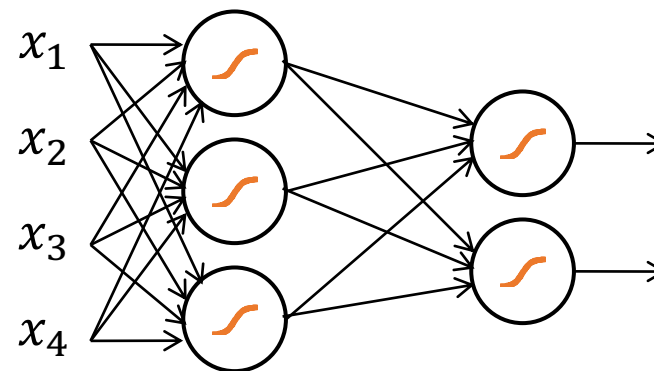
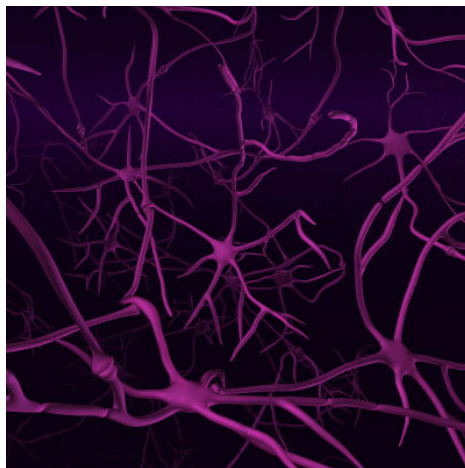
# Biological basis for neural networks

## A neuron



One neuron =  
one logistic  
regression

## Your brain



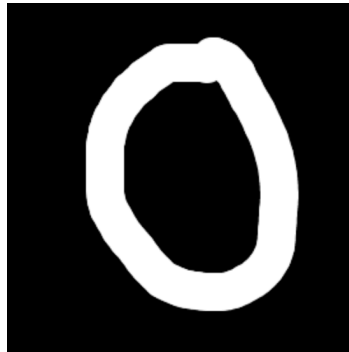
Neural network =  
many logistic  
regressions

(actually, probably someone else's brain)

# Digit recognition example

---

Input image



Input feature vector

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

Output label

$$y^{(i)} = 0$$

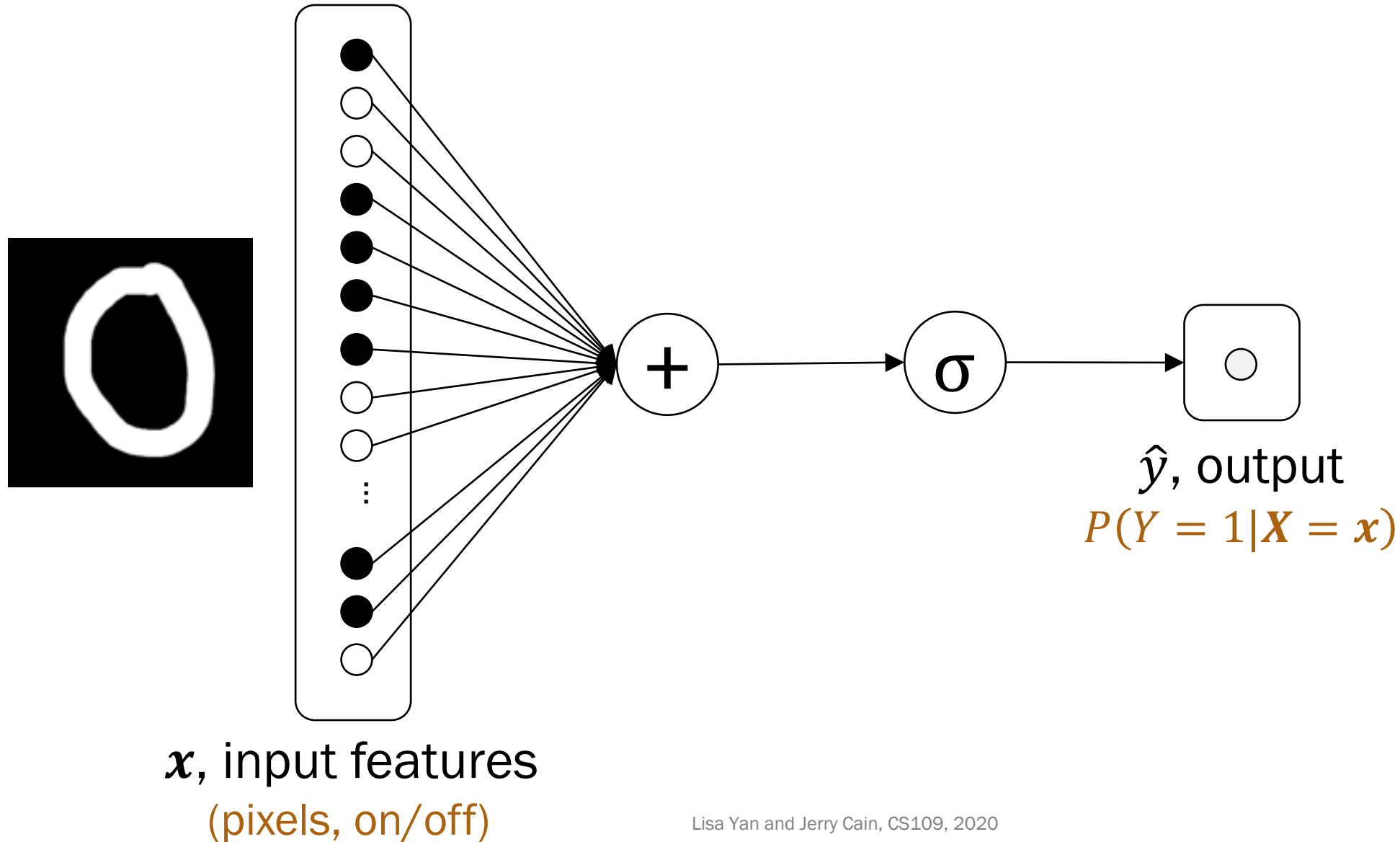


$$\mathbf{x}^{(i)} = [0,0,1,1, \dots, 0,1,1,0, \dots, 0,1,0,0]$$

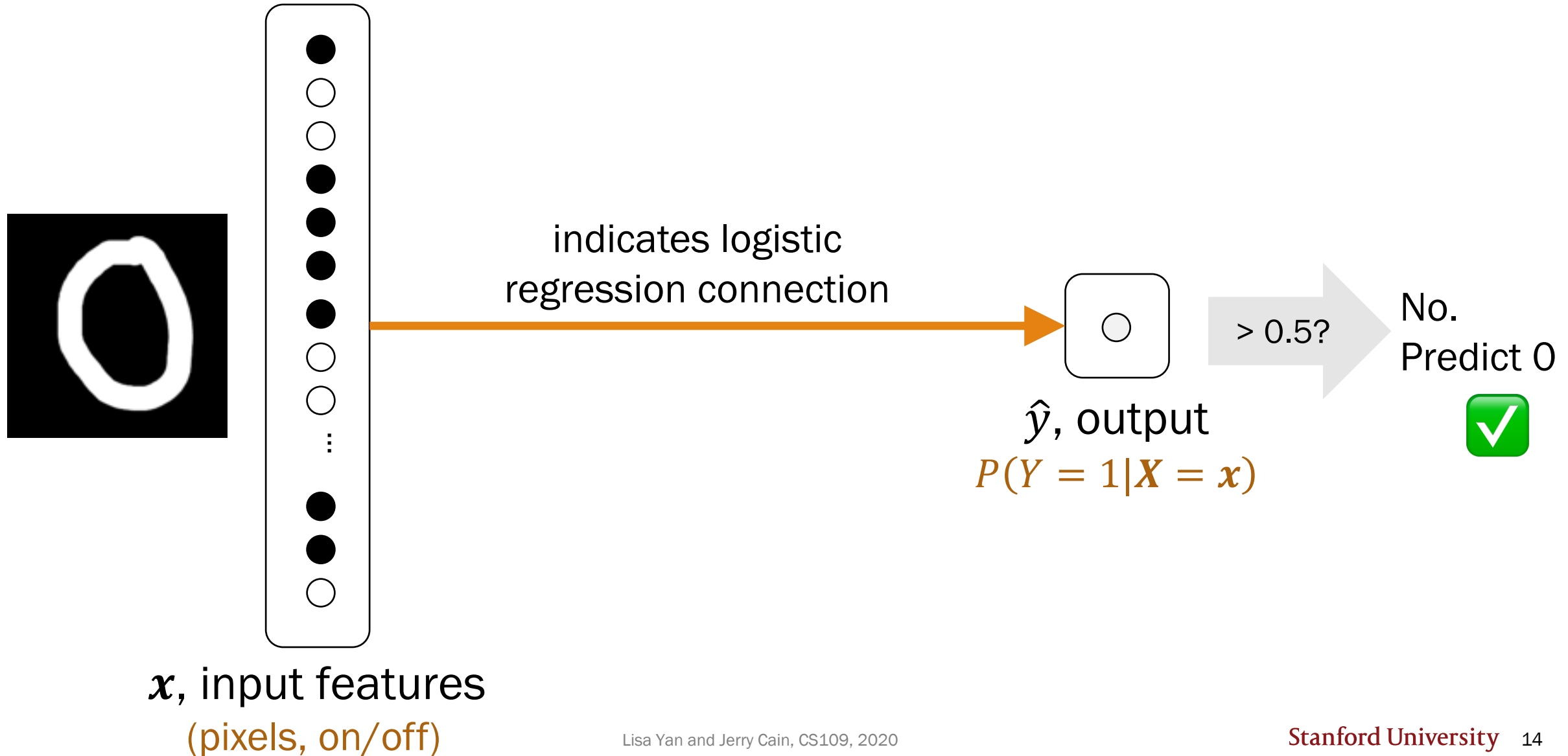
$$y^{(i)} = 1$$

We make feature vectors from (digitized) pictures of numbers.

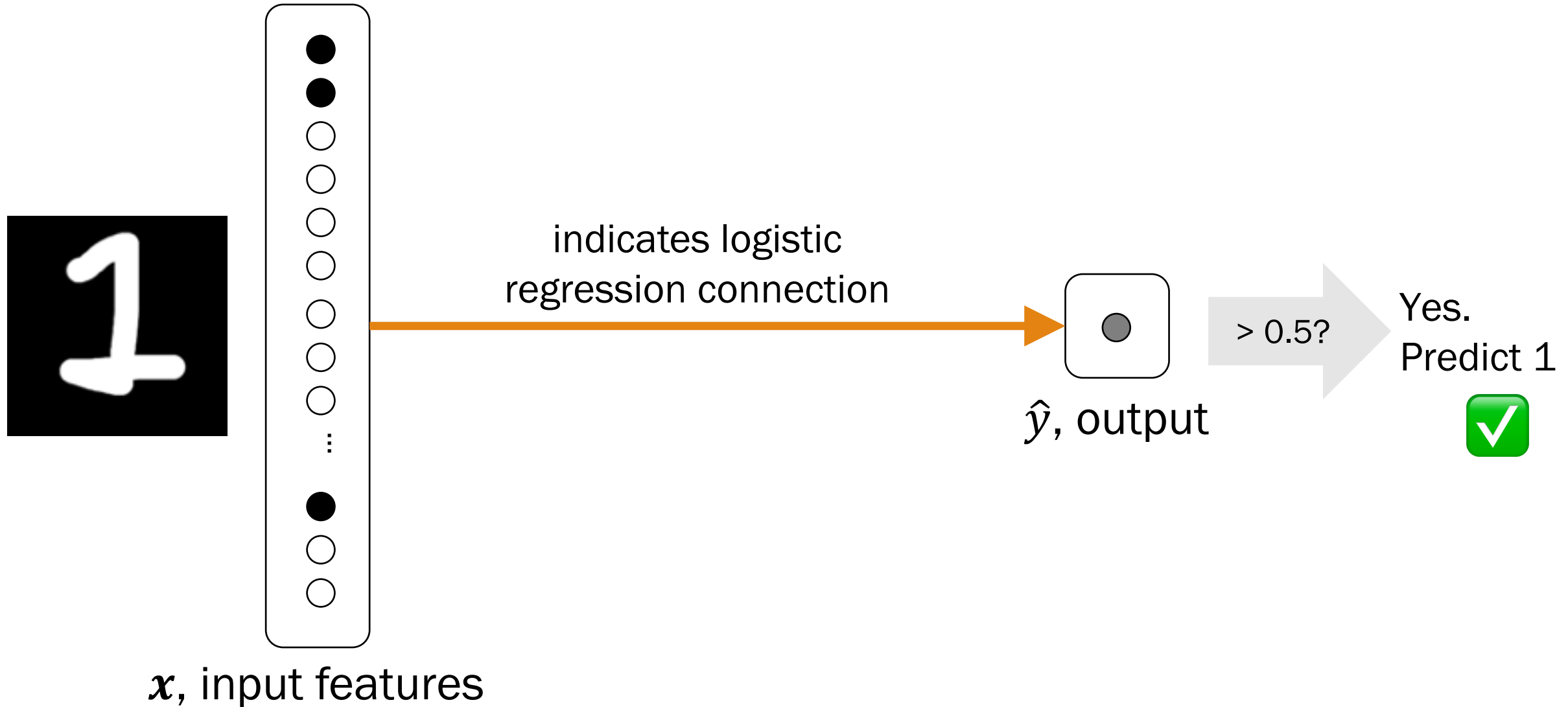
# Logistic Regression



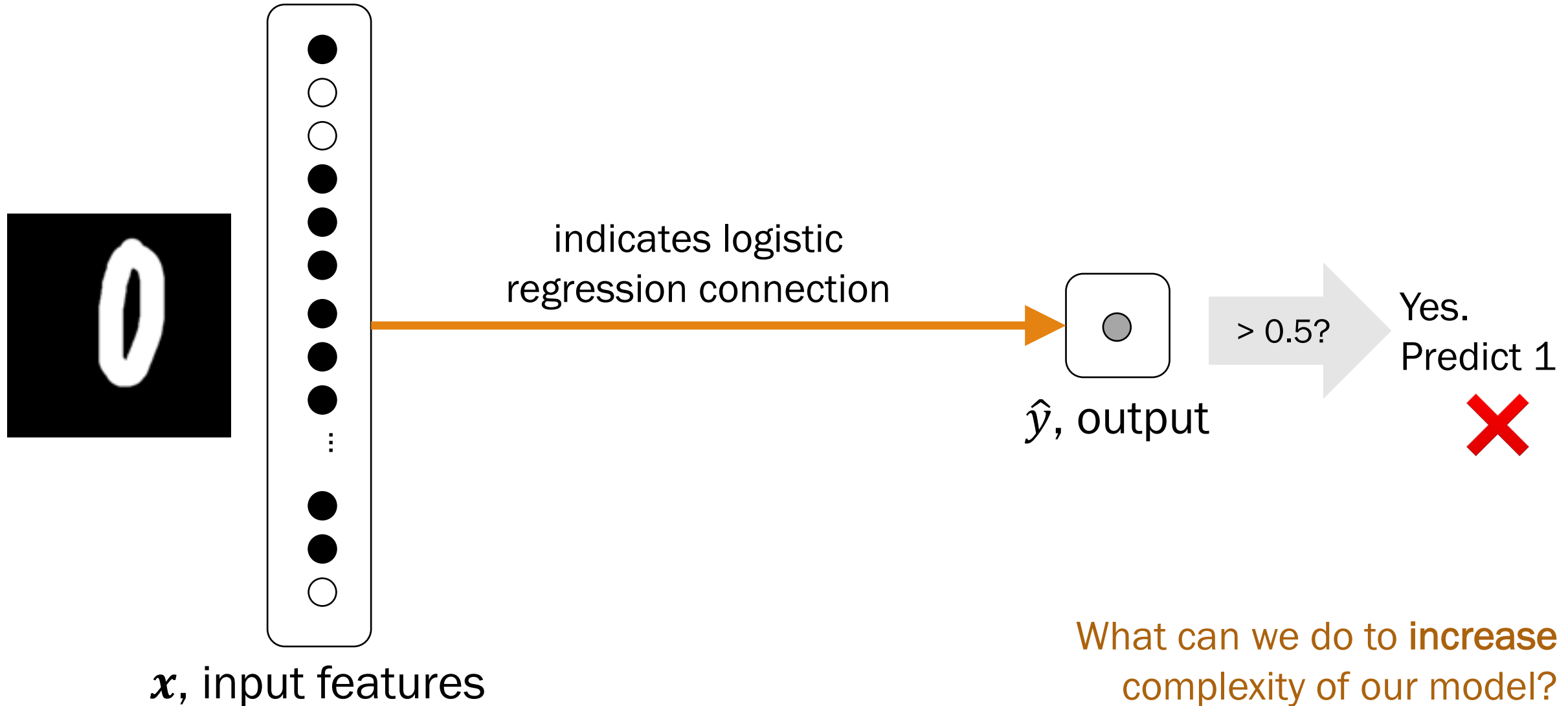
# Logistic Regression



# Logistic Regression

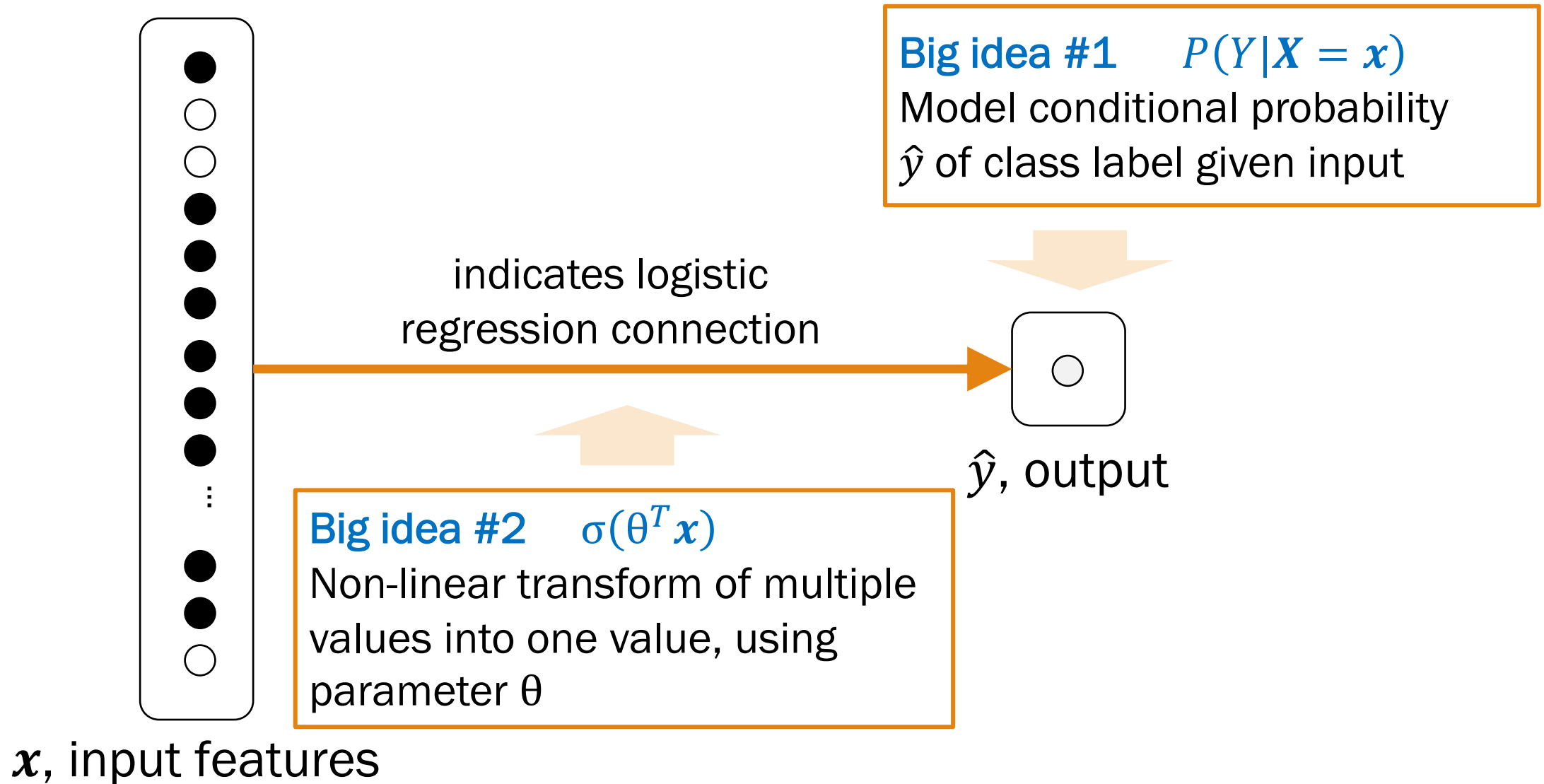


# Logistic Regression

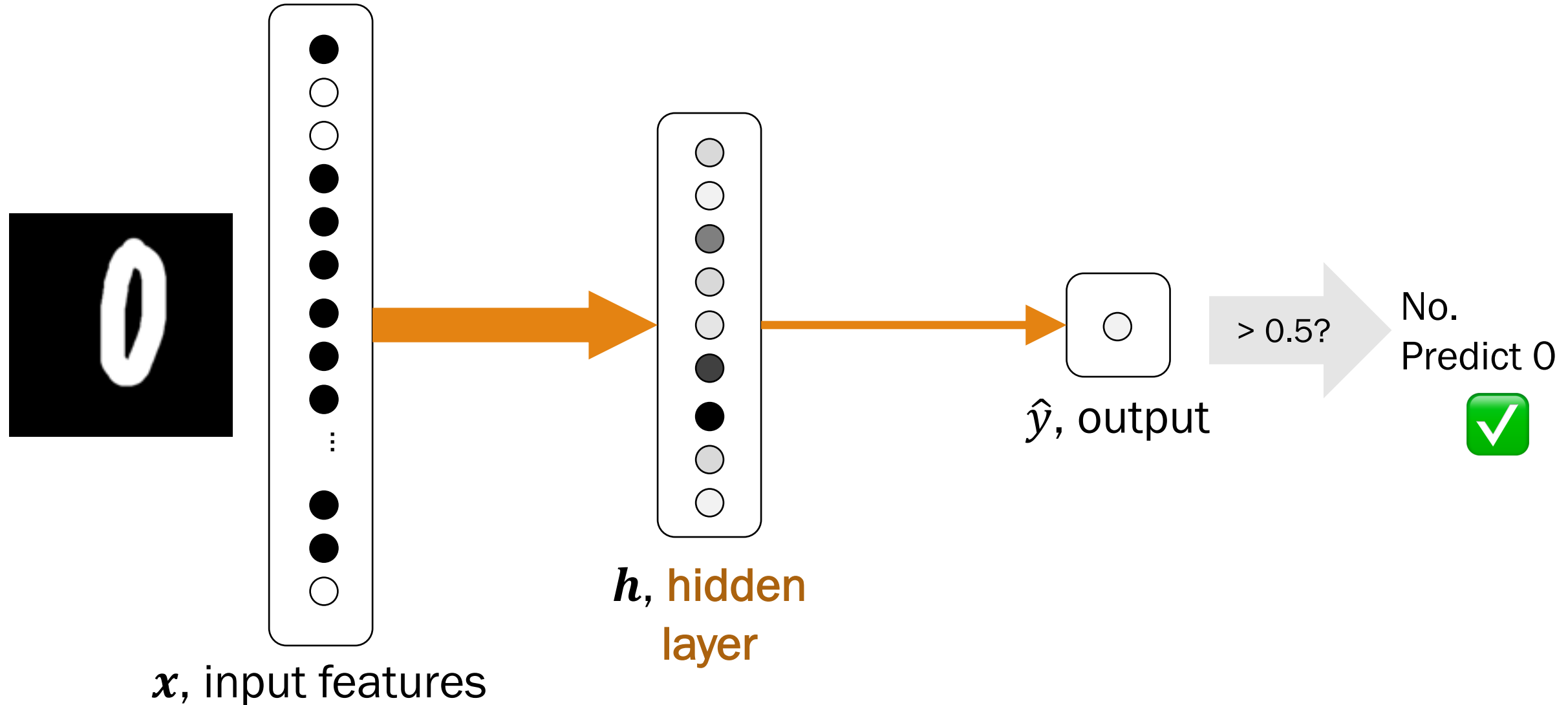




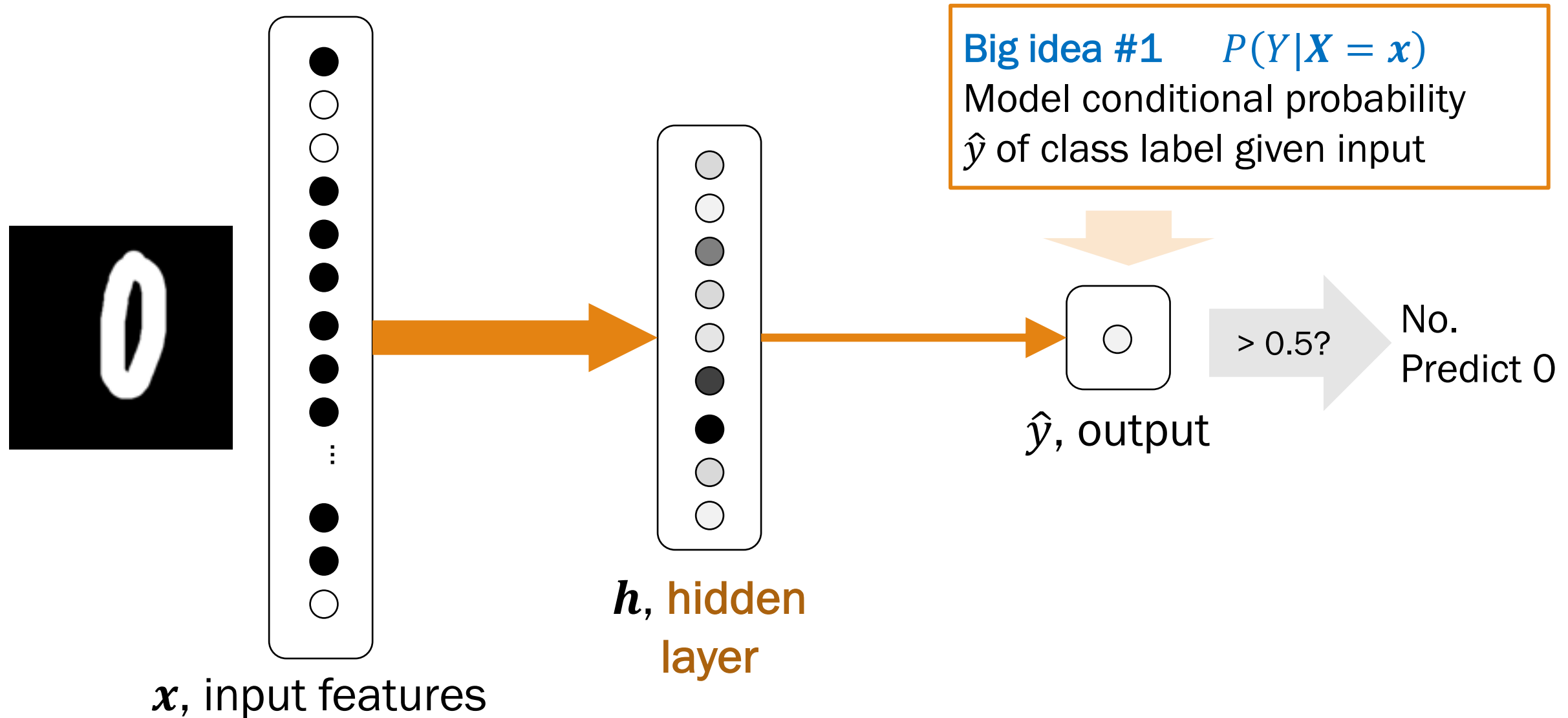
# Take two big ideas from Logistic Regression



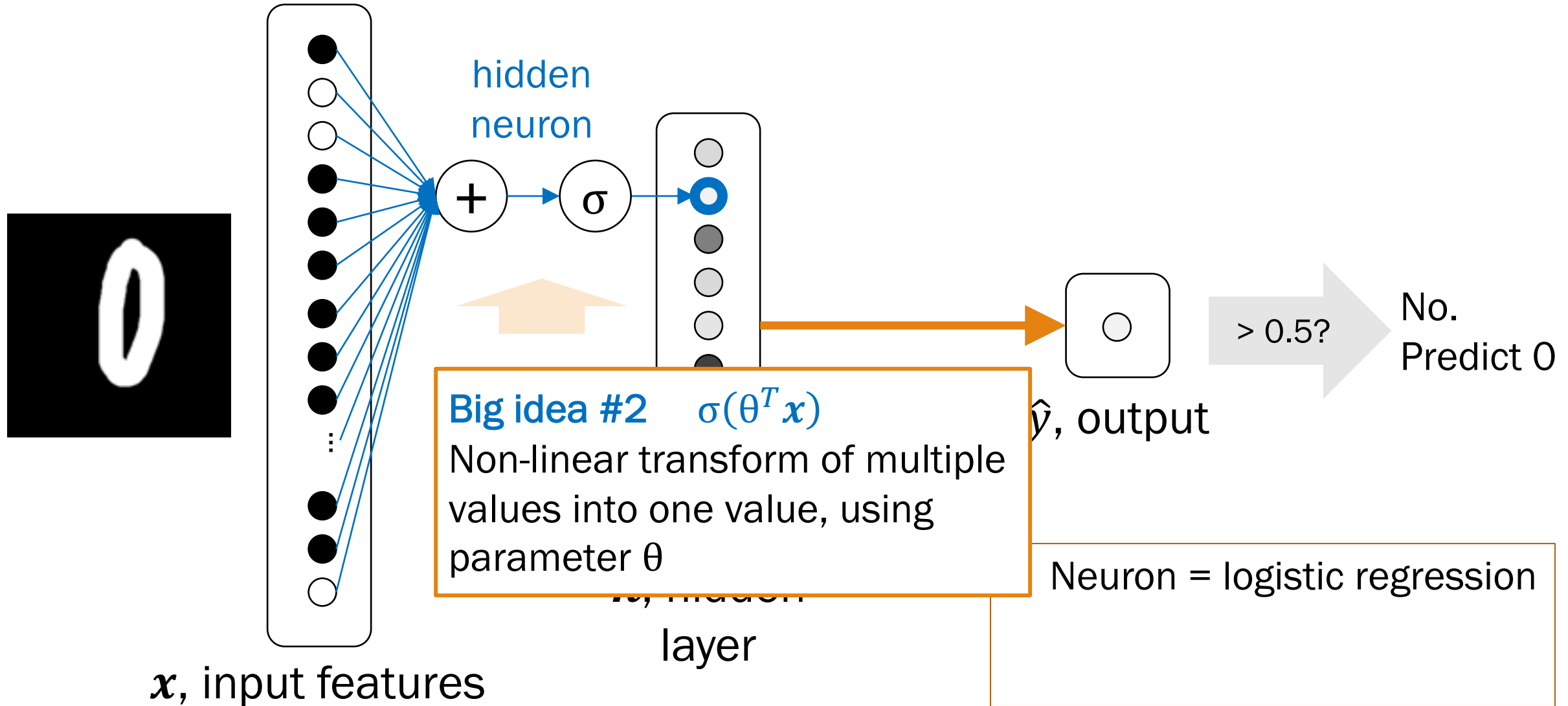
# Introducing: The Neural network



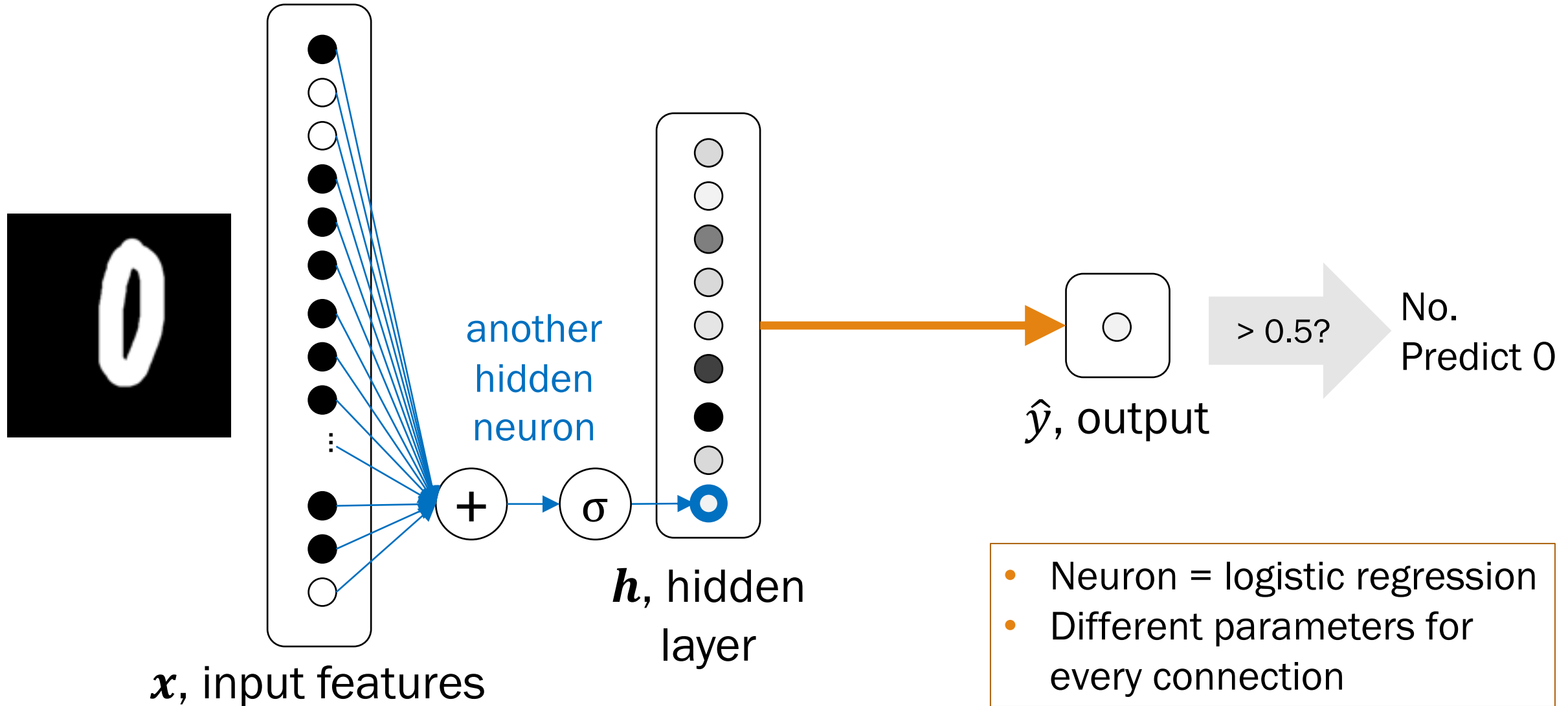
# Neural network



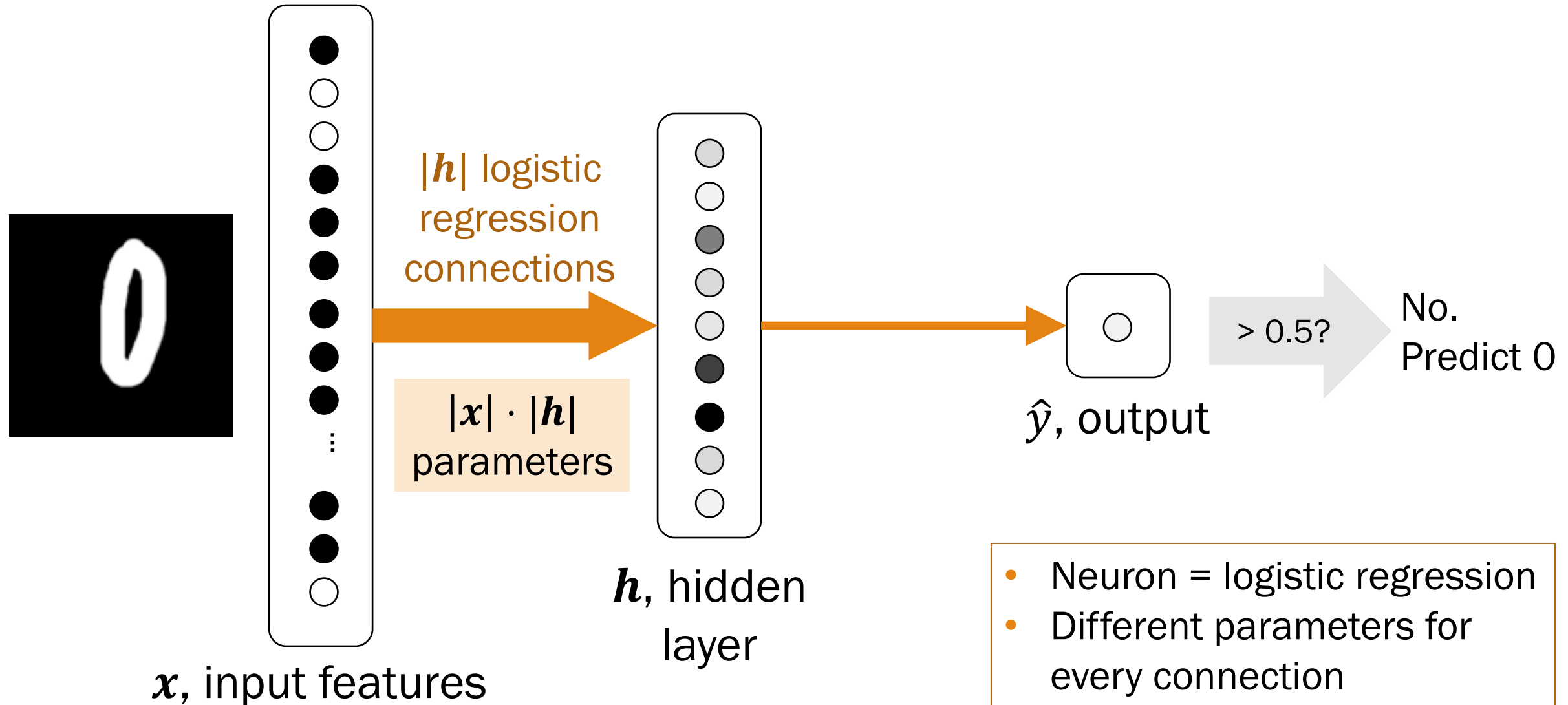
# Feed neurons into other neurons



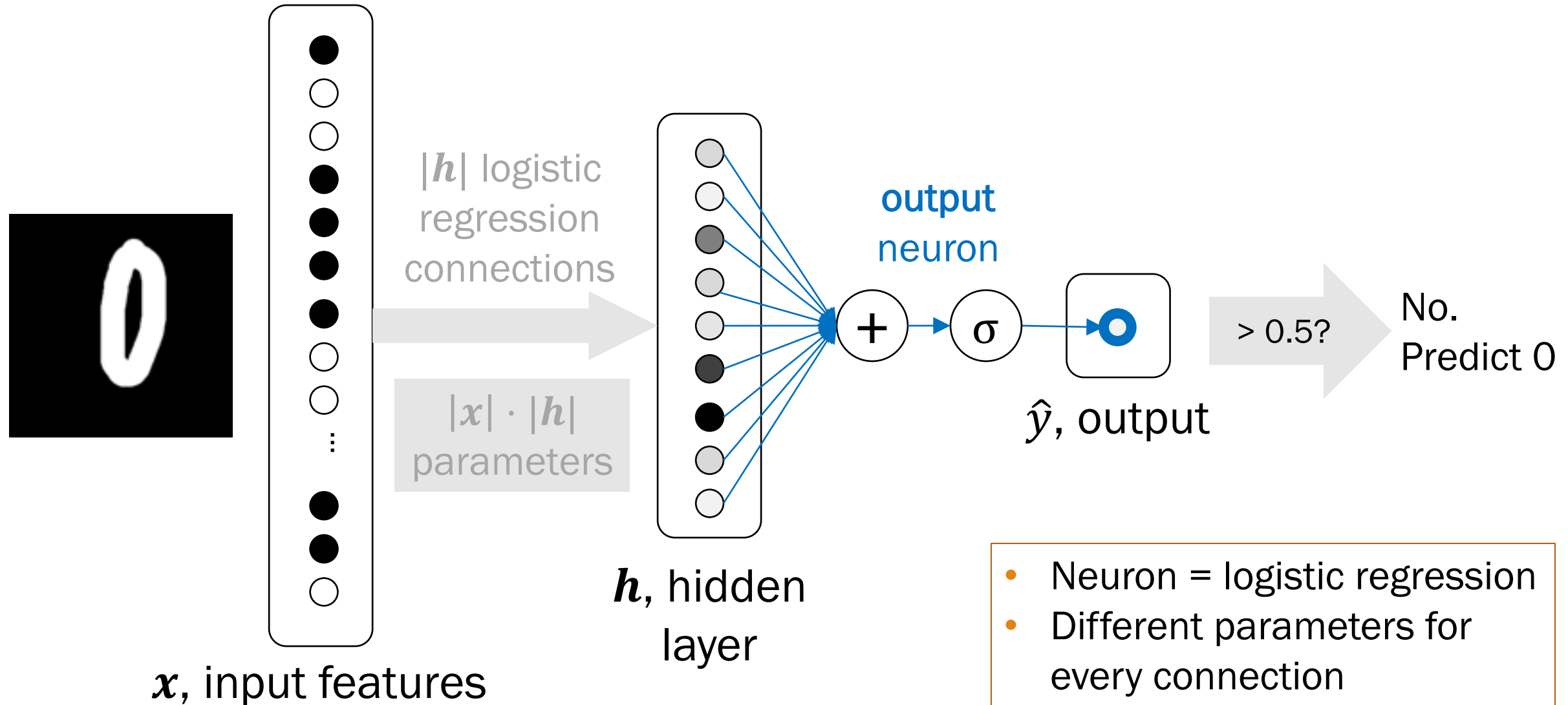
# Feed neurons into other neurons



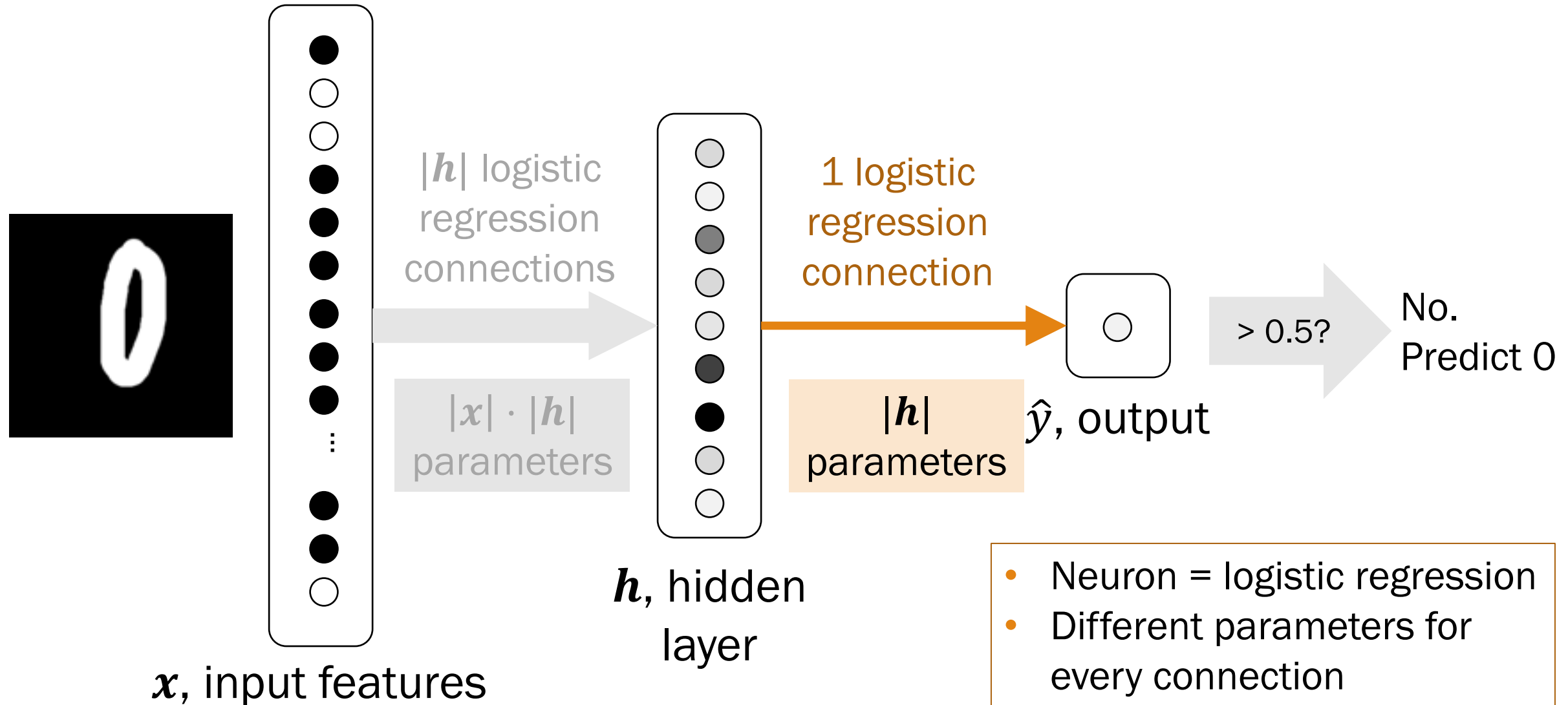
# Feed neurons into other neurons



# Feed neurons into other neurons



# Feed neurons into other neurons





# Think

Slide 26 asks you to think over by yourself.

Post any clarifications in chat!

Think by yourself: 2 min



(by yourself)

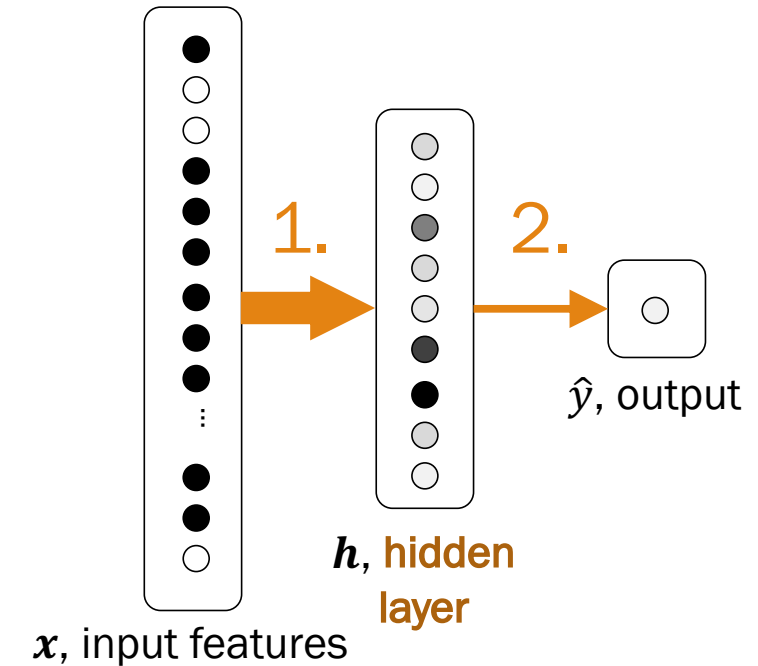
# Why doesn't a linear model introduce "complexity"?

## Neural network:

1. for  $j = 1, \dots, |\mathbf{h}|$ :

$$h_j = \sigma \left( \theta_j^{(h)T} \mathbf{x} \right)$$

2.  $\hat{y} = \sigma \left( \theta^{(\hat{y})T} \mathbf{h} \right) = P(Y = 1 | \mathbf{X} = \mathbf{x})$



## Linear network:

1. for  $j = 1, \dots, |\mathbf{h}|$ :

$$h_j = \theta_j^{(h)T} \mathbf{x}$$

2.  $\hat{y} = \sigma \left( \theta^{(\hat{y})T} \mathbf{h} \right) = P(Y = 1 | \mathbf{X} = \mathbf{x})$



# Why doesn't a linear model introduce "complexity"?

## Neural network:

1. for  $j = 1, \dots, |\mathbf{h}|$ :

$$h_j = \sigma \left( \theta_j^{(h)T} \mathbf{x} \right)$$

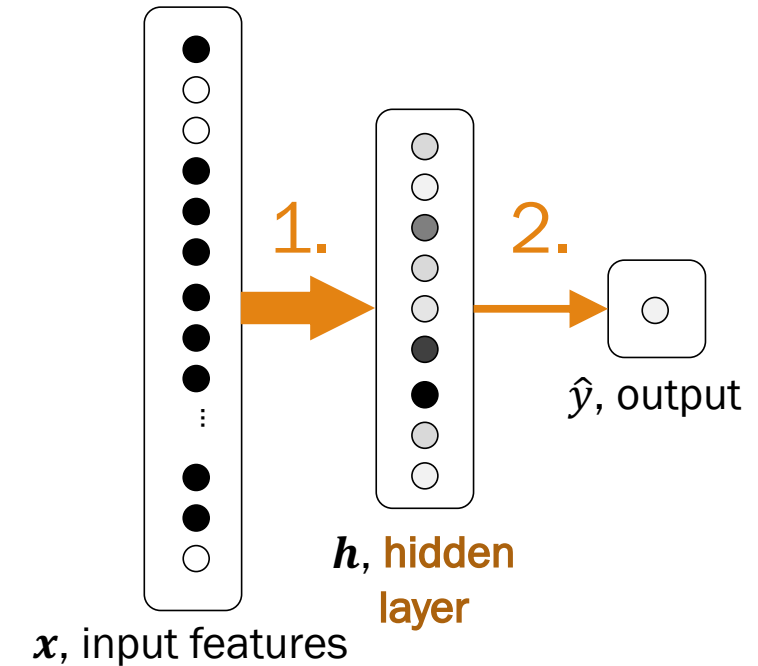
2.  $\hat{y} = \sigma \left( \theta^{(\hat{y})T} \mathbf{h} \right) = P(Y = 1 | \mathbf{X} = \mathbf{x})$

## Linear network:

1. for  $j = 1, \dots, |\mathbf{h}|$ :

$$h_j = \theta_j^{(h)T} \mathbf{x}$$

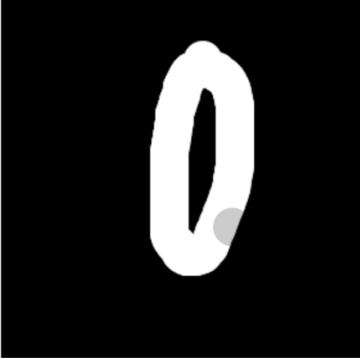
2.  $\hat{y} = \sigma \left( \theta^{(\hat{y})T} \mathbf{h} \right) = P(Y = 1 | \mathbf{X} = \mathbf{x})$



The linear model is effectively a single logistic regression with  $|\mathbf{x}|$  parameters.

# Demonstration

Draw your number here



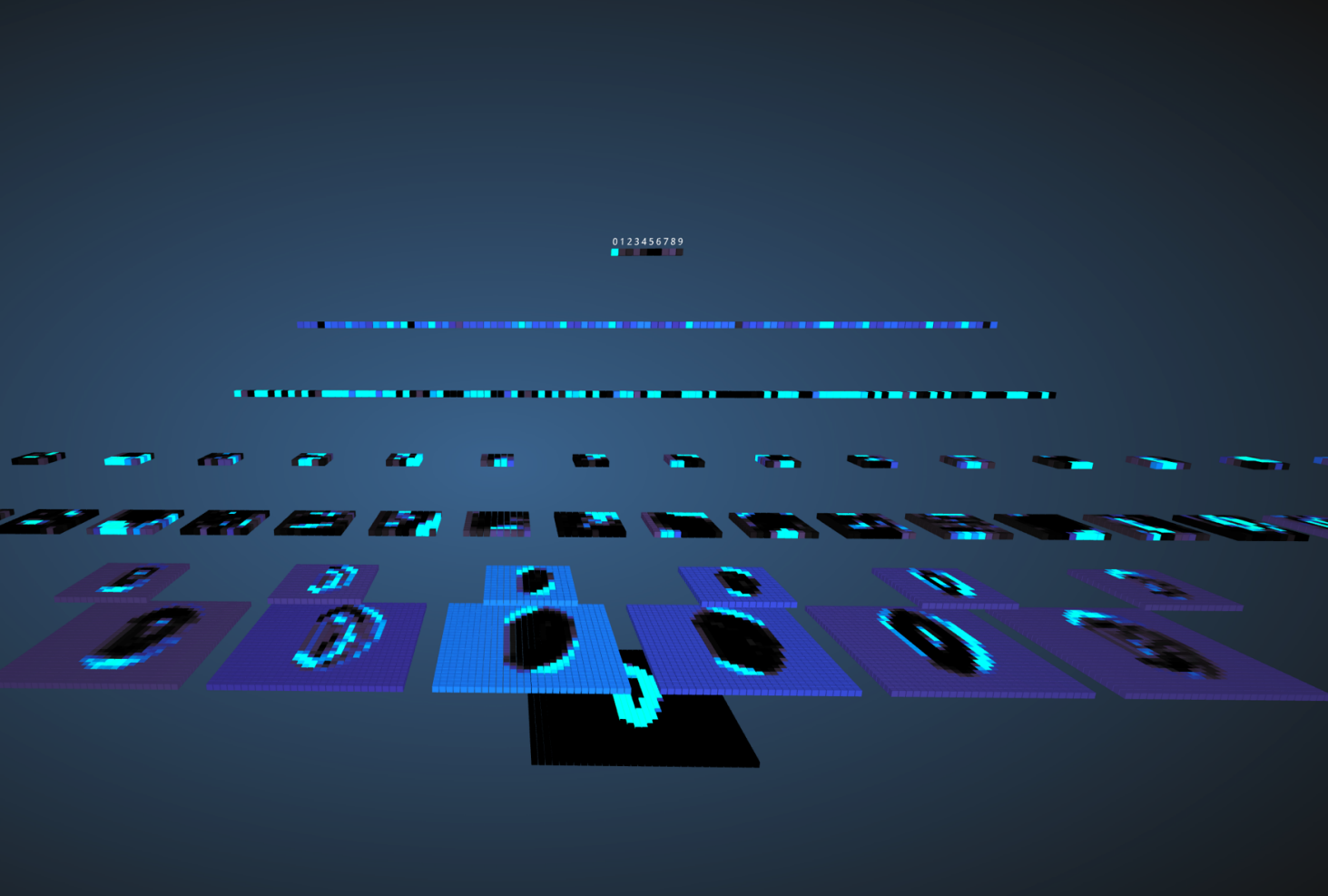
0 1 2 3 4 5 6 7 8 9

X [Pencil] [Eraser]

Downsampled drawing: 0  
First guess: 0  
Second guess: 8

Layer visibility

Input layer	Show
Convolution layer 1	Show
Downsampling layer 1	Show
Convolution layer 2	Show
Downsampling layer 2	Show



<http://scs.ryerson.ca/~aharley/vis/conv/>

# Interlude for jokes

# Probability as college students

## The Six Probability Distributions You'll Meet in Your Sorority

*Gaussian*

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

### The One Who Does It All

You see her everywhere. Physics, math, computer science. How is she in all your classes? And she does amazing in all of them, keeping well ahead of the curve. You'd like to be friends, but despite her popularity, she seems to have been regressing towards mean spirited behavior. At least she seems normal.

*Binomial*

$$p(k) = \binom{n}{k} p^k q^{n-k}$$

### The Confidant

You think she's related with The One Who Does It All. But how did she turn out so sweet? You can tell her anything, and you know she's discreet enough to keep it under wraps. But take care of her. This one will bet all she's got on a handful of coin tosses.

### The Scatterbrain

This girl cannot remember anything. She needs to ask your name every time you meet her. You're pretty sure you were friends during rush, but things have dropped off quickly since then.

*Exponential*

$$f(x) = \lambda e^{-\lambda x}$$

### The Background Boyfriend

He started dating The Confidant last semester, but you can't see what they have in common. He's not obnoxious, but he's not particularly charming either. No matter what the event, he always gives the same response: "sure, k". His flat personality might be mistaken for a chill, laid-back attitude.

*Uniform*

$$f(k) = \frac{1}{b-a}$$

*Poisson*

$$f(k) = e^{-\lambda} \frac{\lambda^k}{k!}$$

### The Ghost

This sister always seems kind of distracted and never shows up to anything. In fact, the last time you saw her was two months ago - counting raindrops outside the science building.

### The Ride or Die

You always know where she's going to be. Your relationship can get convoluted, but she's always got your back when things reset to square 0. Your rock-solid support, you can count on her to never vary in her ~~mountain range~~.

*Delta*

$$\delta(x) = \begin{cases} +\infty, & x = 0 \\ 0, & x \neq 0 \end{cases}$$



(A useful construct that connects discrete PMF to continuous PDF)

# Neural networks

---

A neural network (like logistic regression) gets intelligence from its parameters  $\theta$ .

## Training

- Learn parameters  $\theta$
- Find  $\theta_{MLE}$  that maximizes likelihood of training data (MLE)

## Testing/ Prediction

For input feature vector  $\mathbf{X} = \mathbf{x}$ :

- Use parameters to compute  $\hat{y} = P(Y = 1 | \mathbf{X} = \mathbf{x})$
- Classify instance as: 
$$\begin{cases} 1 & \hat{y} > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

# Neural networks

---

A neural network (like logistic regression) gets intelligence from its parameters  $\theta$ .

## Training

- Learn parameters  $\theta$
- Find  $\theta_{MLE}$  that maximizes likelihood of training data (MLE)

How do we learn the  $|\mathbf{x}| \cdot |\mathbf{h}| + |\mathbf{h}|$  parameters?

Gradient ascent + chain rule!



## 1. Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

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

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

## 2. Compute gradient

Find  $|\mathbf{x}|$  parameters

## 3. Optimize

```
initialize params
repeat many times:
  compute gradient
  params +=  $\eta$  * gradient
```

# Training: Neural networks

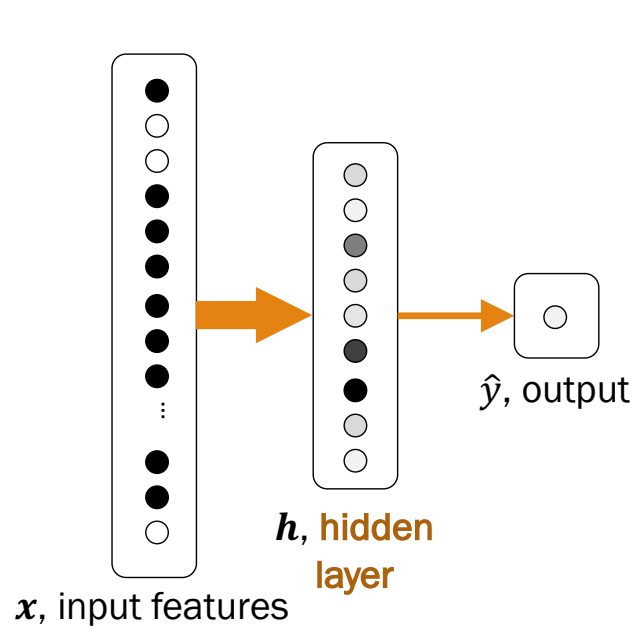
1. Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

2. Compute gradient

3. Optimize

# 1. Same output $\hat{y}$ , same log conditional likelihood



for  $j = 1, \dots, |\mathbf{h}|$ :

$$h_j = \sigma\left(\theta_j^{(h)T} \mathbf{x}\right)$$

$$\hat{y} = \sigma\left(\theta^{(\hat{y})T} \mathbf{h}\right) = P(Y = 1 | \mathbf{X} = \mathbf{x})$$

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

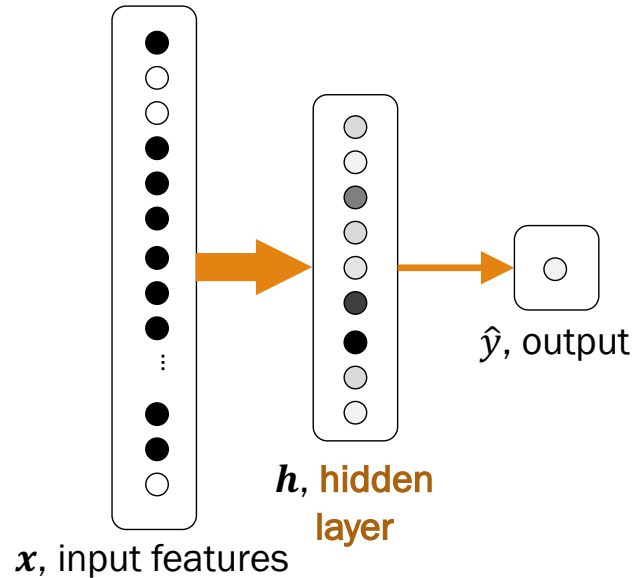
$$L(\theta) = \prod_{i=1}^n P(Y = y^{(i)} | \mathbf{X} = \mathbf{x}^{(i)}, \theta)$$

Binary class labels:  
 $Y \in \{0, 1\}$

$$= \prod_{i=1}^n (\hat{y}^{(i)})^{y^{(i)}} (1 - \hat{y}^{(i)})^{1-y^{(i)}}$$

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

# (model is a little more complicated)



$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$
$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

for  $j = 1, \dots, |\mathbf{h}|$ :

$$h_j = \sigma \left( \theta_j^{(h)T} \mathbf{x} \right) \quad \text{dimension } |\mathbf{x}|$$

$$\hat{y} = \sigma \left( \theta^{(\hat{y})T} \mathbf{h} \right) = P(Y = 1 | \mathbf{X} = \mathbf{x}) \quad \text{dimension } |\mathbf{h}|$$

To optimize for  
log conditional likelihood,  
we now need to find:

$$|\mathbf{h}| \cdot |\mathbf{x}| + |\mathbf{x}| \quad \text{parameters}$$

## 2. Compute gradient

1. Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

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

$$h_j = \sigma \left( \theta_j^{(h)T} \mathbf{x} \right) \quad \text{for } j = 1, \dots, |\mathbf{h}| \quad \hat{y} = \sigma \left( \theta^{(\hat{y})T} \mathbf{h} \right)$$

2. Compute gradient

Take gradient with respect to all  $\theta$  parameters

3. Optimize

**Calculus refresher #1:**

Derivative(sum) =  
sum(derivative)

**Calculus refresher #2:**

Chain rule 

# 3. Optimize

1. Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

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

$$h_j = \sigma \left( \theta_j^{(h)T} \mathbf{x} \right) \quad \text{for } j = 1, \dots, |\mathbf{h}| \quad \hat{y} = \sigma \left( \theta^{(\hat{y})T} \mathbf{h} \right)$$

2. Compute gradient

Take gradient with respect to all  $\theta$  parameters

3. Optimize

```
initialize params
repeat many times:
  compute gradient
  params +=  $\eta$  * gradient
```

# Training a neural net

1. Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

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

Wait, did we just skip something difficult?

2. Compute

3. Optimize

```
initialize params
repeat many times:
  compute gradient
  params += η * gradient
```

## 2. Compute gradient via backpropagation

1. Optimization problem:

$$\theta_{MLE} = \arg \max_{\theta} \prod_{i=1}^n f(y^{(i)} | \mathbf{x}^{(i)}, \theta) = \arg \max_{\theta} LL(\theta)$$

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

$$h_j = \sigma \left( \theta_j^{(h)T} \mathbf{x} \right) \quad \text{for } j = 1, \dots, |\mathbf{h}| \quad \hat{y} = \sigma \left( \theta^{(\hat{y})T} \mathbf{h} \right)$$

2. Compute gradient

Take gradient with respect to all  $\theta$  parameters

3. Optimize

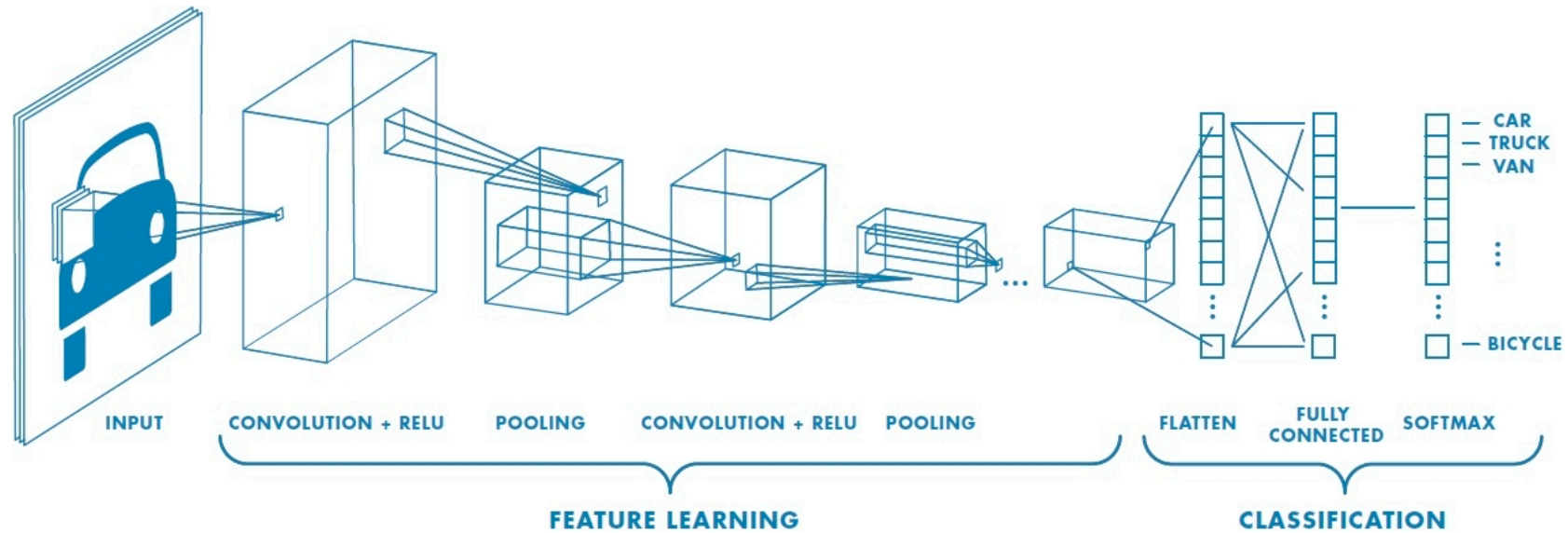
initial  
repeat  
compute  
parameters

Learn the tricks behind **backpropagation** in CS229, CS231N, CS224N, etc.



# Beyond the basics

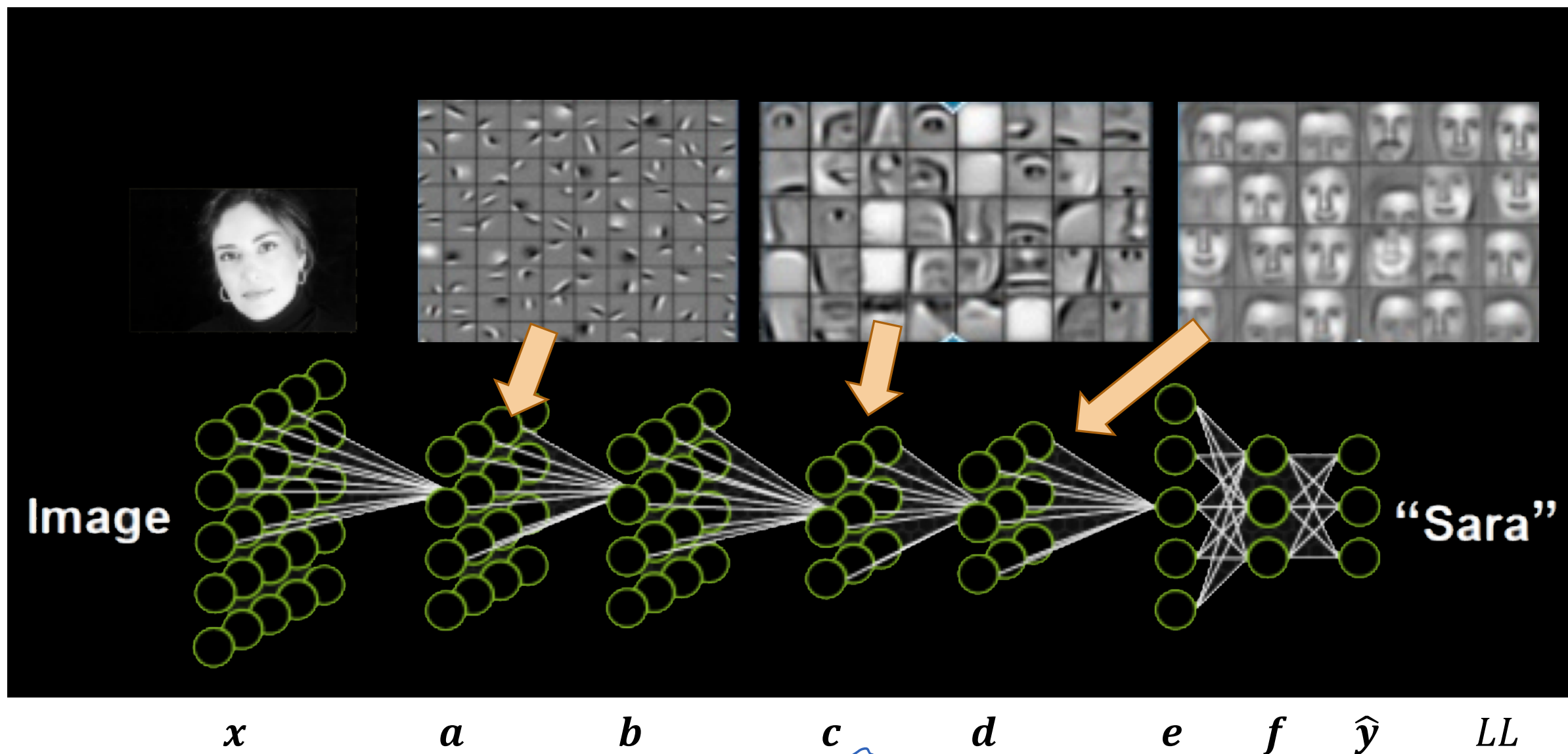
# Shared weights?



It turns out if you want to force some of your weights to be shared over different neurons, the math isn't much harder.

**Convolution** is an example of such weight-sharing and is used a lot for vision (Convolutional Neural Networks, CNN).

# Neural networks with multiple layers



$x$

$a$

$b$

$c$

$d$

$e$

$f$

$\hat{y}$

$LL$

# Neurons learn features of the dataset



Neurons in later layers will respond strongly to high-level features of your **training data**.

If your training data is faces, you will get lots of face neurons.

If your training data is all of YouTube...



...you get a cat neuron.



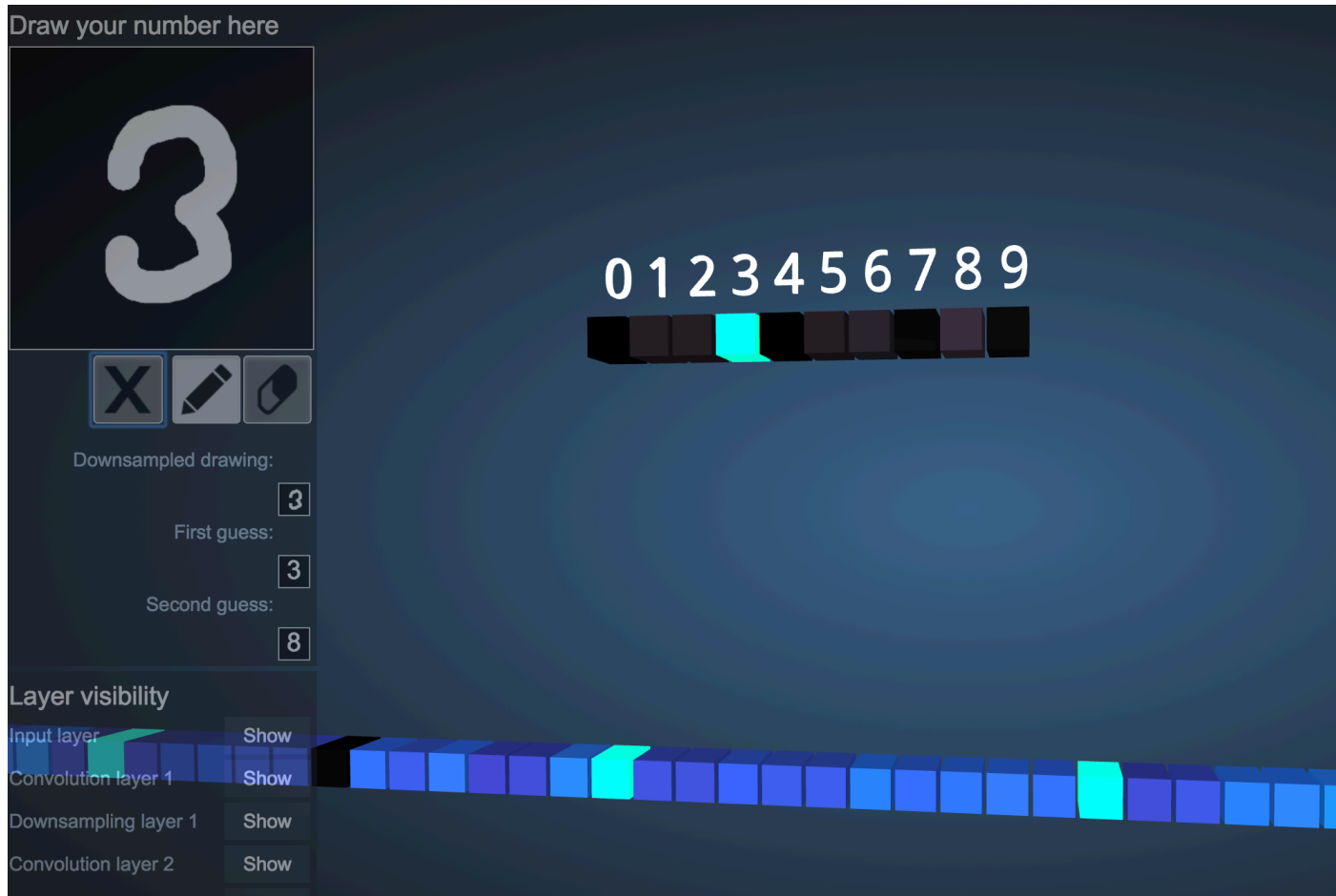
Top stimuli in test set



Optimal stimulus found by numerical optimization



# Multiple outputs?



**Softmax** is a generalization of the sigmoid function.

sigmoid( $z$ ): value in range  $[0, 1]$

$z \in \mathbb{R}$ :

$$P(Y = 1 | X = \mathbf{x}) = \sigma(z)$$

(equivalent: **Bernoulli**  $p$ )



softmax( $z$ ):  $k$ -dimensional values in range  $[0, 1]$  that add up to 1

$\mathbf{z} \in \mathbb{R}^k$ :

$$P(Y = i | X = \mathbf{x}) = \text{softmax}(\mathbf{z})_i$$

(equivalent: **Multinomial**  $p_1, \dots, p_k$ )

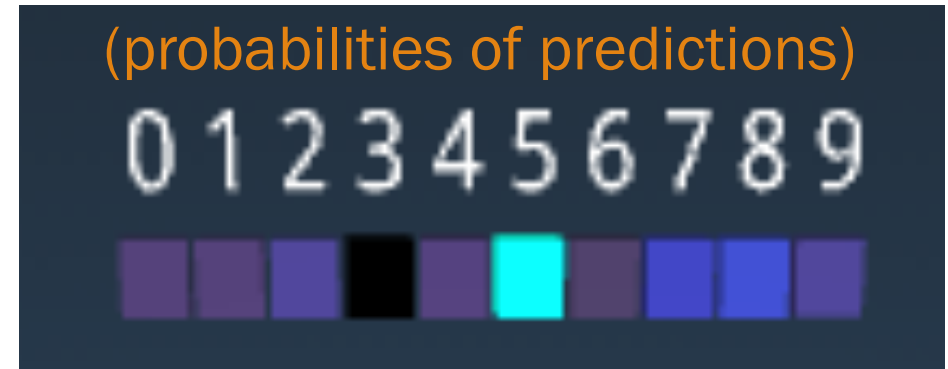
# Softmax test metric: Top-5 error

$Y = y$	$P(Y = y X = \mathbf{x})$
5	0.14
8	0.13
7	0.12
2	0.10
9	0.10
4	0.09
1	0.09
0	0.09
6	0.08
3	0.05



## Top-5 classification error

What % of datapoints did *not* have the correct class label in the top-5 predictions?



# ImageNet classification

22,000 categories

14,000,000 images

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

...  
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





# ImageNet classification challenge

~~22,000 categories~~

1000 categories

smoothhound shark, *Mustelus mustelus*  
dogfish, *Mustelus canis*

Florida smoothhound, *Mustelus norrisi*

14,000,000 images

1,200,000 images in train set

codon obseus

200,000 images in test set

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

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  
rougthead 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*

...

# ImageNet challenge: Top-5 classification error

(lower is better)

# 99.5%

Random guess

$$P(\text{true class label not in 5 guesses}) = \frac{\binom{999}{5}}{\binom{1000}{5}} = \frac{995}{1000}$$

# ImageNet challenge: Top-5 classification error

(lower is better)

99.5%

Random guess

25.8%

Pre-Neural Networks

5.1%

Humans  
(2014)

16.4%

GoogLe Net  
(2015)

Russakovsky et al., ImageNet Large Scale Visual Recognition Challenge. IJCV 2015

Szegedy et al., Going Deeper With Convolutions. CVPR 2015

Hu et al., Squeeze-and-Excitation Networks. Preprint arXiv 2017

# ImageNet challenge: Top-5 classification error

(lower is better)

99.5%

Random guess

25.8%

Pre-Neural Networks

5.1%

Humans  
(2014)

16.4%

GoogLe Net  
(2015)

2.25%

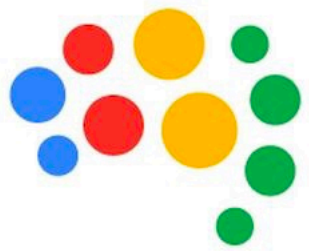
SENet  
(2017)

Russakovsky et al., ImageNet Large Scale Visual Recognition Challenge. IJCV 2015

Szegedy et al., Going Deeper With Convolutions. CVPR 2015

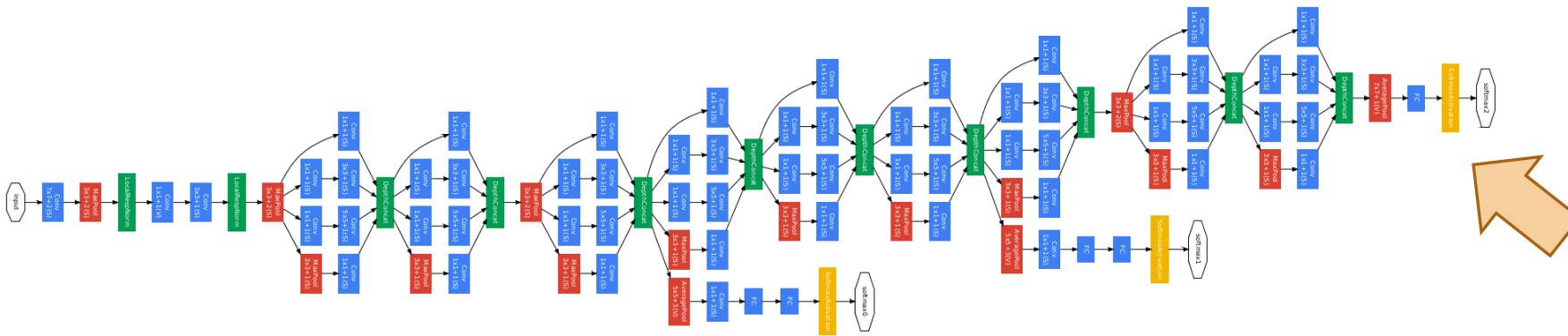
Hu et al., Squeeze-and-Excitation Networks. Preprint arXiv 2017

# GoogLeNet (2015)



## Google Brain

1 Trillion Artificial Neurons  
(btw human brains have 1 billion neurons)



Multiple,  
Multi class output

22 layers deep!

# Speeding up gradient descent

minimizes loss (a function of prediction error)

```
initialize  $\theta_j = 0$  for  $0 \leq j \leq m$   
repeat many times:
```

```
  gradient[j] = 0 for  $0 \leq j \leq m$ 
```

```
  for each training example  $(x, y)$ :
```

```
    for each  $0 \leq j \leq m$ :
```

```
      compute gradient
```

```
   $\theta_j -= \eta * \text{gradient}[j]$  for all  $0 \leq j \leq m$ 
```

1. What if we have 1,200,000 images in our training set?

2. How can we speed up the update?

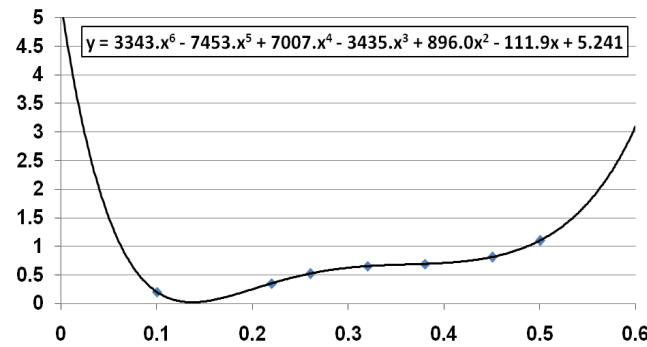
Our **batch gradient descent** (over the entire training set) will be slow + expensive.

1. Use **stochastic gradient descent** (randomly select training examples with replacement).
2. **Momentum update** (Incorporate “acceleration” or “deceleration” of gradient updates so far)

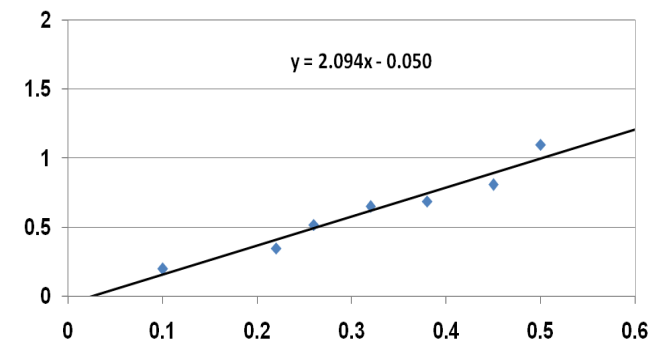
# Good ML = Generalization

## Overfitting

Fitting the training data too well, such that we lose generality of model for predicting new data



perfect fit, but bad predictor for new data

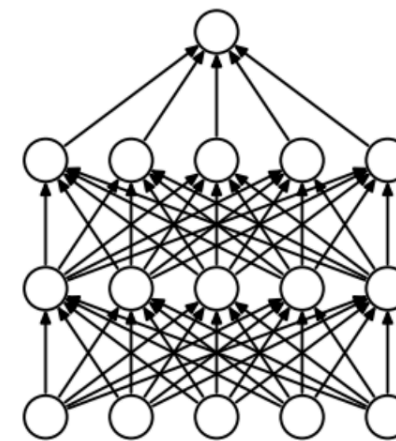


more general fit + better predictor for new data

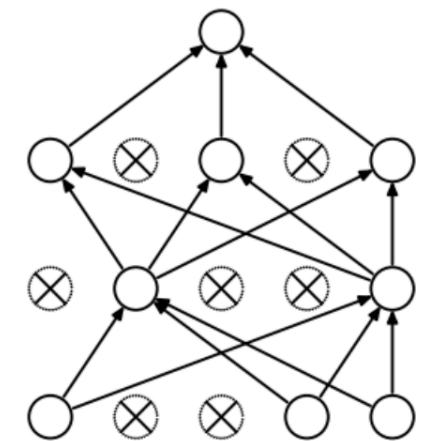
## Dropout

During training, randomly leave out some neurons each training step.

It will make your network more robust.



(a) Standard Neural Net



(b) After applying dropout.

# Making decisions?

---



Not everything is classification.

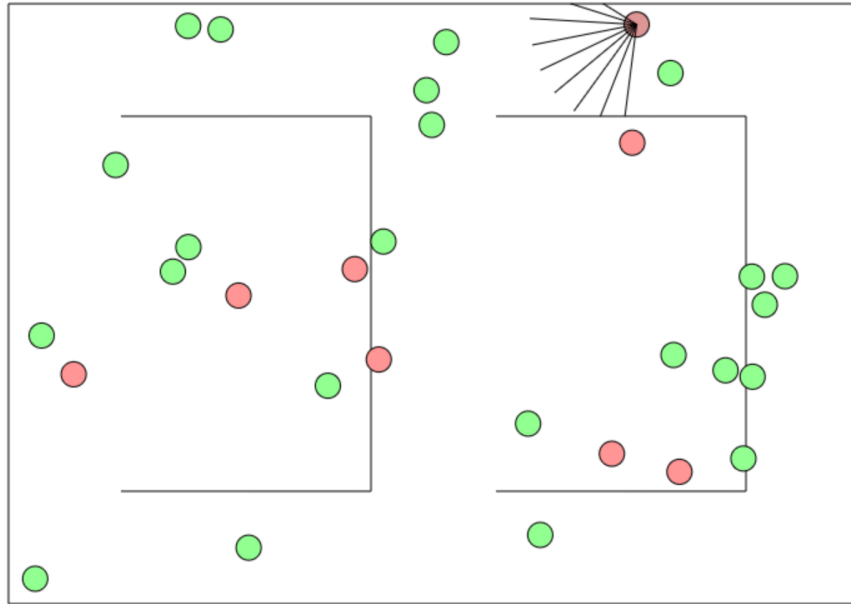
## Deep Reinforcement Learning

Instead of having the output of a model be a probability, you make output an expectation.

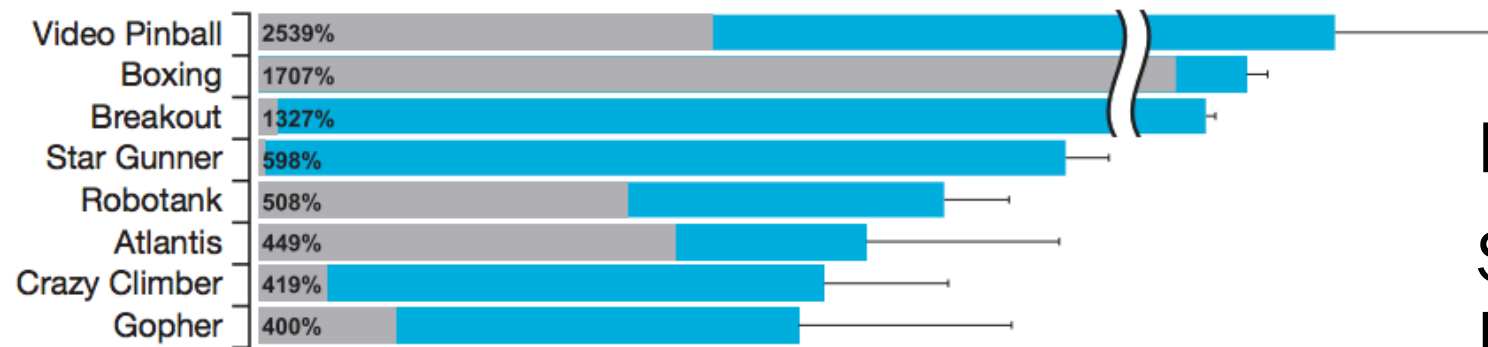
<http://cs.stanford.edu/people/karpathy/convnetjs/demo/rldemo.html>



# Deep Reinforcement Learning



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

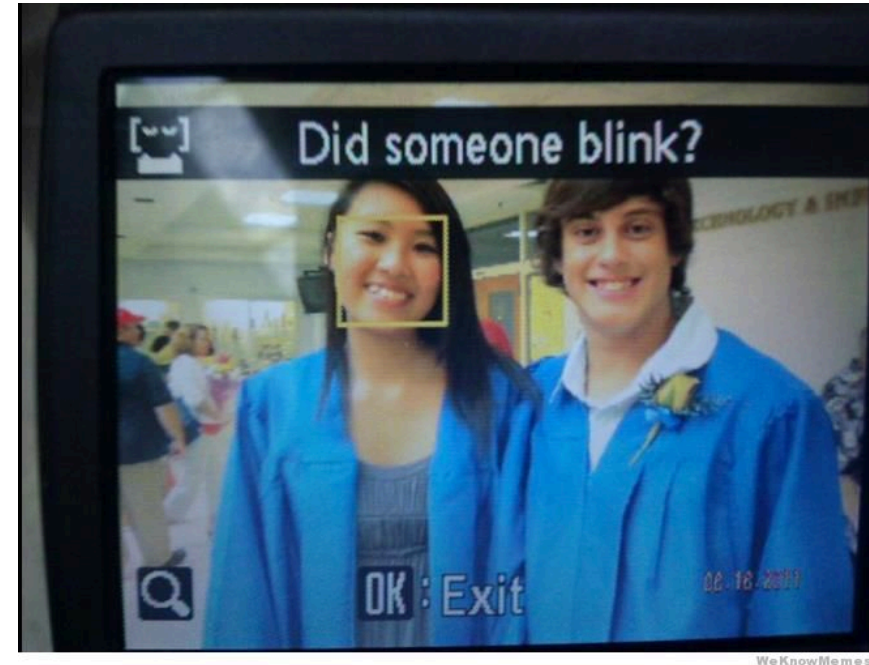


Deep Mind Atari Games  
Score compared to best  
human

What's missing?

How are you getting your data?

# Ethics and datasets



Sometimes machine learning feels universally unbiased.

We can even prove our estimators are “unbiased” (mathematically).

Google/Nikon/HP had biased datasets.

# Should your data be unbiased?

---

Dataset: Google News

$$\vec{\text{man}} - \vec{\text{woman}} \approx \vec{\text{king}} - \vec{\text{queen}}$$

$$\vec{\text{man}} - \vec{\text{woman}} \approx \vec{\text{computer programmer}} - \vec{\text{homemaker}}.$$

Should our unbiased data collection reflect society's systemic bias?

# How can we explain decisions?

---



If your task is **image classification**, reasoning about high-level features is relatively easy.

Everything can be visualized.

What if you are trying to classify social outcomes?

- Criminal recidivism
- Job performance
- Policing
- Terrorist risk
- At-risk kids

Ethics in Machine Learning  
is a whole new field. 😊