# CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2023

HAVERFORD
COLLEGE

# Admin

- Exam will be handed back on Tuesday

- Last candidate talk at **4:15pm TODAY**
  - Tea at 4pm
  - Student lunch Friday 12:30-1:30pm

# Outline for November 30

- Gaussian Mixture Models (GMMs)

- Kernel Density Estimation (KDE)

- Missing data

- Begin: neural networks

# Outline for November 30

- Gaussian Mixture Models (GMMs)


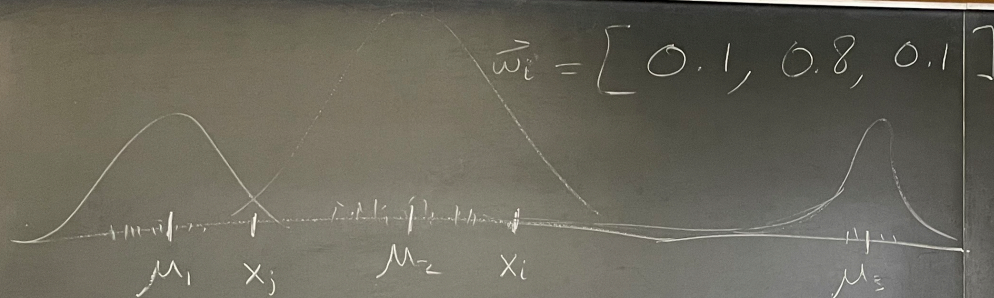- Kernel Density Estimation (KDE)


- Missing data


- Begin: neural networks

# Gaussian Mixture Models

- $\pi_k$ = prob of class $k$

  $= \dfrac{1}{K}$ for all $k$

- $\vec{\mu}_k$ = mean of cluster $k$

  $=$ random points

- $\sigma_k^2$ = variance of cluster $k$

  $=$ sample variance of pts closest to each mean.

initialization

---

# E-step   "soft" assignment

$w_{ik}$ = prob that $\vec{x}_i$ came from cluster $k$

$$= p(k|\vec{x}_i) = \frac{p(k)\, p(\vec{x}_i|k)}{p(\vec{x}_i)}$$

$$= \frac{\pi_k\, N(\vec{x}_i\,;\, \vec{\mu}_k, \sigma_k^2)}{\sum\limits_{j=1}^{K} \pi_j\, N(\vec{x}_i\,;\, \vec{\mu}_j, \sigma_j^2)}$$

$\vec{w}_i = [0.1, 0.8, 0.1]$



$M_1 \quad X_j \qquad M_2 \quad X_i \qquad\qquad M_3$

$\vec{w}_j = [0.49, 0.49, 0.02]$

$$W = \begin{bmatrix} M_1 & M_2 & M_3 \\ 0.1 & 0.8 & 0.1 \\ 0.49 & 0.49 & 0.02 \\ \vdots & & \end{bmatrix} \begin{matrix} \to 1 \\ \to 1 \end{matrix}$$

$n \times K$

---

$\boxed{M\text{-step}}$ param update

$M_k = \sum_{i=1}^{n} w_{ik} \quad \Big\}$ # pts assigned to cluster $k$

① $\pi_k = \dfrac{M_k}{n}$

② $\vec{M}_k = \dfrac{1}{M_k} \sum_{i=1}^{n} w_{ik} \vec{x}_i$

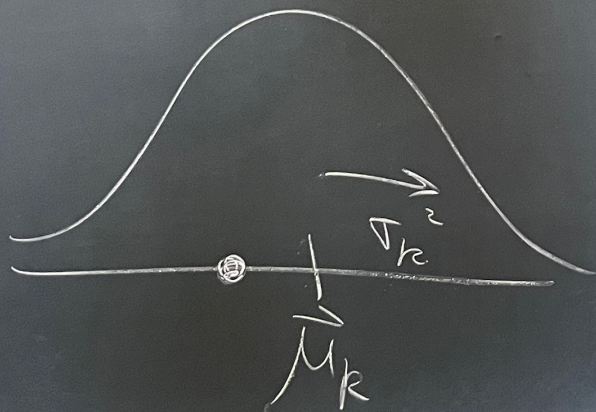③ $\sigma_k^2 =$ weighted sample variance.

K-means

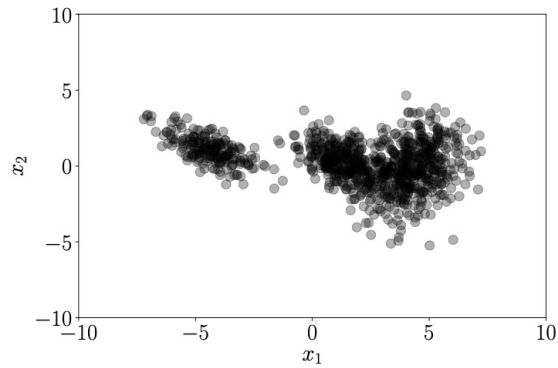$\boxed{\dfrac{1}{|c_k|} \sum_{x_i \in c_k} \vec{x}_i}$

# Generative process

① choose a cluster $k$
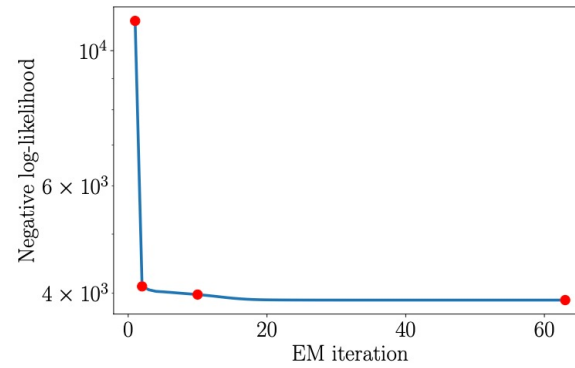   using $[\pi_1, \pi_2 \cdots \pi_k]$

② use $\bar{\mu}_k$ & $\sigma_k^2$ to sample
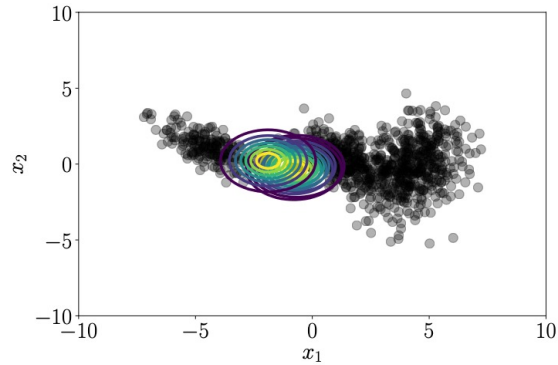   from $N(\bar{\mu}_k, \sigma_k^2)$

$$\sigma_k^2 \qquad \frac{x_i}{\sum_{j=1}^{N} x_j}$$

$\bar{\mu}_k$

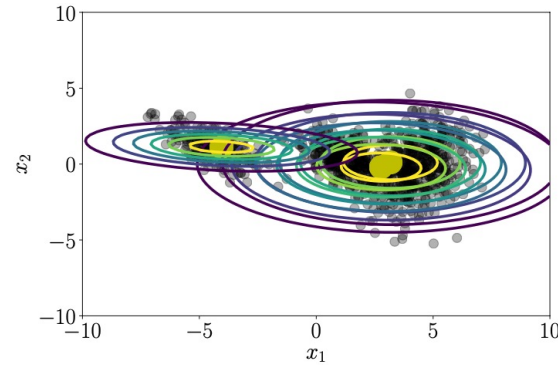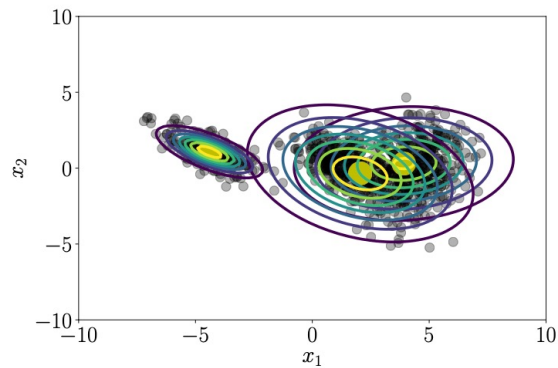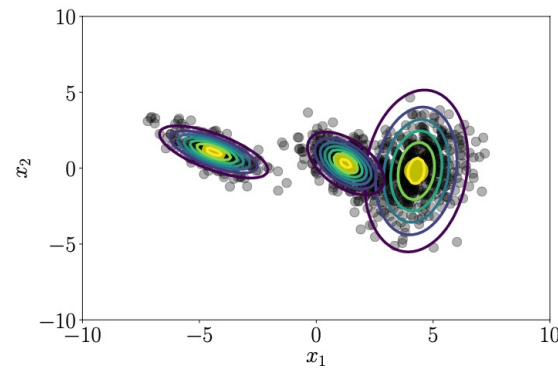*Density Estimation with Gaussian Mixture Models*



(a) Dataset.

(b) Negative log-likelihood.

(c) EM initialization.

(d) EM after one iteration.

(e) EM after 10 iterations.

(f) EM after 62 iterations.

Figure 11.7 from MML textbook
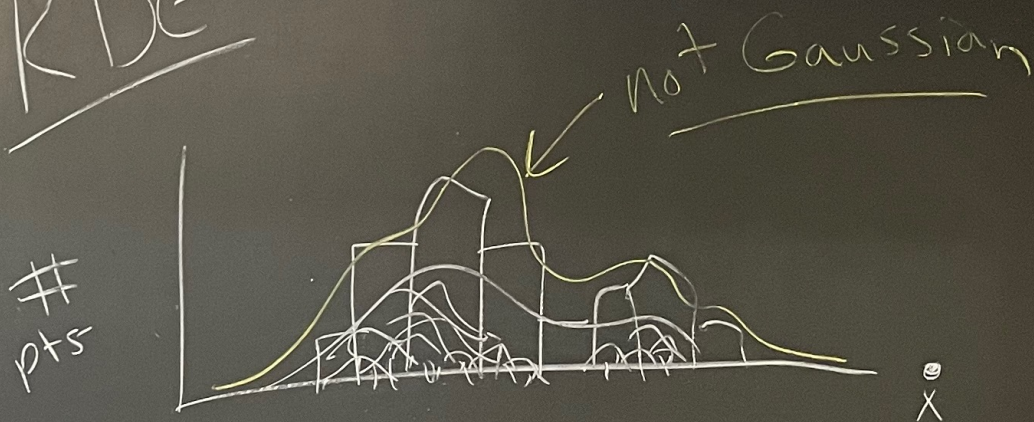
# Outline for November 30

- Gaussian Mixture Models (GMMs)

- Kernel Density Estimation (KDE)

- Missing data

- Begin: neural networks
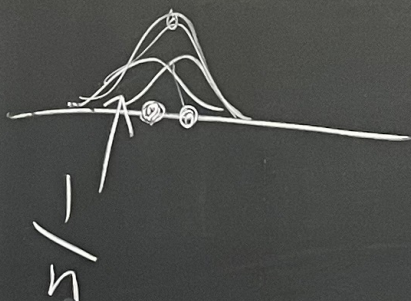
# KDE

not Gaussian



\# pts

$$p(x) = \frac{1}{nh} \sum_{i=1}^{n} k\left(\frac{x - x_i}{h}\right)$$

kernel

width

n/1
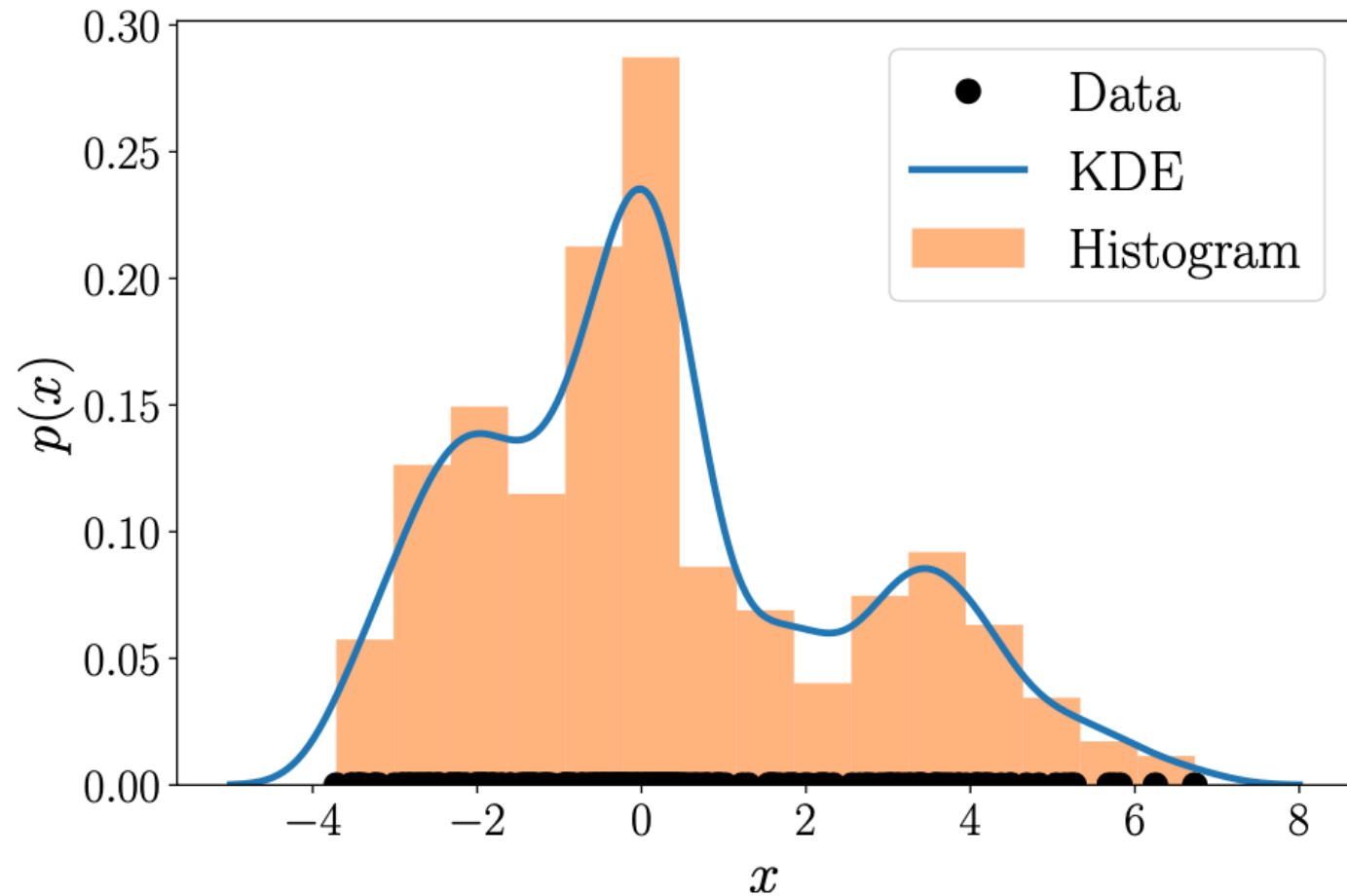
# KDE (Kernel Density Estimation)



Figure 11.9 from MML textbook

# Outline for November 30

- Gaussian Mixture Models (GMMs)

- Kernel Density Estimation (KDE)

- **Missing data**

- Begin: neural networks

# Types of missing data

- MCAR: Missing Completely At Random. Not related to:
    - Specific values
    - Observed responses


- MAR: Missing At Random. Not related to:
    - Specific values


- MNAR: Missing Not At Random

Reference: "The prevention and handling of the missing data" Kang (2013)

# Techniques for handling missing data

- Try to prevent the problem in the first place
  - Careful study design, follow-up with participants, etc

- Omit rows with missing data (reduces $n$)

- Omit only when value is needed
  - i.e. Naïve Bayes, per-feature estimates

- Mean substitution (per feature)

Reference: "The prevention and handling of the missing data" Kang (2013)
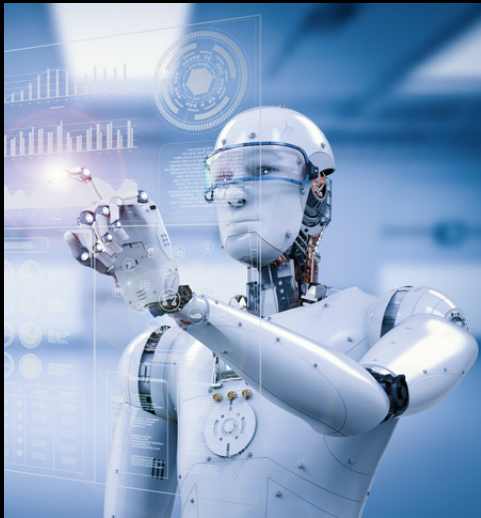
# Techniques for handling missing data

- Imputation
  - Use similar examples to guess the missing values
  - Can be done locally or globally


- Last observation carried forward
  - Useful for time-series data

Reference: "The prevention and handling of the missing data" Kang (2013)

# Outline for November 30

- Gaussian Mixture Models (GMMs)

- Kernel Density Estimation (KDE)

- Missing data

- Begin: neural networks

# Biological Inspiration

# Goal: learn from complicated inputs



$X_1$

$X_2$

$X_3$

$X_4$

$X_5$

$X_6$

input data

?

$Y_1$ glasses?

$Y_2$ smiling?

$Y_3$ identity?

parameters

Image: Labeled Faces in the Wild (UMass)

# Idea: transform data into lower dimension



Image: Labeled Faces in the Wild (UMass)

# Multi-layer networks = "deep learning"



input data     hidden layer 1     hidden layer 2     parameters

$Y_1$ glasses?

$Y_2$ smiling?

$Y_3$ identity?

# History of Neural Networks

- Perceptron can be interpreted as a simple neural network

- Misconceptions about the weaknesses of perceptrons contributed to declining funding for NN research

- Difficulty of training multi-layer NNs contributed to second setback

- Mid 2000's: breakthroughs in NN training contribute to rise of "deep learning"

# Number of papers that mention "deep learning" over time

# Big picture for today

- Neural networks can approximate any function!

# Big picture for today

- Neural networks can approximate any function!

- For our purposes in ML, we want to use them to approximate a function from our inputs to our outputs

# Big picture for today

- Neural networks can approximate any function!

- For our purposes in ML, we want to use them to approximate a function from our inputs to our outputs

- We will train our network by asking it to minimize the loss between its output and the true output
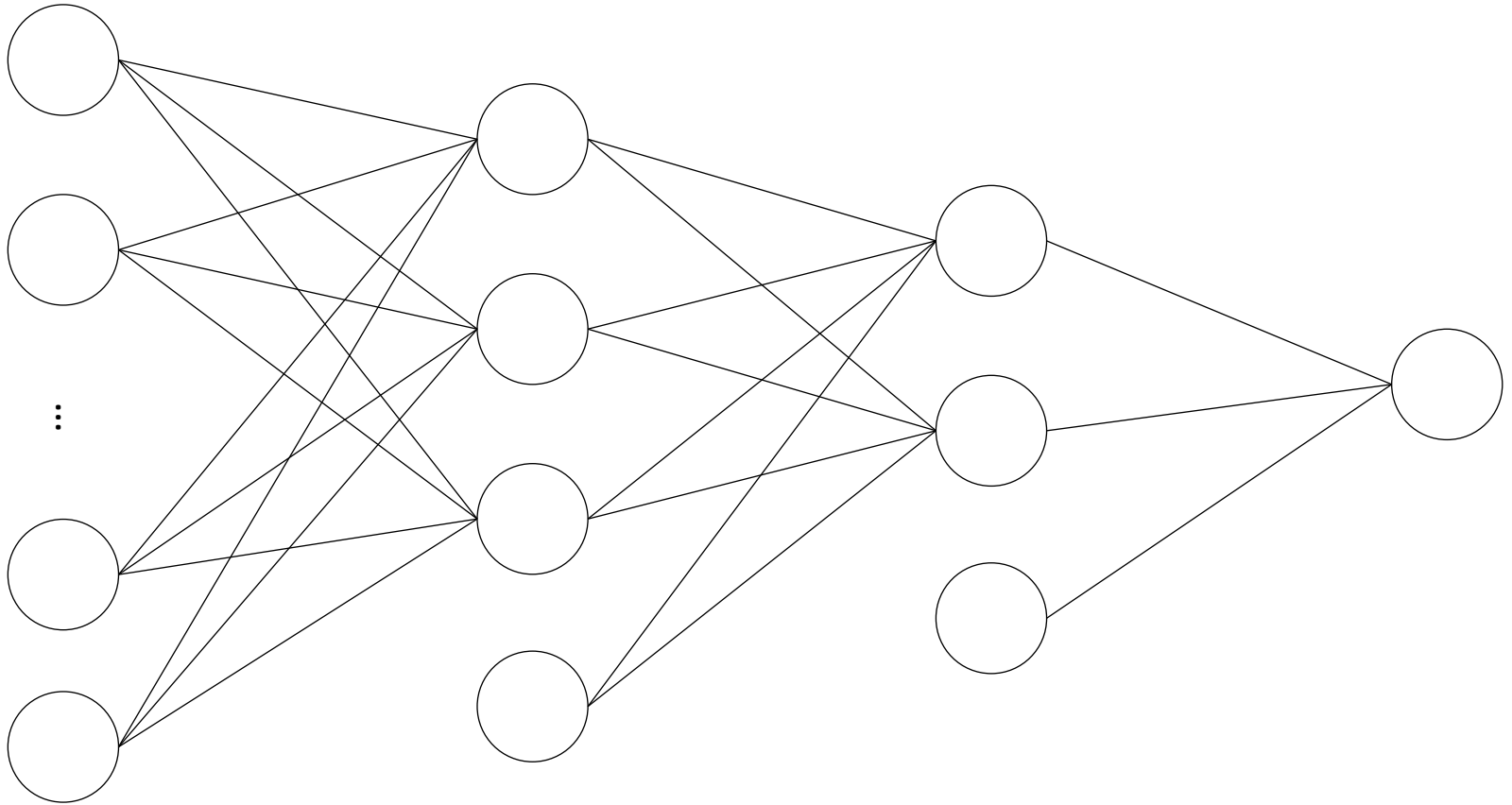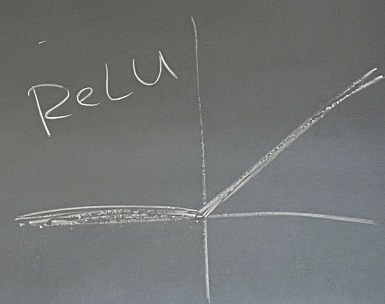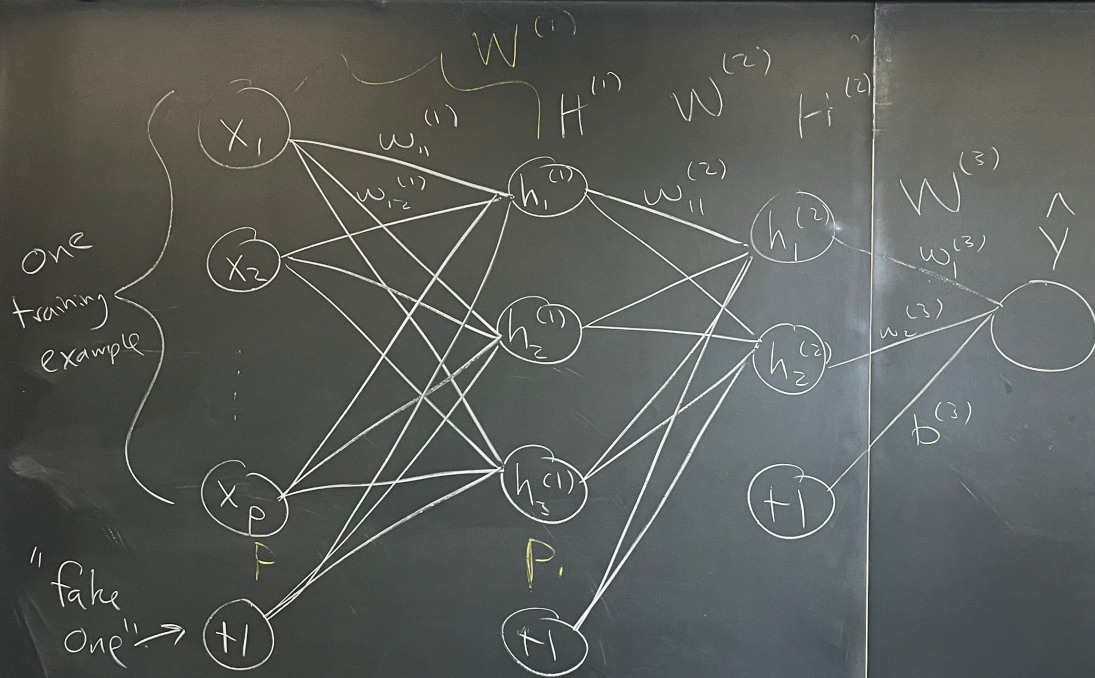
# Big picture for today

- Neural networks can approximate any function!

- For our purposes in ML, we want to use them to approximate a function from our inputs to our outputs

- We will train our network by asking it to minimize the loss between its output and the true output

- We will use SGD-like approaches to minimize loss

Fully Connected Neural Network Architecture

$W^{(1)}$

$X_1$

$W_{11}^{(1)}$

$W_{12}^{(1)}$

$H^{(1)}$ $W^{(2)}$ $H^{(2)}$

$h_1^{(1)}$ $W_{11}^{(2)}$

one training example

$X_2$

$W^{(3)}$ $\hat{y}$

$h_1^{(2)}$ $W_1^{(3)}$

$h_2^{(1)}$

$h_2^{(2)}$ $W_2^{(3)}$

$X_P$

$h_3^{(1)}$ $b^{(3)}$

+1

"fake one" → +1

+1

$P$

ReLU

$a$

$\begin{bmatrix} -1 & 2 & 5 \\ 0 & -2 & 1 \\ -3 & -4 & 8 \end{bmatrix}$

$\begin{bmatrix} 0 & 2 & 5 \\ 0 & 0 & 1 \\ 0 & 0 & 8 \end{bmatrix}$

$$H^{(1)} = a\left(W^{(1)}X + \vec{b}^{(1)}\right)$$

activation function

$p_1 \times p$    $p \times n$     $p_1 \times 1$

$p_1 \times n$

$$\begin{bmatrix} | & | \\ x_1 & x_2 \cdots \\ | & | \end{bmatrix}$$
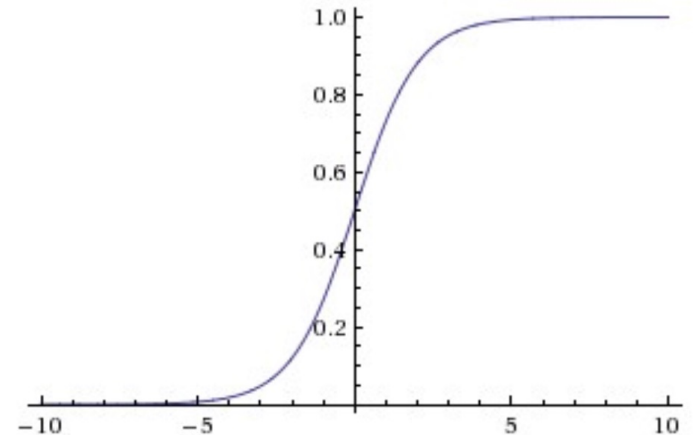
$p \times n$

$$H^{(2)} = a\left(W^{(2)}H^{(1)} + \vec{b}^{(2)}\right)$$

$$\vec{y} = a\left(W^{(3)}H^{(2)} + b^{(3)}\right)$$

# Option 1: sigmoid function
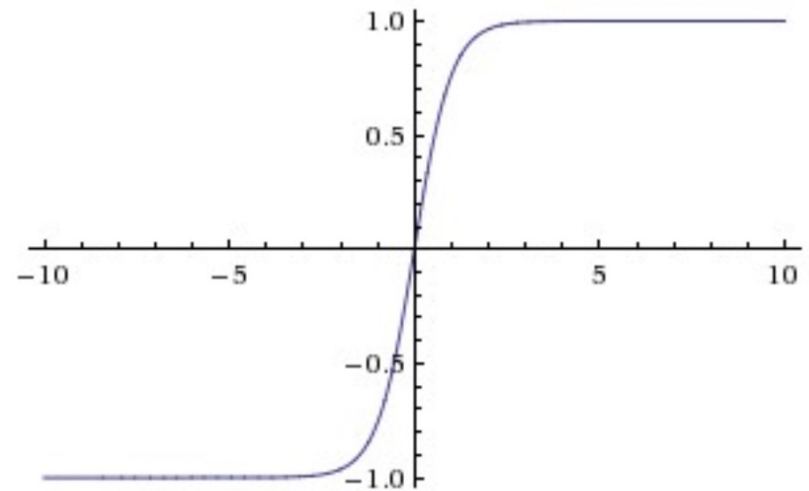
- Input: all real numbers, output: [0, 1]

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

# Option 2: hyperbolic tangent

- Input: all real numbers, output: [-1, 1]

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

# Option 3: Rectified Linear Unit (ReLU)

- Return *x* if *x* is positive (i.e. threshold at 0)

$$f(x) = \max(0, x)$$

# Pros and Cons of Activation Functions

1) Sigmoid

- (-) When input becomes very positive or very negative, gradient approaches 0 (saturates and stops gradient descent)
- (-) Not zero-centered, so gradient on weights can end up all positive or all negative (zig-zag in gradient descent)
- (+) Derivative is easy to compute given function value!

# Pros and Cons of Activation Functions

## 1) Sigmoid

- (-) When input becomes very positive or very negative, gradient approaches 0 (saturates and stops gradient descent)
- (-) Not zero-centered, so gradient on weights can end up all positive or all negative (zig-zag in gradient descent)
- (+) Derivative is easy to compute given function value!

## 2) Tanh

- (-) Still has a tendency to prematurely kill the gradient
- (+) Zero-centered so we get a range of gradients
- (+) Rescaling of sigmoid function so derivative is also not too difficult

# Pros and Cons of Activation Functions

## 1) Sigmoid

- (-) When input becomes very positive or very negative, gradient approaches 0 (saturates and stops gradient descent)
- (-) Not zero-centered, so gradient on weights can end up all positive or all negative (zig-zag in gradient descent)
- (+) Derivative is easy to compute given function value!
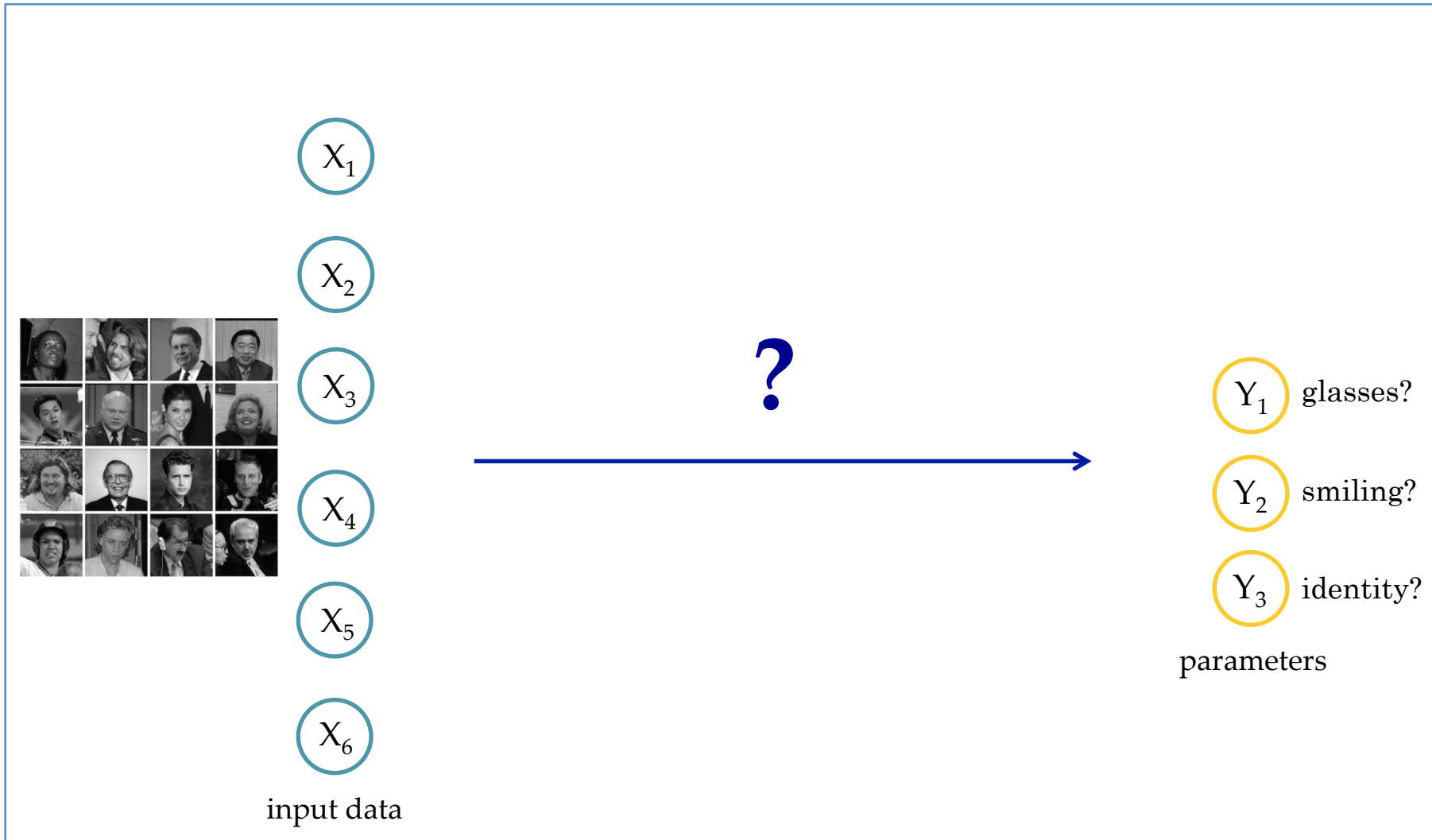
## 2) Tanh

- (-) Still has a tendency to prematurely kill the gradient
- (+) Zero-centered so we get a range of gradients
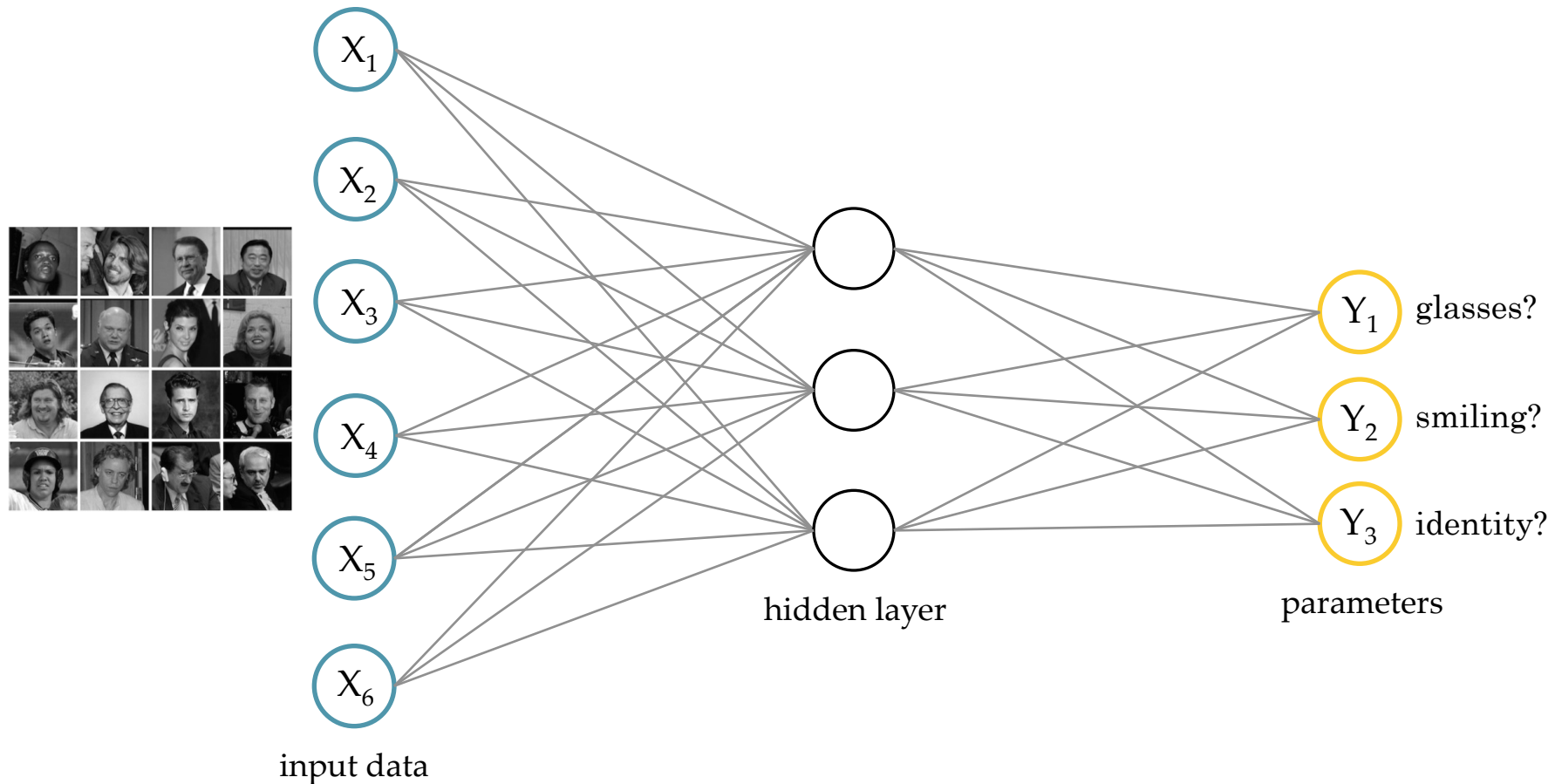- (+) Rescaling of sigmoid function so derivative is also not too difficult

## 3) ReLU

- (+) Works well in practice (accelerates convergence)
- (+) Function value very easy to compute! (no exponentials)
- (-) Units can have no signal if input becomes too negative throughout gradient descent
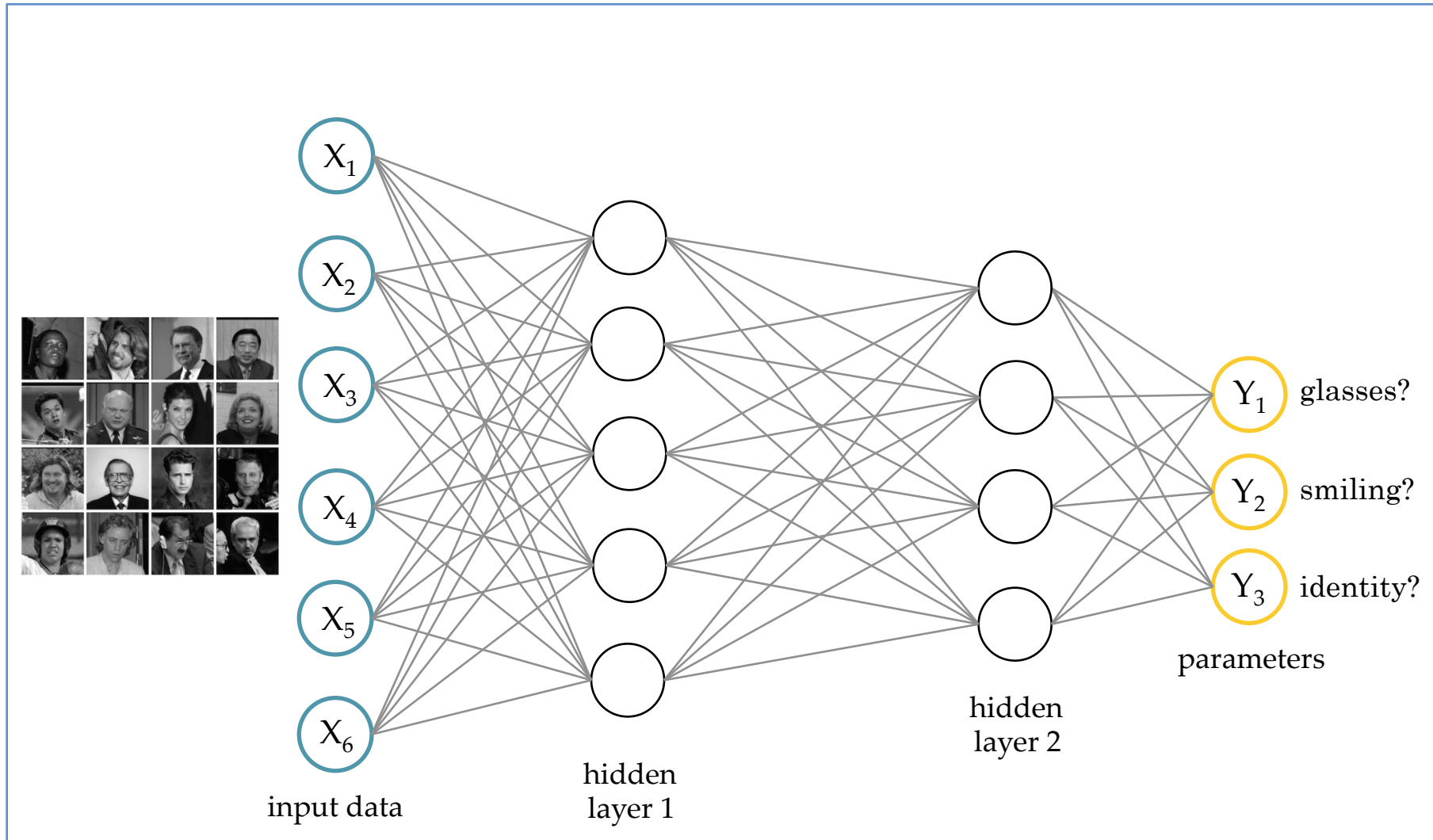
More info:
http://cs231n.github.io/neural-networks-1/

# Goal: find a function between input and output



X₁

X₂

X₃

X₄

X₅

X₆

input data

?

Y₁  glasses?

Y₂  smiling?

Y₃  identity?

parameters

# First idea: one hidden layer



input data      hidden layer      parameters

$X_1$ $X_2$ $X_3$ $X_4$ $X_5$ $X_6$

$Y_1$ glasses?
$Y_2$ smiling?
$Y_3$ identity?

# Second idea: more hidden layers ("deep" learning)



input data      hidden layer 1      hidden layer 2      parameters

$Y_1$ glasses?

$Y_2$ smiling?

$Y_3$ identity?

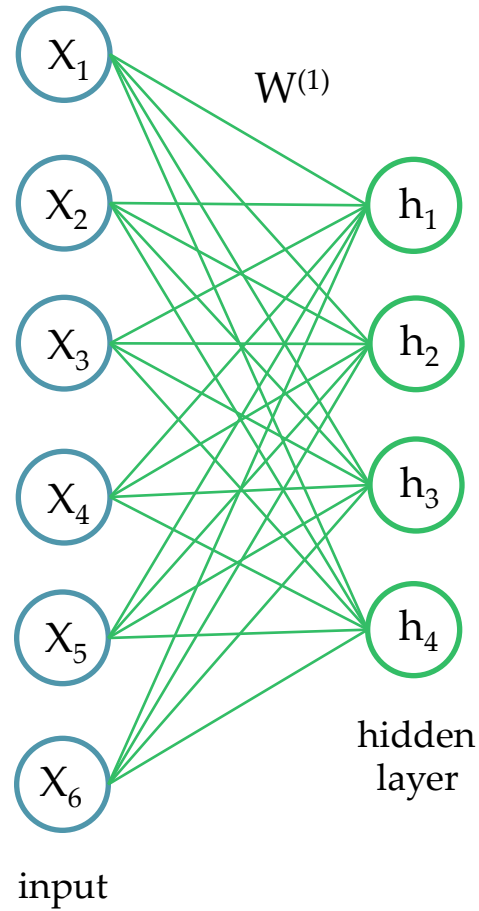# Flatten pixels of image into a single vector



input data

# Detour to autoencoders



$X_1$

$X_2$

$X_3$

$X_4$

$X_5$

$X_6$

input

# Detour to autoencoders



$X_1$

$X_2$

$X_3$

$X_4$

$X_5$

$X_6$

input

$W^{(1)}$

$h_1$

$h_2$

$h_3$

$h_4$

hidden
layer

# Detour to autoencoders



$W^{(1)}$     $W^{(2)}$

$X_1$   $h_1$   $X_1^*$

$X_2$   $h_2$   $X_2^*$

$X_3$   $h_3$   $X_3^*$

$X_4$   $h_4$   $X_4^*$

$X_5$    $X_5^*$

$X_6$    $X_6^*$

hidden layer

input

reconstructed input

# Use <u>unsupervised pre-training</u> to find a function from the input to itself



input data      hidden layer 1      reconstructed input

# Hidden units can be interpreted as edges



input data                hidden layer 1          reconstructed input

# Now: throw away reconstruction and input



input data

hidden layer 1

# Now: throw away reconstruction and input



hidden
layer 1
second layer

# Then repeat the entire process for each layer



$g_1$

$g_2$

$g_3$

$g_4$

$h_1^*$

$h_2^*$

$h_3^*$

$h_4^*$

$h_5^*$

hidden
layer 2

reconstructed
input

hidden
layer 1

# Then repeat the entire process for each layer



hidden layer 1

hidden layer 2

$h_1^*$  $h_2^*$  $h_3^*$  $h_4^*$  $h_5^*$

reconstructed input

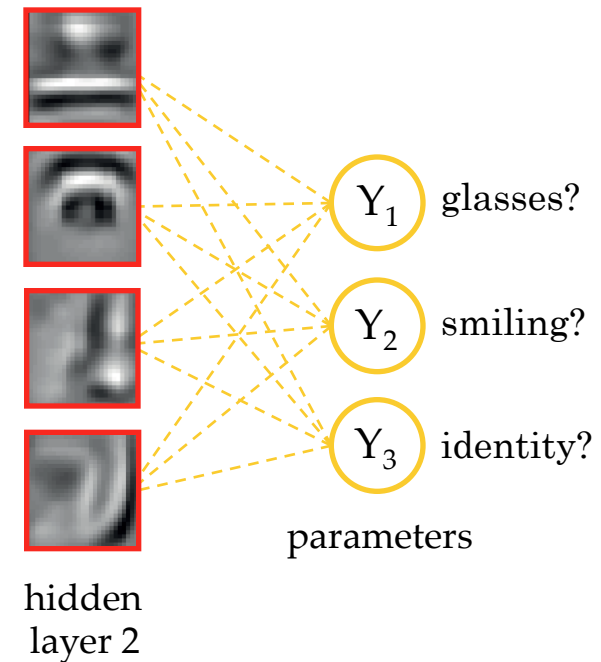# Then repeat the entire process for each layer



hidden layer 1

hidden layer 2

# Then repeat the entire process for each layer



hidden
layer 2

# In the last layer, use the outputs (supervised)



hidden
layer 2

$Y_1$ glasses?

$Y_2$ smiling?

$Y_3$ identity?

parameters

# In the last layer, use the outputs (supervised)



hidden
layer 2

$Y_1$  glasses?

$Y_2$  smiling?

$Y_3$  identity?

parameters

# Finally, "fine-tune" the entire network!



input data

hidden layer 1

hidden layer 2

$Y_1$ glasses?

$Y_2$ smiling?

$Y_3$ identity?

parameters