



# 10-601 Introduction to Machine Learning

Machine Learning Department  
School of Computer Science  
Carnegie Mellon University

## Neural Networks and Backpropagation

### Neural Net Readings:

Murphy --  
Bishop 5  
HTF 11  
Mitchell 4

Matt Gormley  
Lecture 19  
March 29, 2017

# Reminders

- **Homework 6: Unsupervised Learning**
  - Release: Wed, Mar. 22
  - Due: Mon, Apr. 03 at 11:59pm
- **Homework 5 (Part II): Peer Review**
  - Release: Wed, Mar. 29
  - Due: Wed, Apr. 05 at 11:59pm
- **Peer Tutoring**

Expectation: You should spend at most 1 hour on your reviews

# Neural Networks Outline

- **Logistic Regression (Recap)**
  - Data, Model, Learning, Prediction
- **Neural Networks**
  - A Recipe for Machine Learning
  - Visual Notation for Neural Networks
  - Example: Logistic Regression Output Surface
  - 2-Layer Neural Network
  - 3-Layer Neural Network
- **Neural Net Architectures**
  - Objective Functions
  - Activation Functions
- **Backpropagation**
  - Basic Chain Rule (of calculus)
  - Chain Rule for Arbitrary Computation Graph
  - Backpropagation Algorithm
  - Module-based Automatic Differentiation (Autodiff)

# **RECALL: LOGISTIC REGRESSION**

# Using gradient ascent for linear classifiers

Recall...

Key idea behind today's lecture:

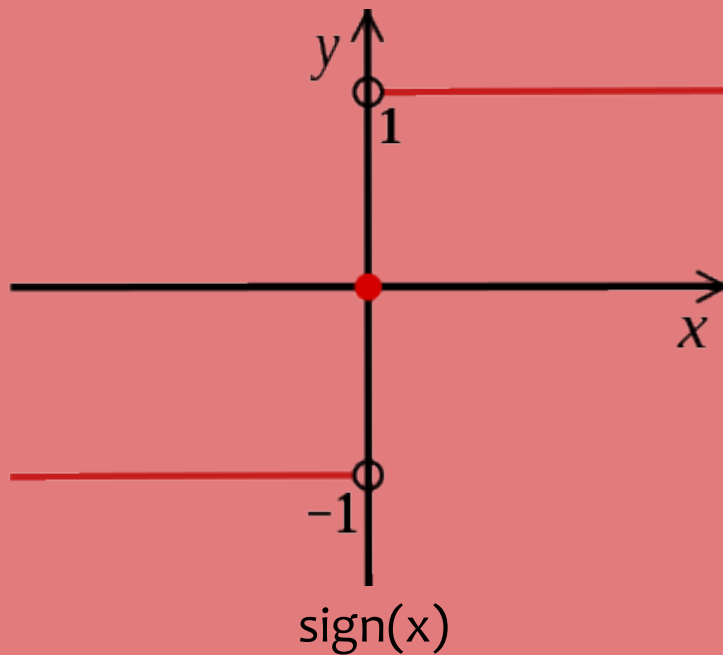
1. Define a linear classifier (logistic regression)
2. Define an objective function (likelihood)
3. Optimize it with gradient descent to learn parameters
4. Predict the class with highest probability under the model

# Using gradient ascent for linear classifiers

Recall...

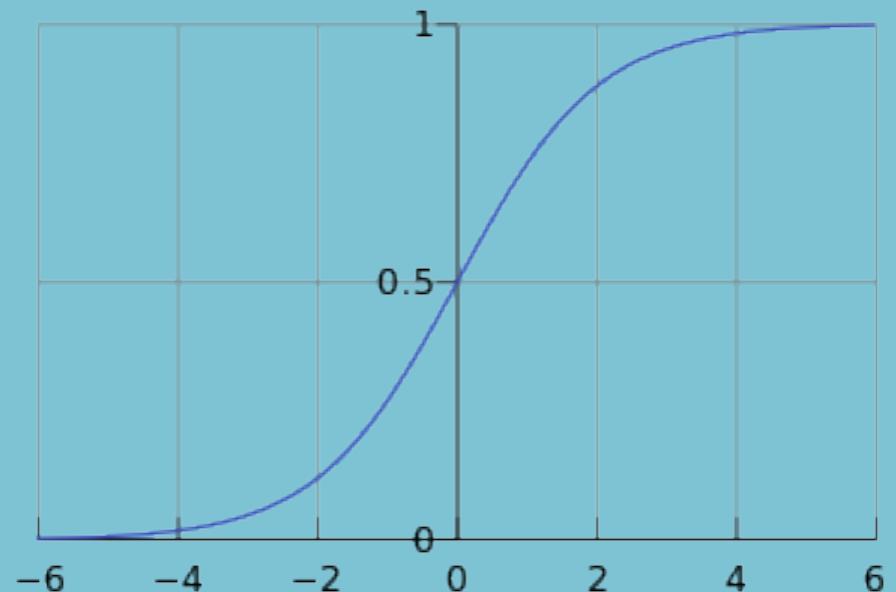
This decision function isn't differentiable:

$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x})$$



Use a differentiable function instead:

$$p_{\boldsymbol{\theta}}(y = 1 | \mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$



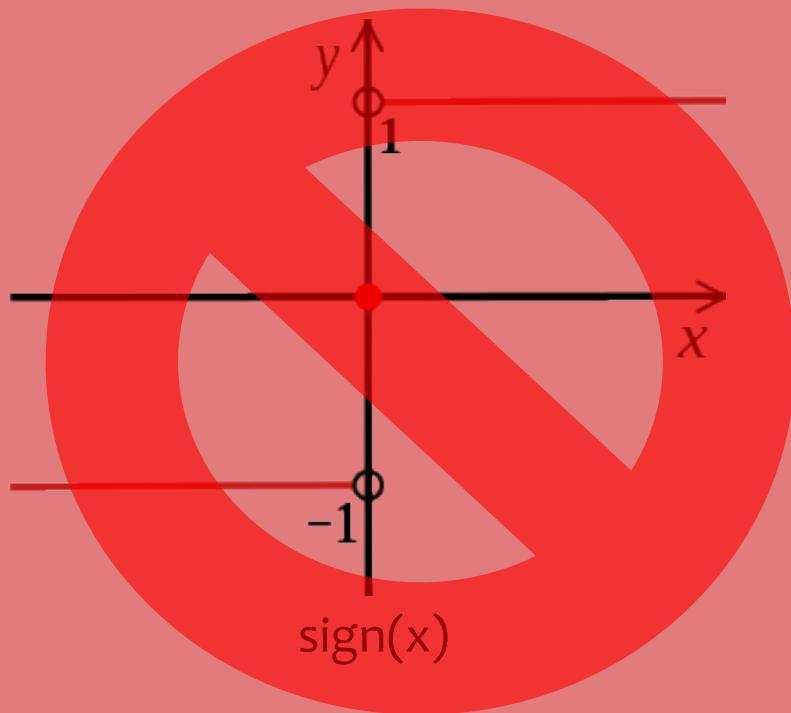
$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$

# Using gradient ascent for linear classifiers

Recall...

This decision function isn't differentiable:

$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x})$$



Use a differentiable function instead:

$$p_{\boldsymbol{\theta}}(y = 1 | \mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$



$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$

# Logistic Regression

**Data:** Inputs are continuous vectors of length  $K$ . Outputs are discrete.

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N \text{ where } \mathbf{x} \in \mathbb{R}^K \text{ and } y \in \{0, 1\}$$

**Model:** Logistic function applied to dot product of parameters with input vector.

$$p_{\boldsymbol{\theta}}(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$

**Learning:** finds the parameters that minimize some objective function.  $\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta})$

**Prediction:** Output is the most probable class.

$$\hat{y} = \underset{y \in \{0, 1\}}{\operatorname{argmax}} p_{\boldsymbol{\theta}}(y|\mathbf{x})$$



# NEURAL NETWORKS

## Background

# A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\theta}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

*Face*



*Face*



*Not a face*



**Examples:** Linear regression,  
Logistic regression, Neural Network

**Examples:** Mean-squared error,  
Cross Entropy

## Background

# A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

4. Train with SGD:

(take small steps  
opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

## Background

1. Given training data

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of the

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

## A Recipe for Gradients

**Backpropagation** can compute this gradient!

And it's a **special case of a more general algorithm** called reverse-mode automatic differentiation that can compute the gradient of any differentiable function efficiently!

opposite the gradient)


$$\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

# Goals for Today's Lecture

1. Explore a **new class of decision functions** (Neural Networks)
2. Consider **variants of this recipe** for training

– Decision function

$$\hat{y} = f_{\theta}(x_i)$$

– Loss function

$$\ell(\hat{y}, y_i) \in \mathbb{R}$$

4. Train with SGD:

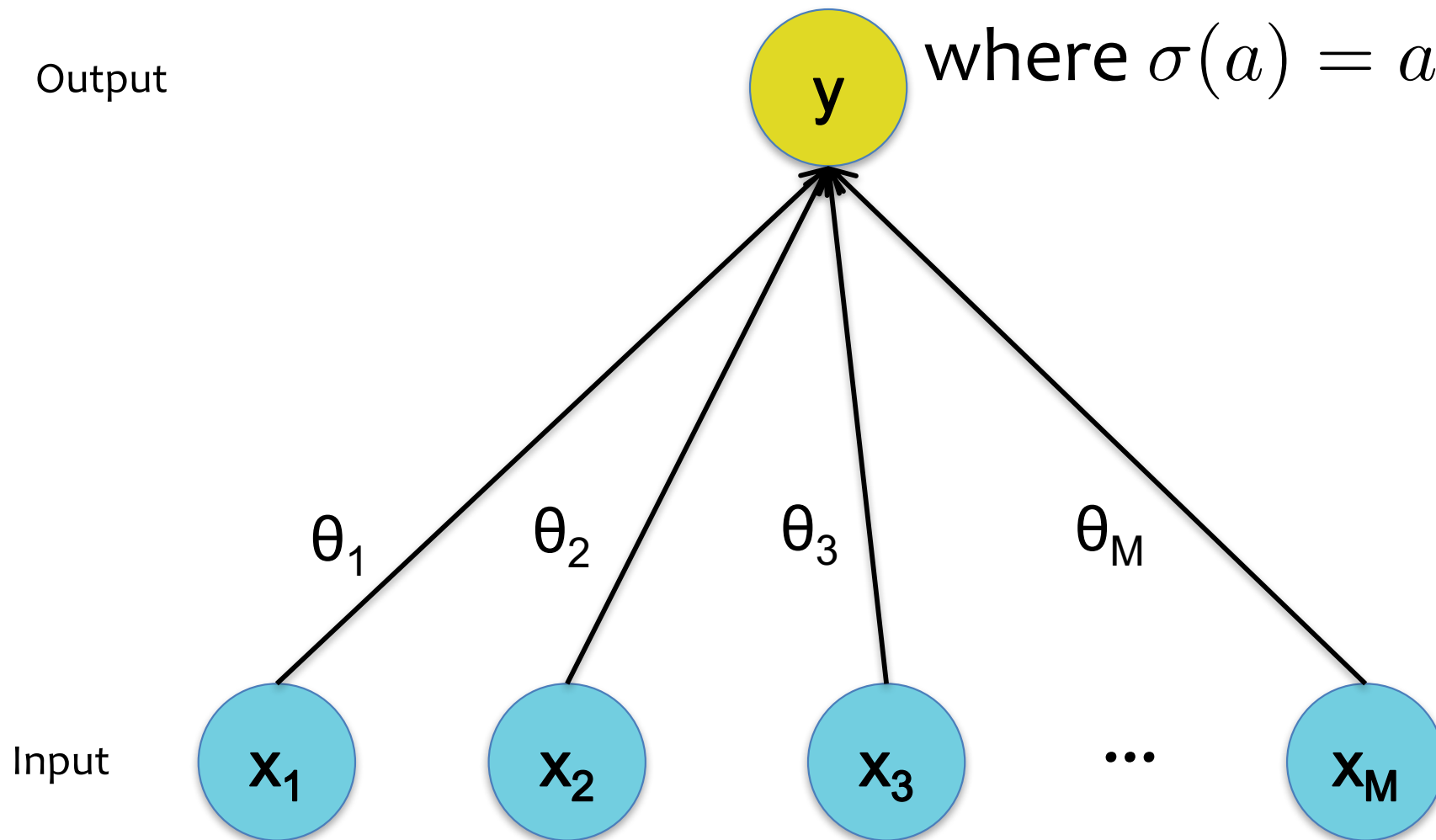
– Take small steps opposite the gradient)

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(x_i), y_i)$$

$$y = h_{\theta}(\mathbf{x}) = \sigma(\boldsymbol{\theta}^T \mathbf{x})$$

where  $\sigma(a) = a$

Output

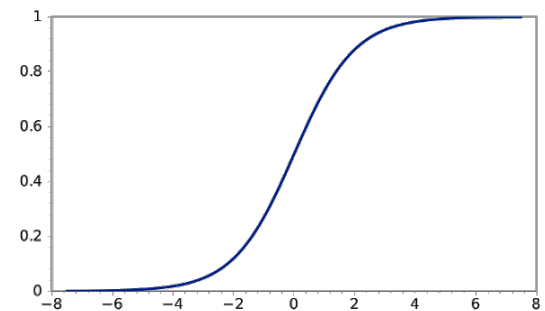
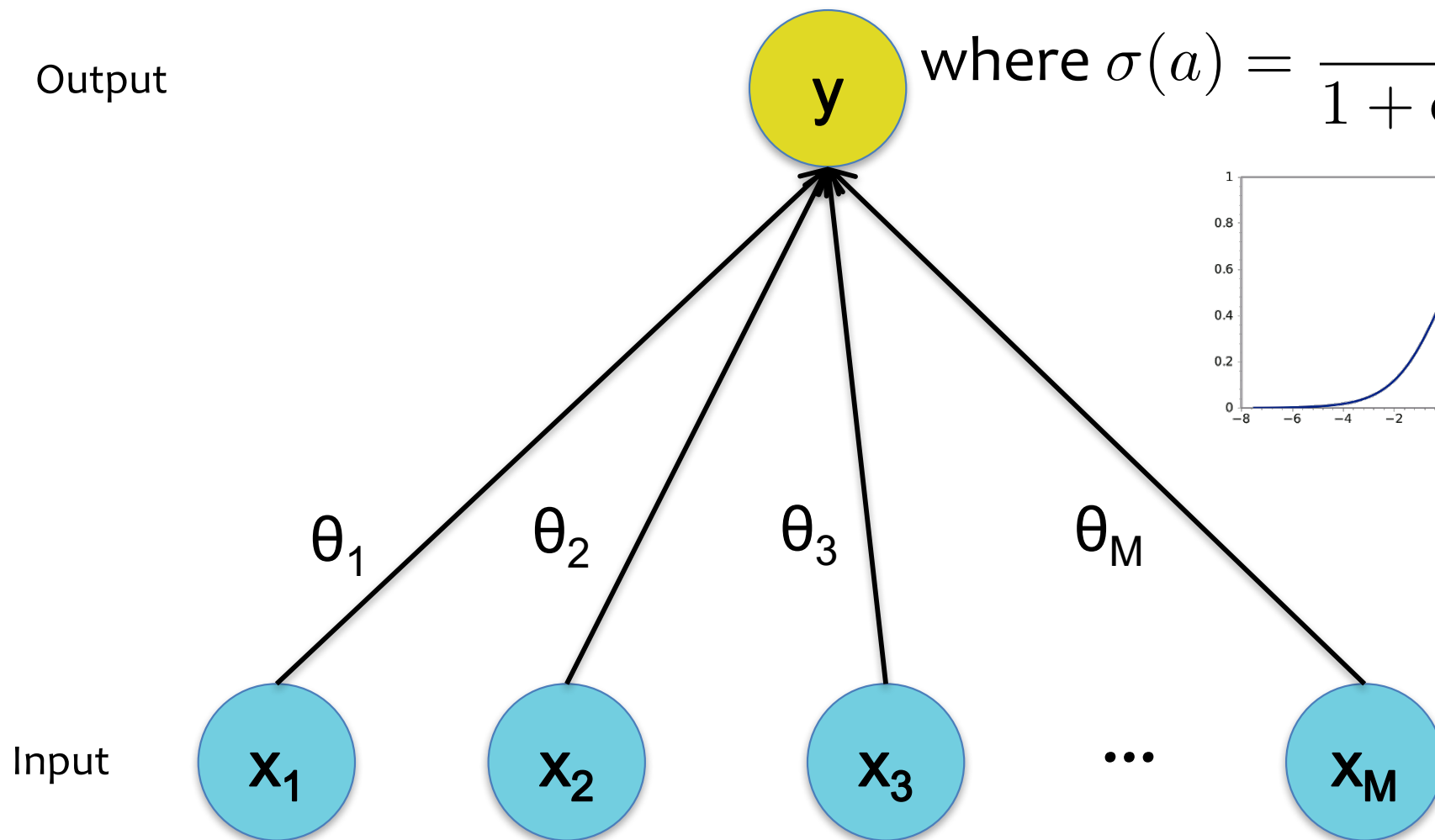


# Logistic Regression

$$y = h_{\theta}(x) = \sigma(\theta^T x)$$

$$\text{where } \sigma(a) = \frac{1}{1 + \exp(-a)}$$

Output



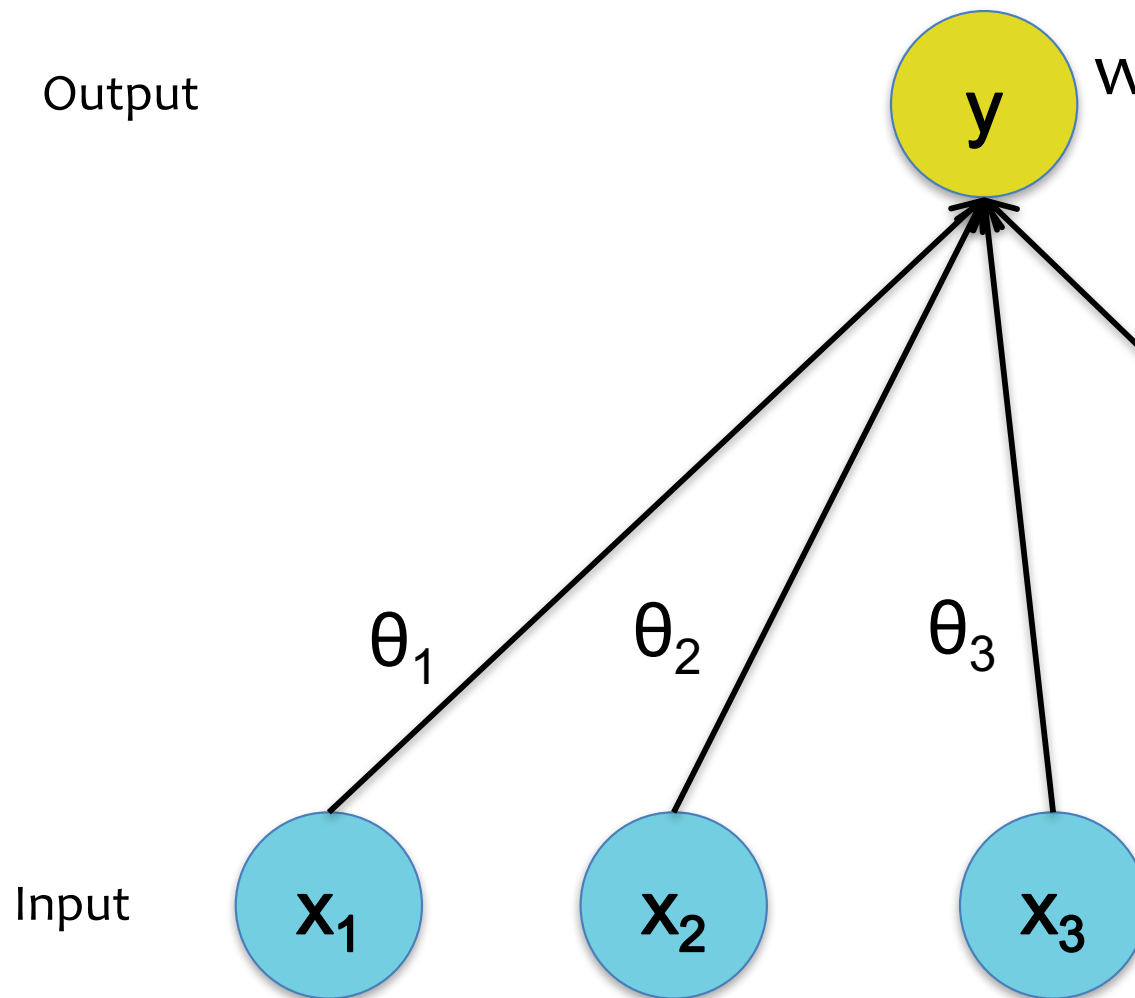
## Decision Functions

# Logistic Regression

$$y = h_{\theta}(x) = \sigma(\theta^T x)$$

$$\text{where } \sigma(a) = \frac{1}{1 + \exp(-a)}$$

Output

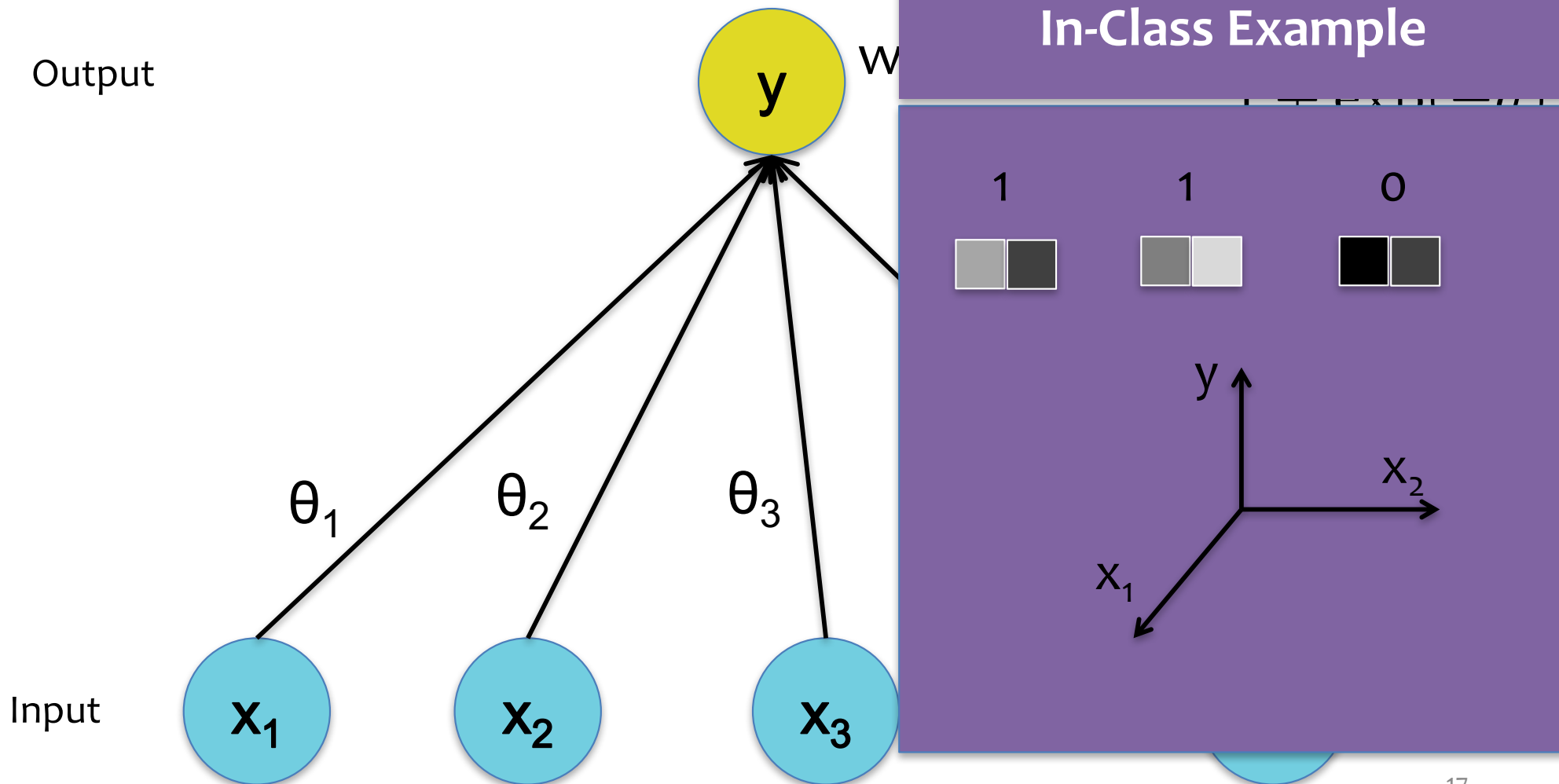




## Decision Functions

# Logistic Regression

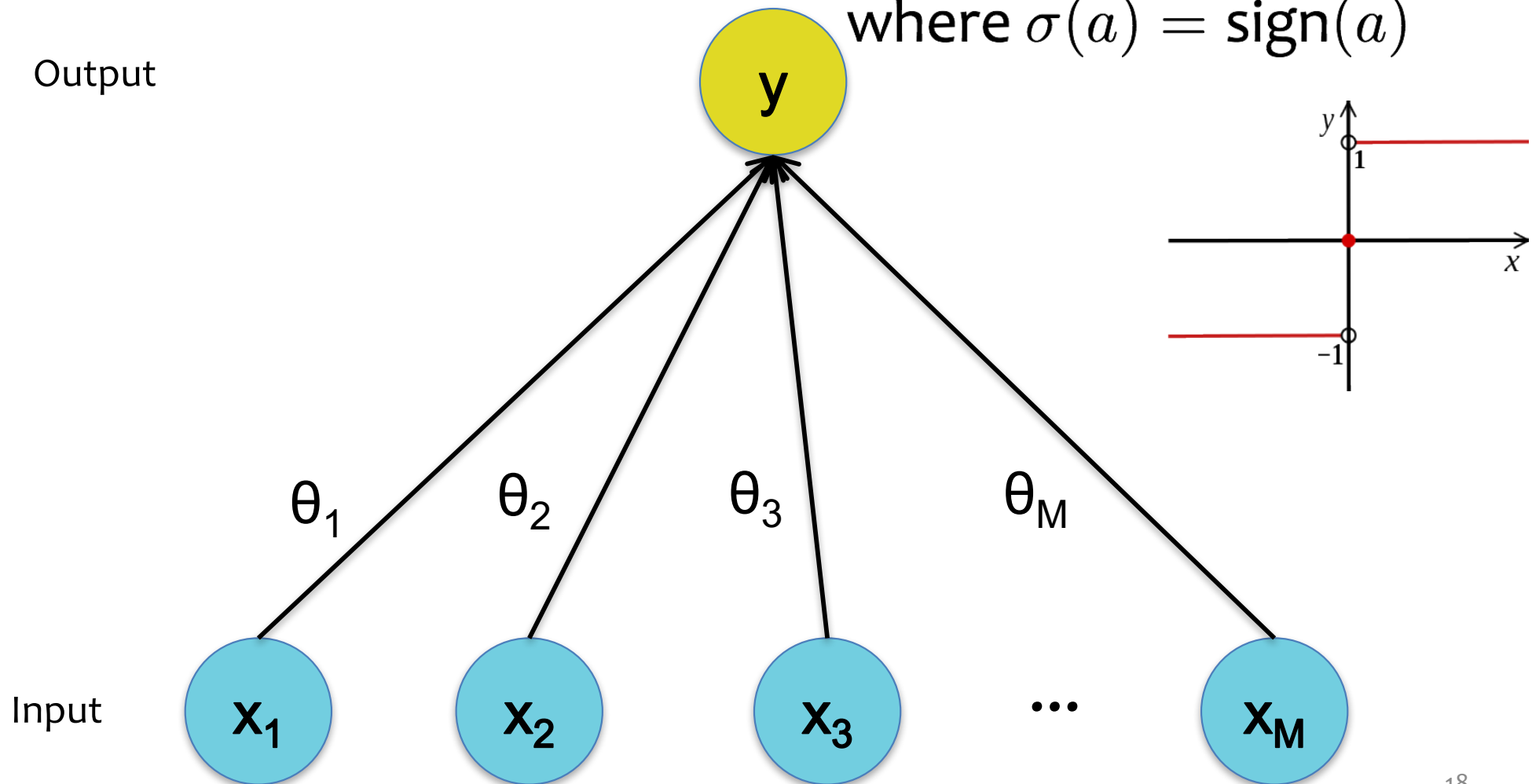
$$y = h_{\theta}(x) = \sigma(\theta^T x)$$



$$y = h_{\theta}(\mathbf{x}) = \sigma(\theta^T \mathbf{x})$$

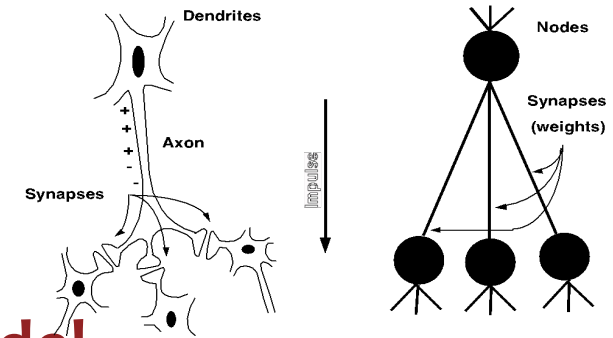
$$\text{where } \sigma(a) = \text{sign}(a)$$

Output



# From Biological to Artificial

The motivation for Artificial Neural Networks comes from biology...



## Biological “Model”

- **Neuron:** an excitable cell
- **Synapse:** connection between neurons
- A neuron sends an **electrochemical pulse** along its synapses when a sufficient voltage change occurs
- **Biological Neural Network:** collection of neurons along some pathway through the brain

## Biological “Computation”

- Neuron switching time :  $\sim 0.001$  sec
- Number of neurons:  $\sim 10^{10}$
- Connections per neuron:  $\sim 10^{4-5}$
- Scene recognition time:  $\sim 0.1$  sec

## Artificial Model

- **Neuron:** node in a directed acyclic graph (DAG)
- **Weight:** multiplier on each edge
- **Activation Function:** nonlinear thresholding function, which allows a neuron to “fire” when the input value is sufficiently high
- **Artificial Neural Network:** collection of neurons into a DAG, which define some differentiable function

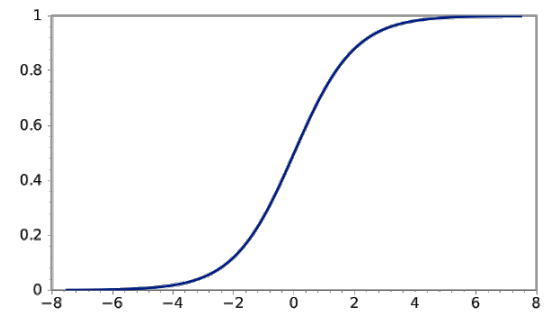
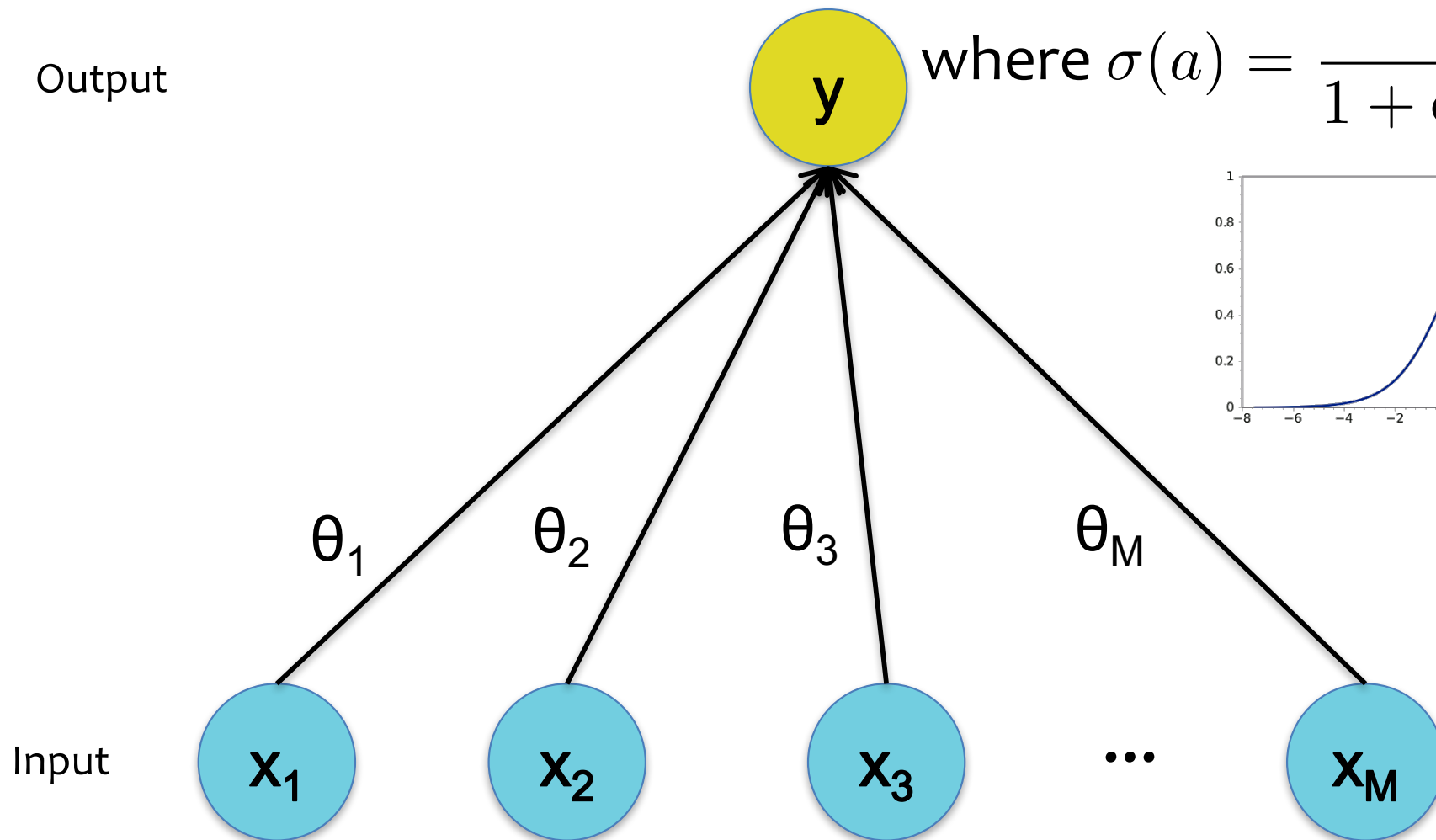
## Artificial Computation

- Many neuron-like threshold switching units
- Many weighted interconnections among units
- Highly parallel, distributed processes

$$y = h_{\theta}(x) = \sigma(\theta^T x)$$

$$\text{where } \sigma(a) = \frac{1}{1 + \exp(-a)}$$

Output

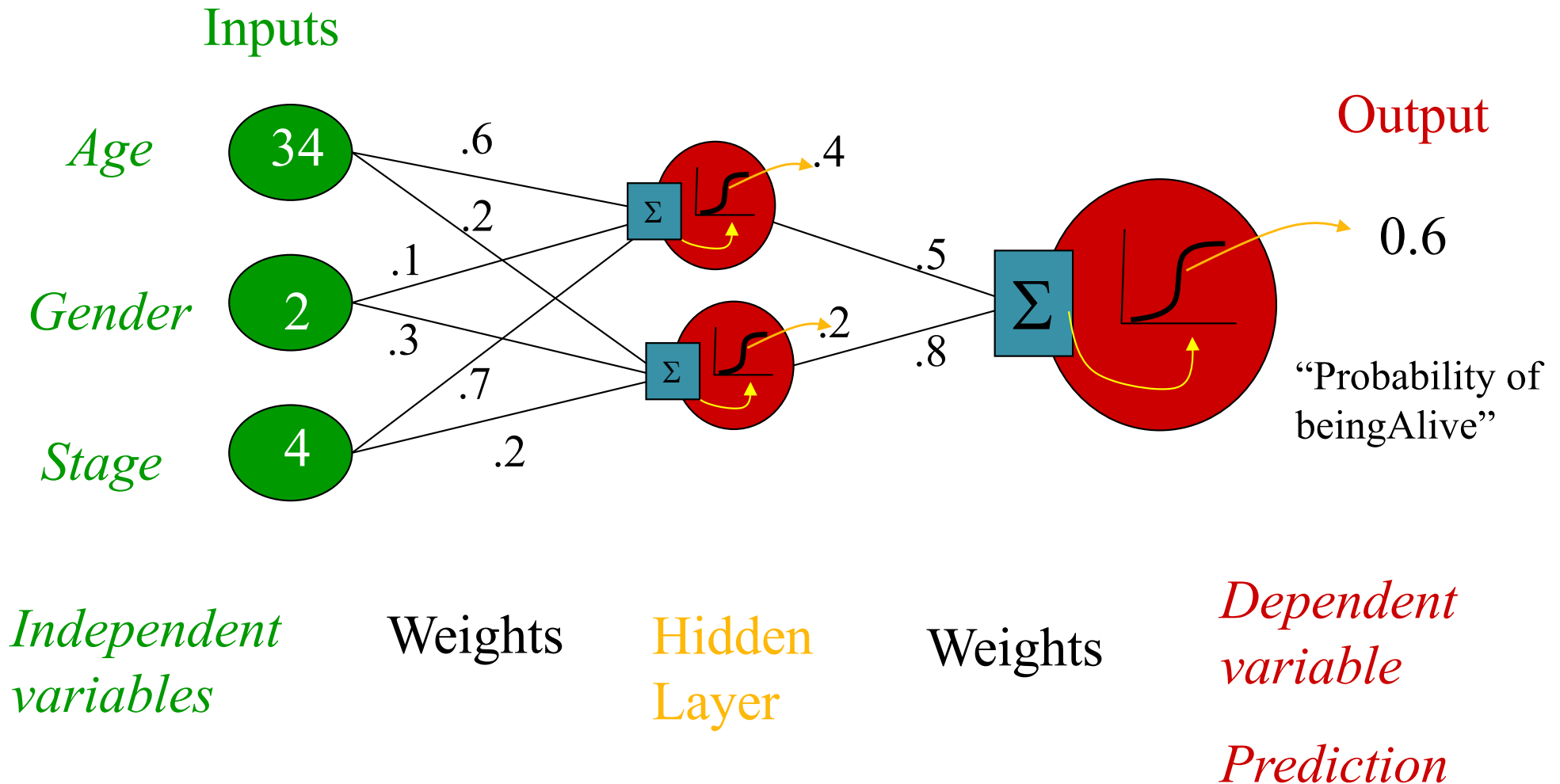


# Neural Networks

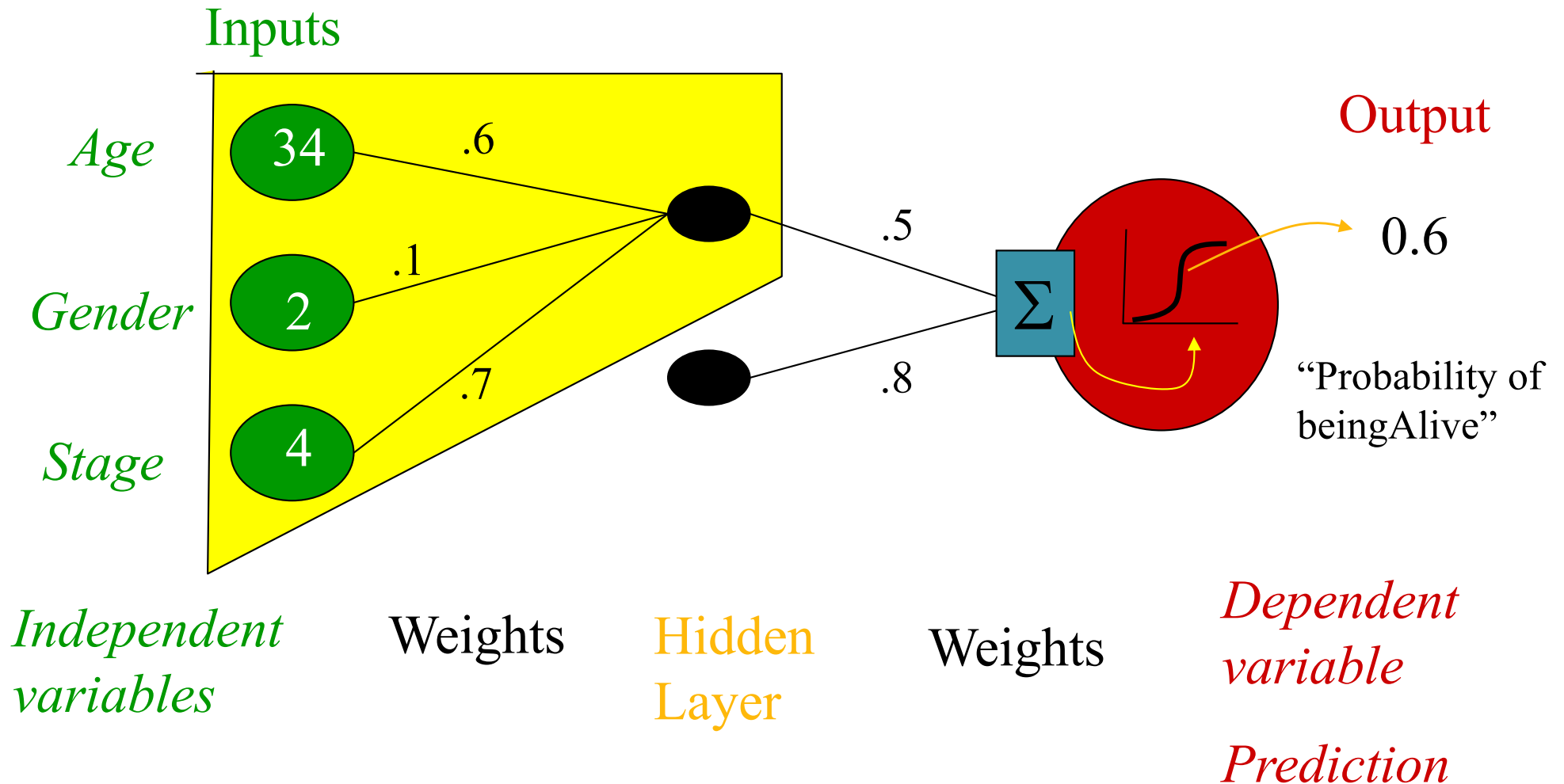
## *Whiteboard*

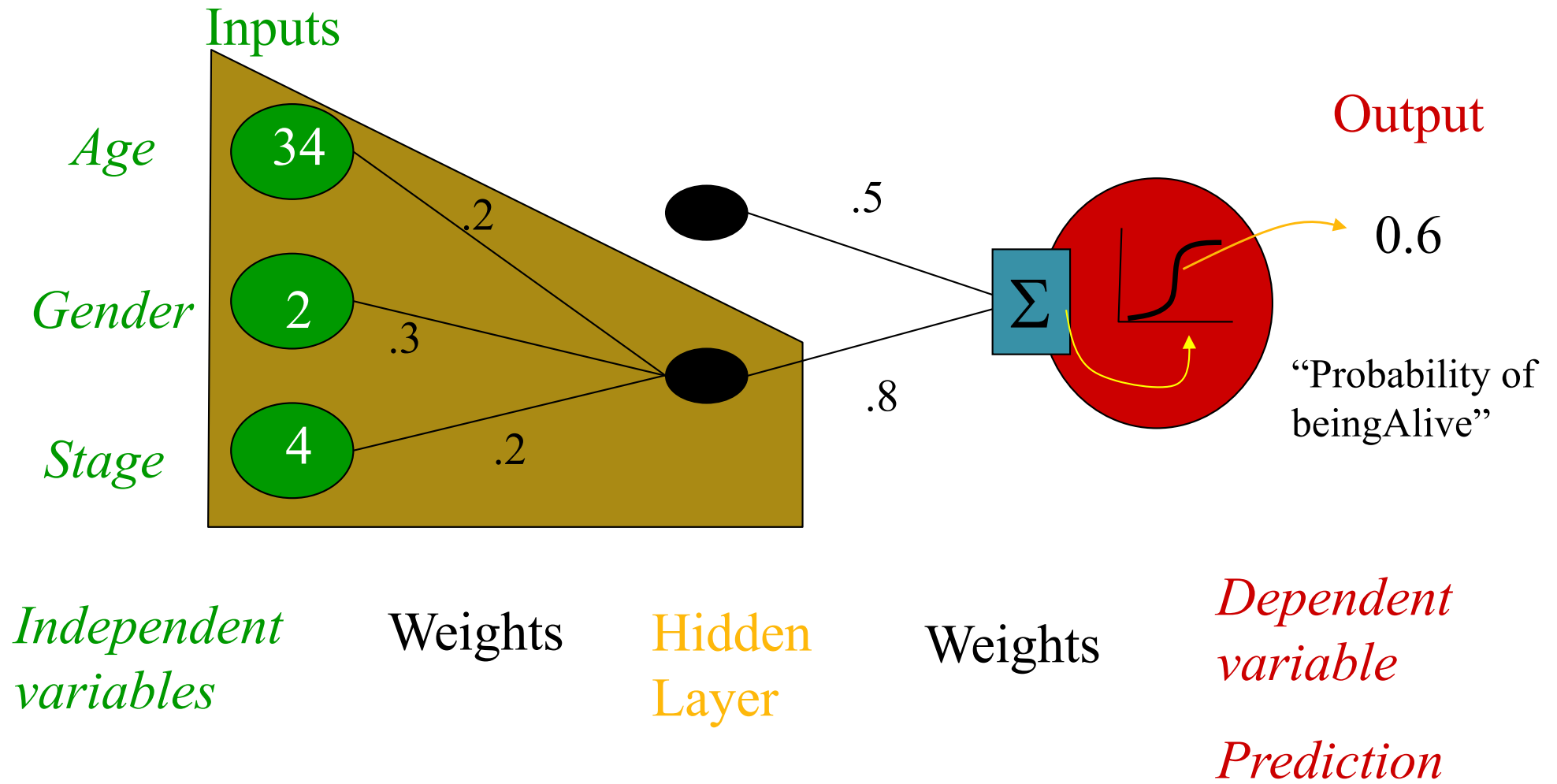
- Example: Neural Network w/1 Hidden Layer
- Example: Neural Network w/2 Hidden Layers
- Example: Feed Forward Neural Network

# Neural Network Model

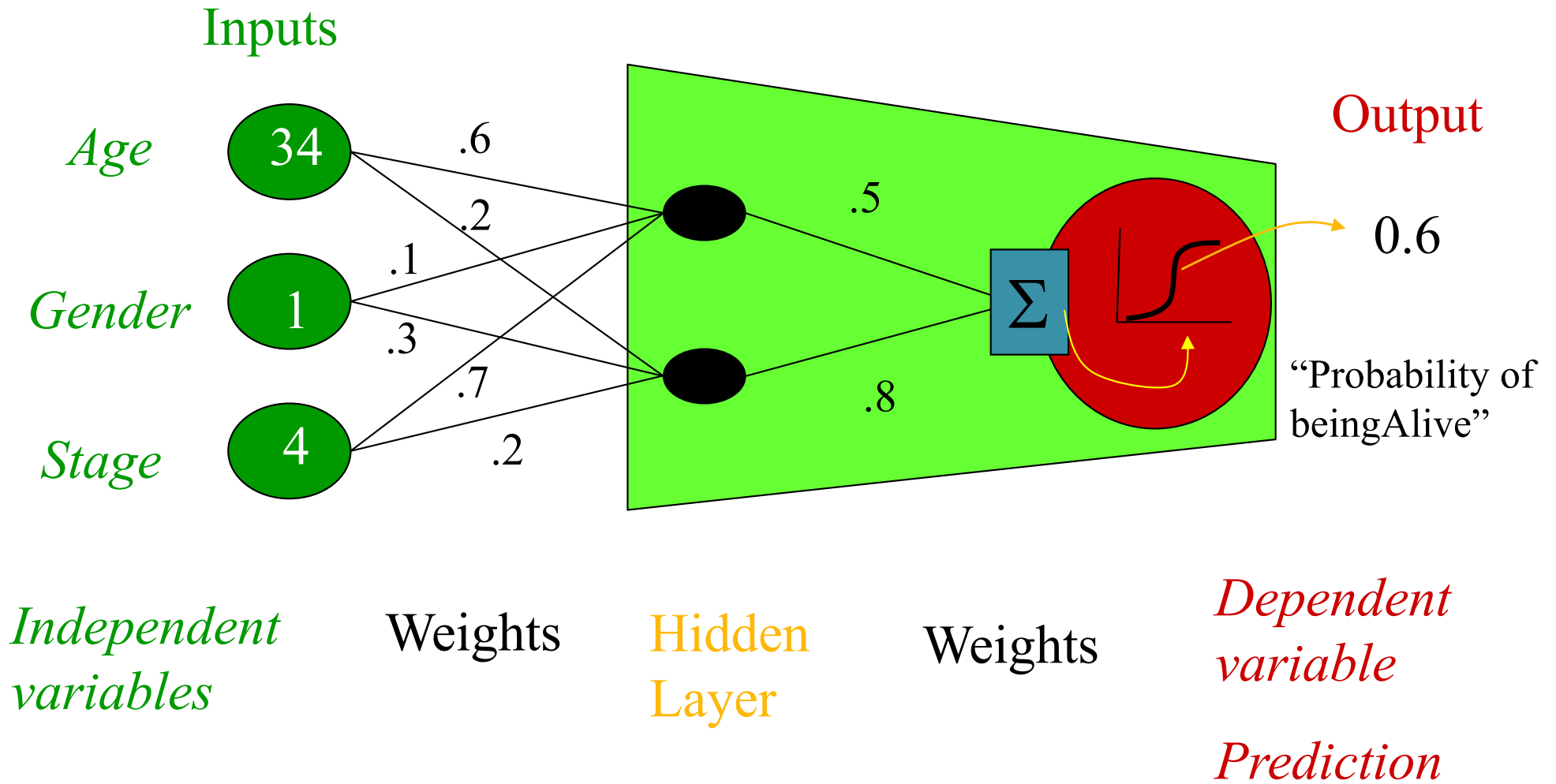


# “Combined logistic models”

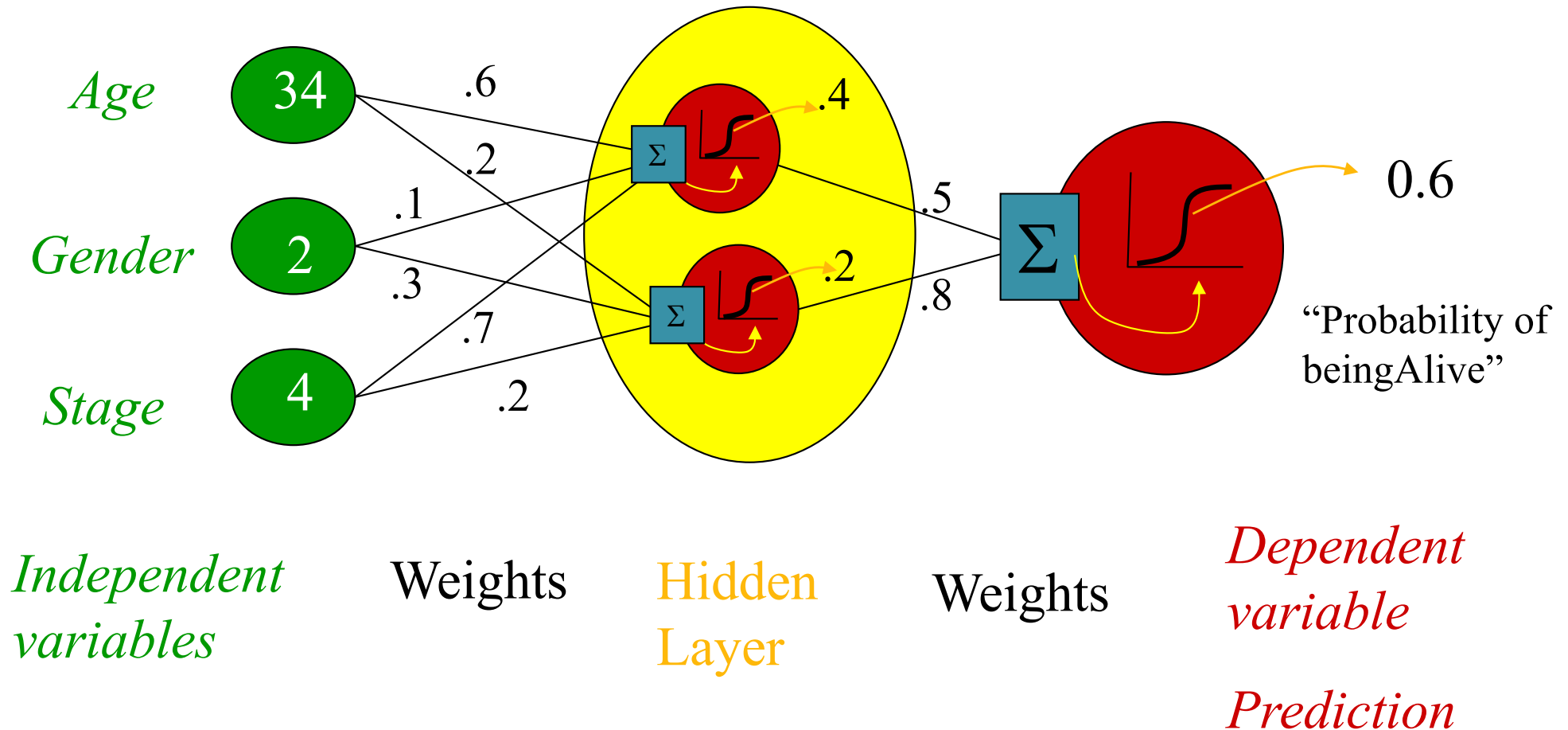






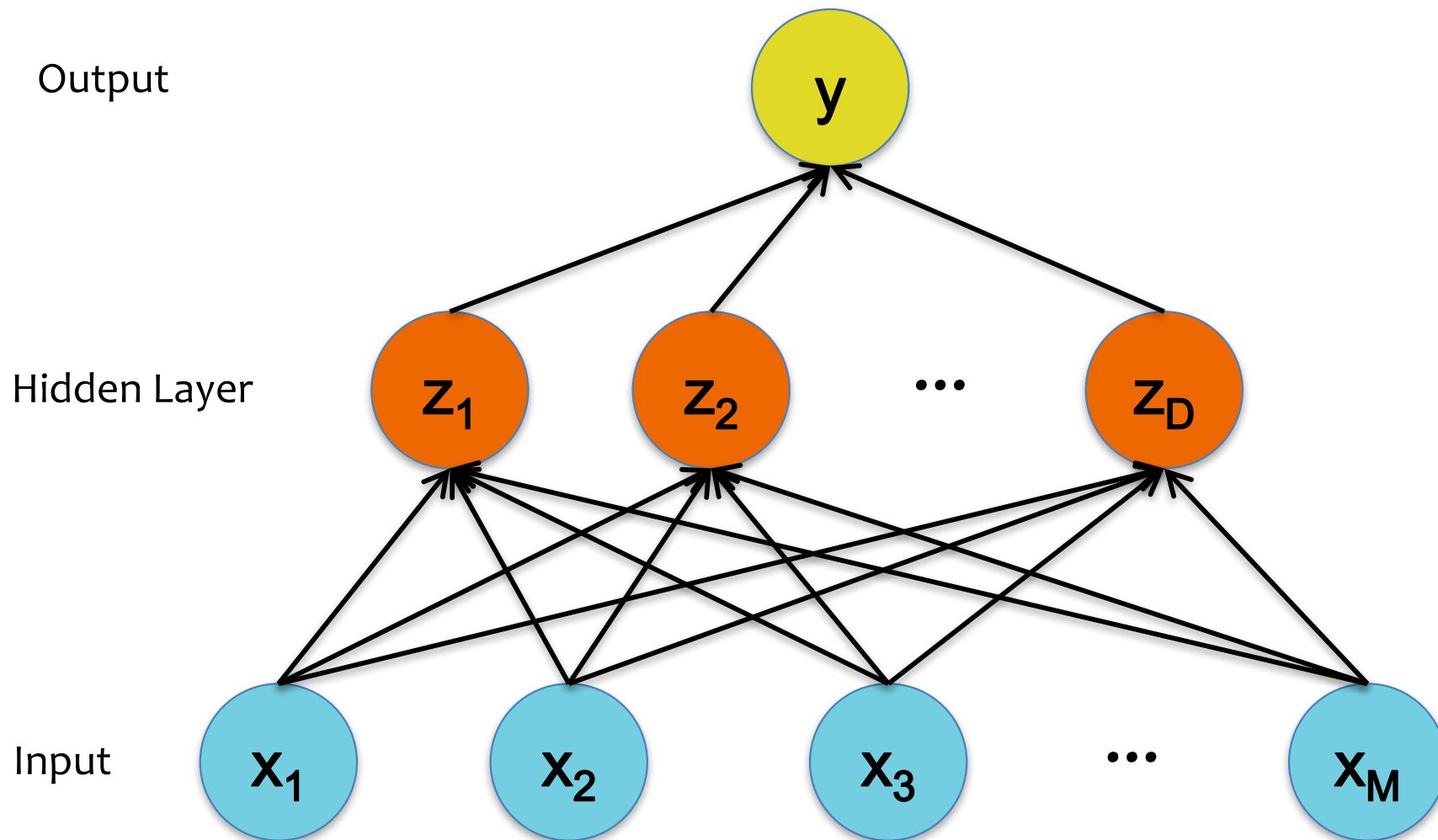


# Not really, no target for hidden units...



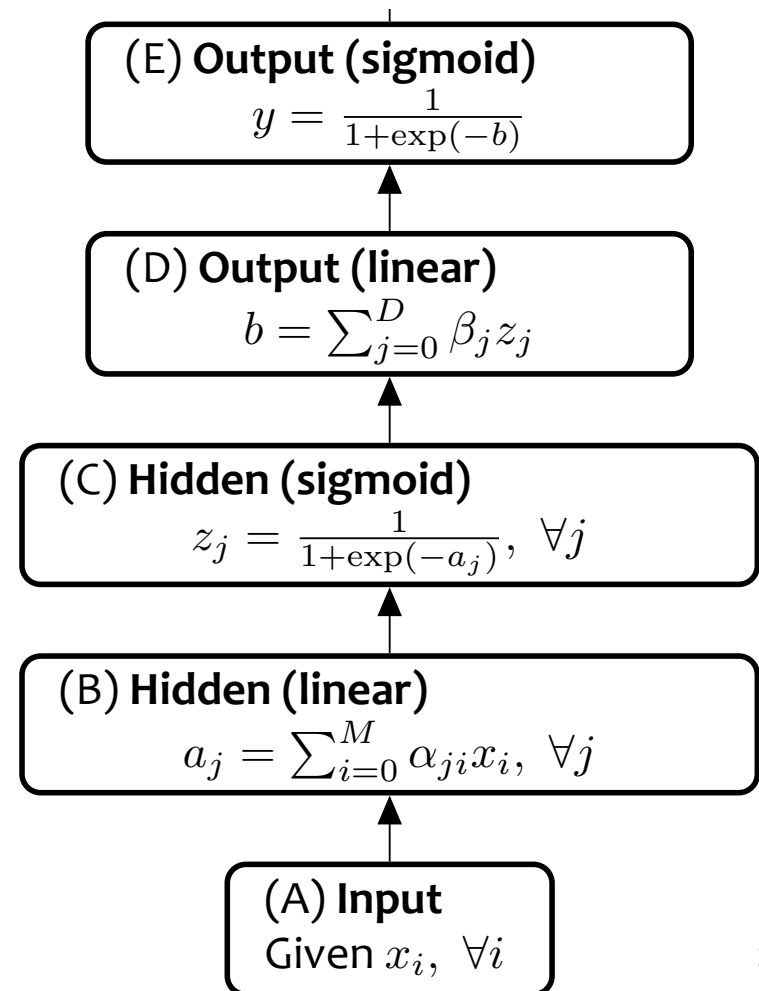
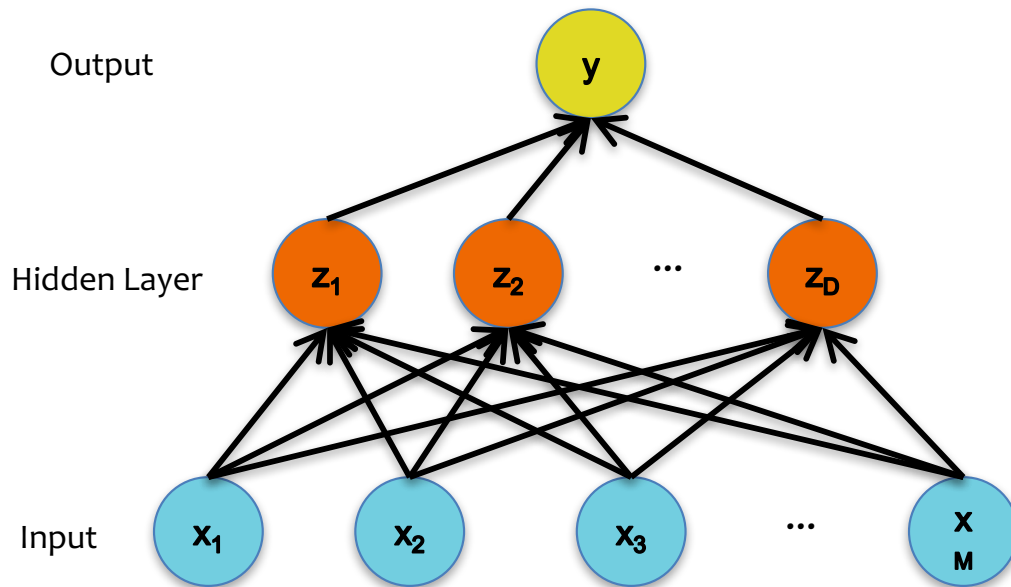
# Decision Functions

# Neural Network



# Decision Functions

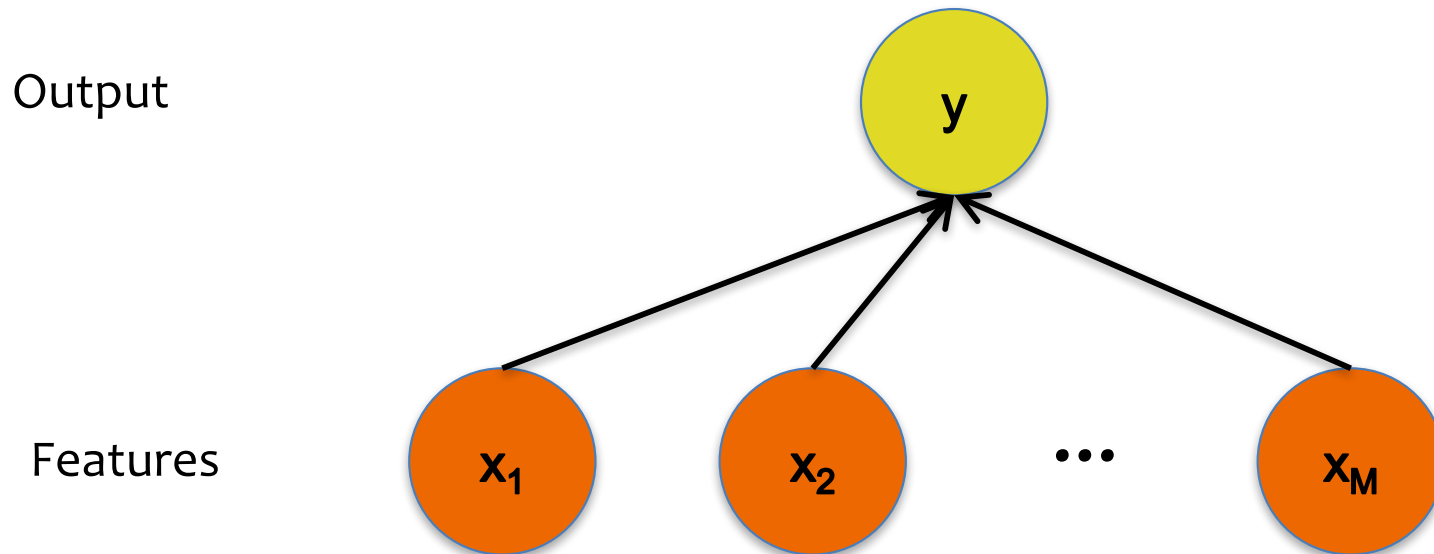
# Neural Network



# **NEURAL NETWORKS: REPRESENTATIONAL POWER**

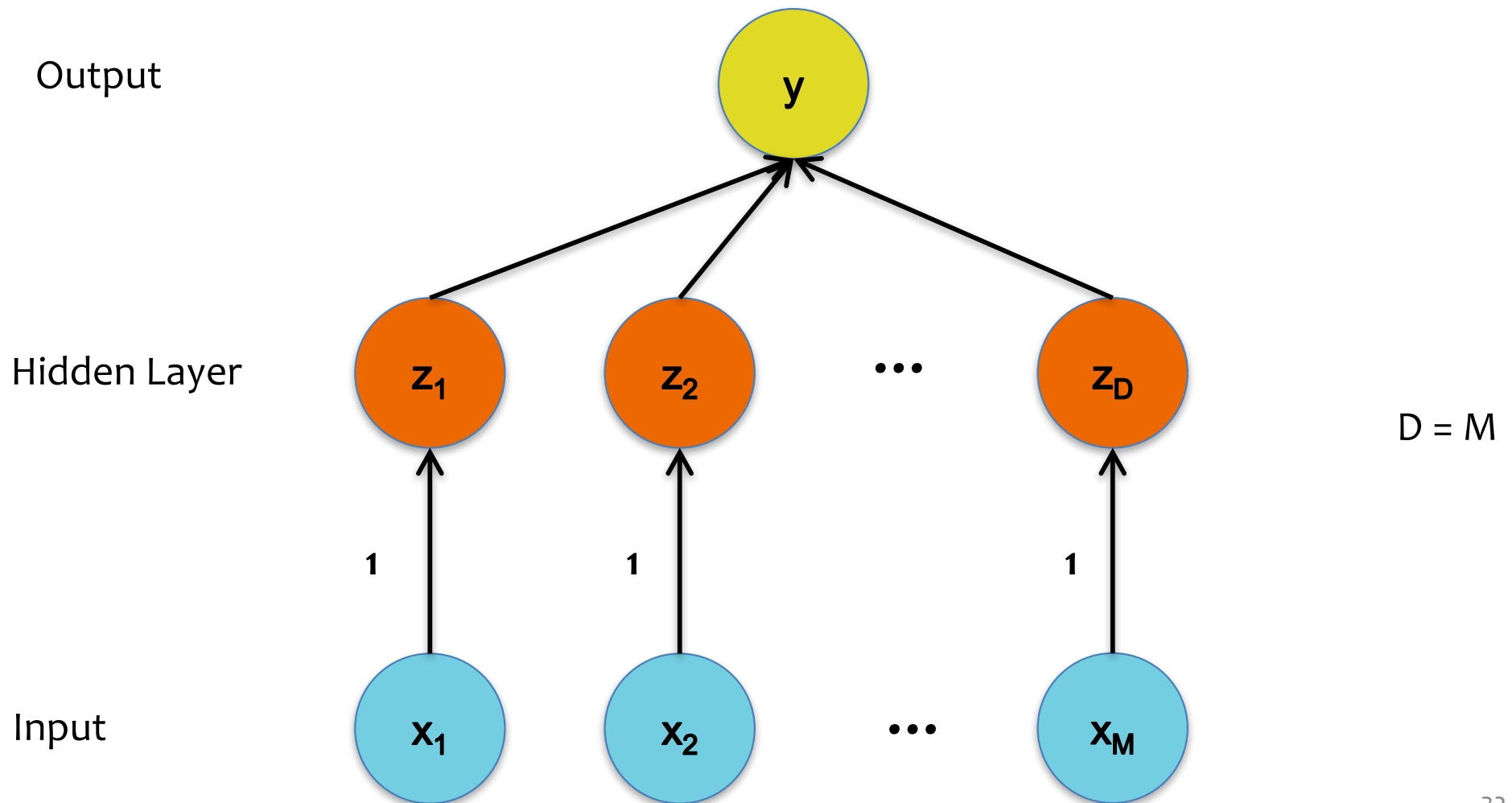
# Building a Neural Net

*Q: How many hidden units,  $D$ , should we use?*



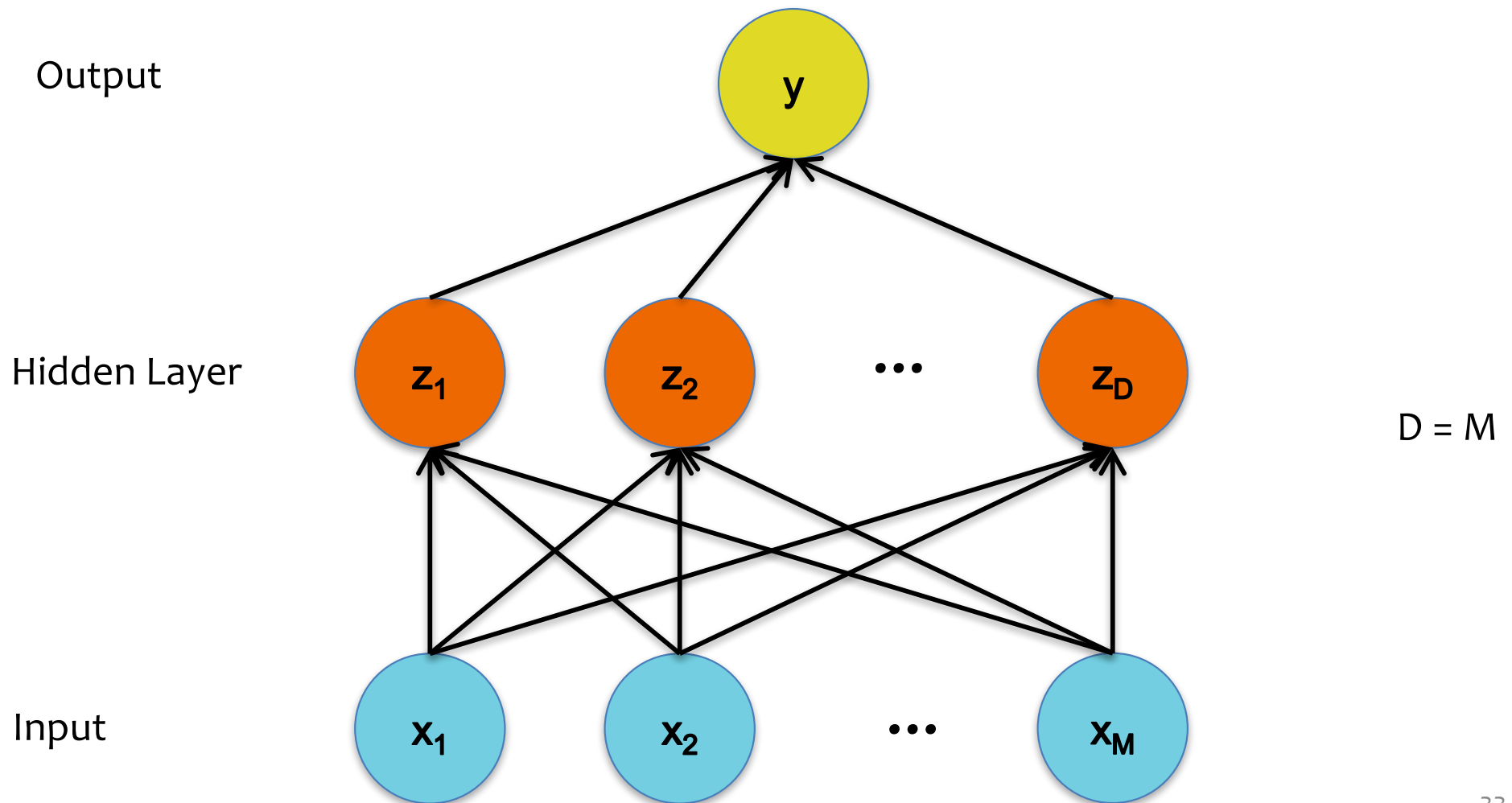
# Building a Neural Net

*Q: How many hidden units,  $D$ , should we use?*



# Building a Neural Net

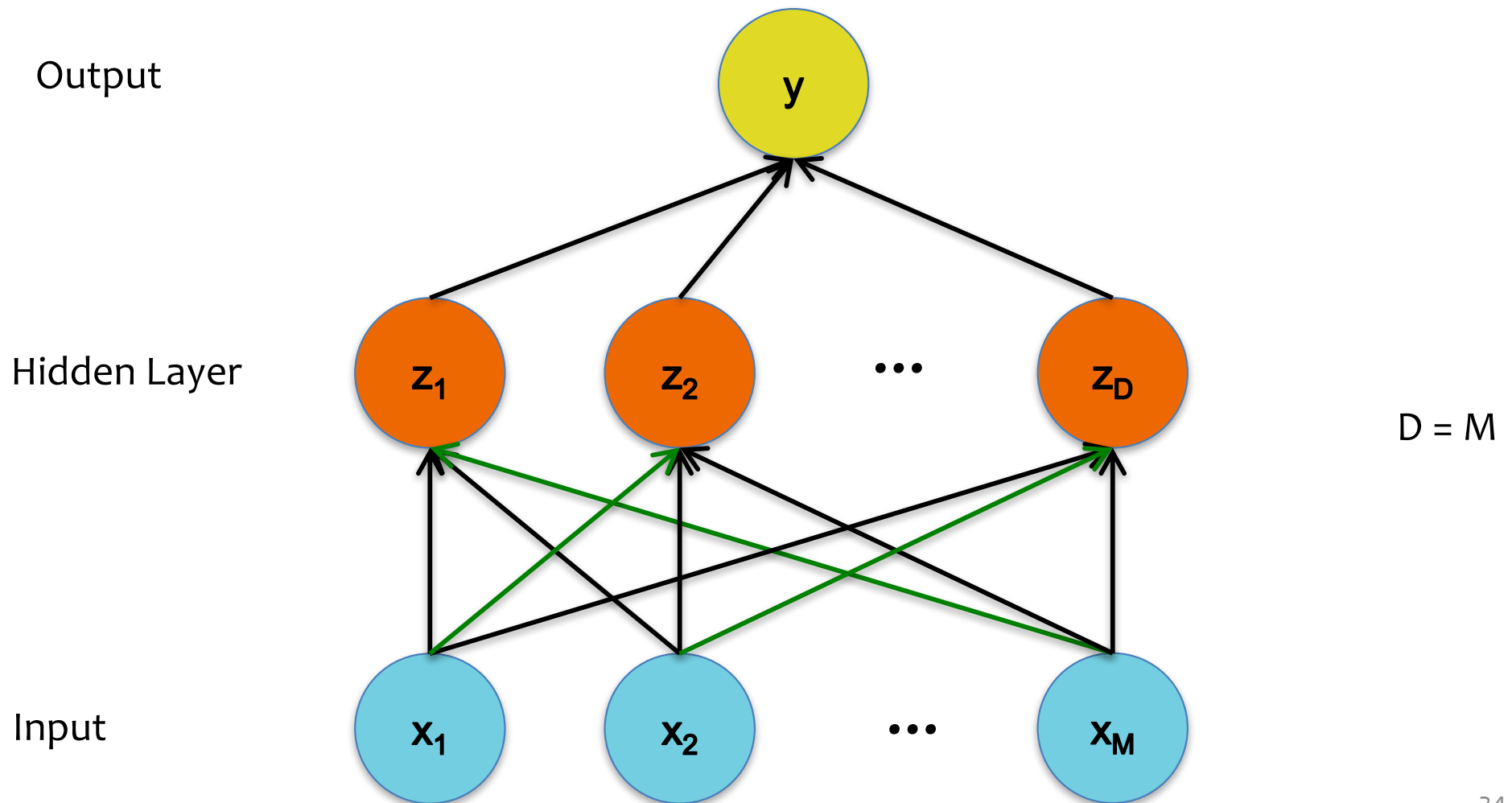
*Q: How many hidden units,  $D$ , should we use?*





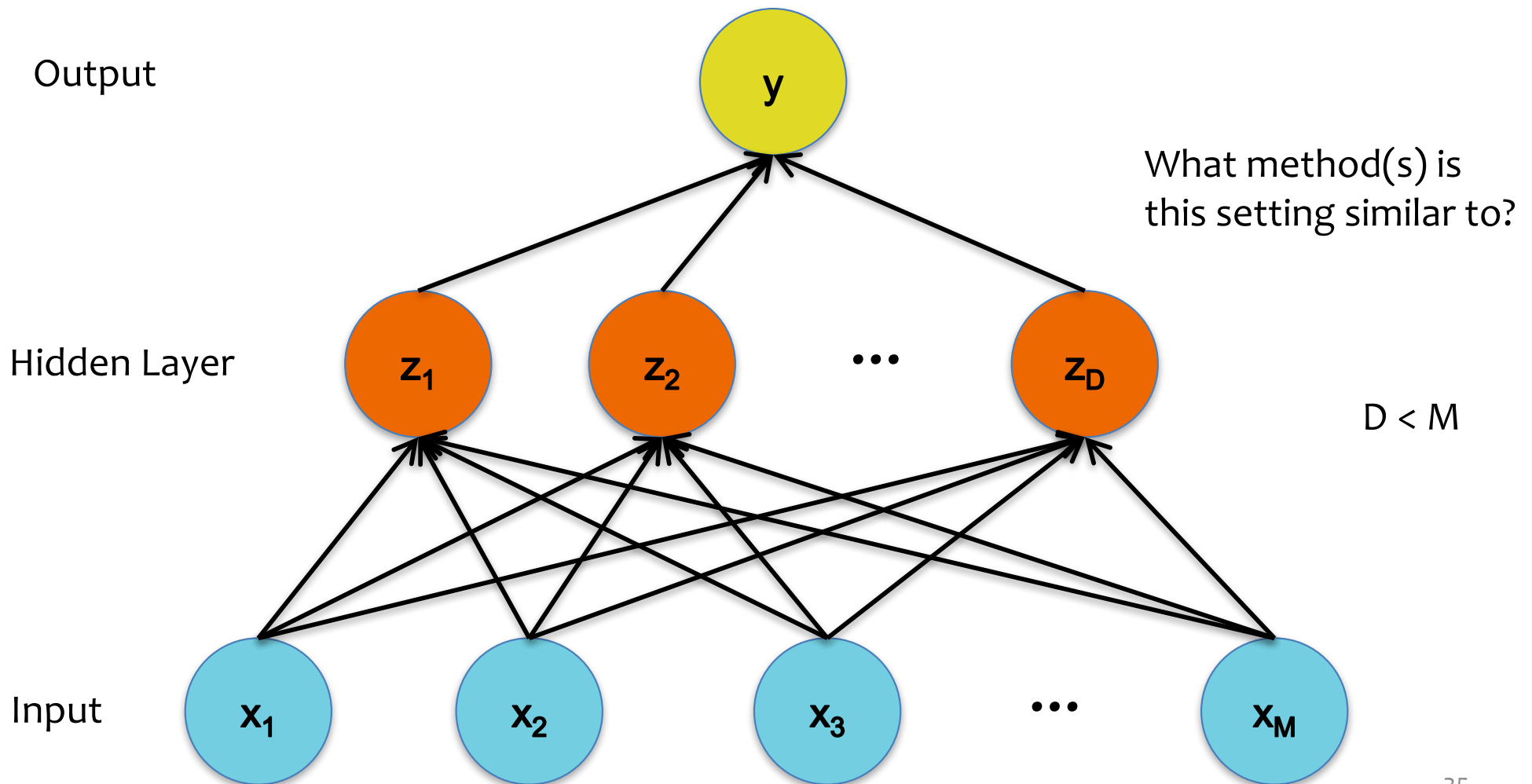
# Building a Neural Net

*Q: How many hidden units,  $D$ , should we use?*



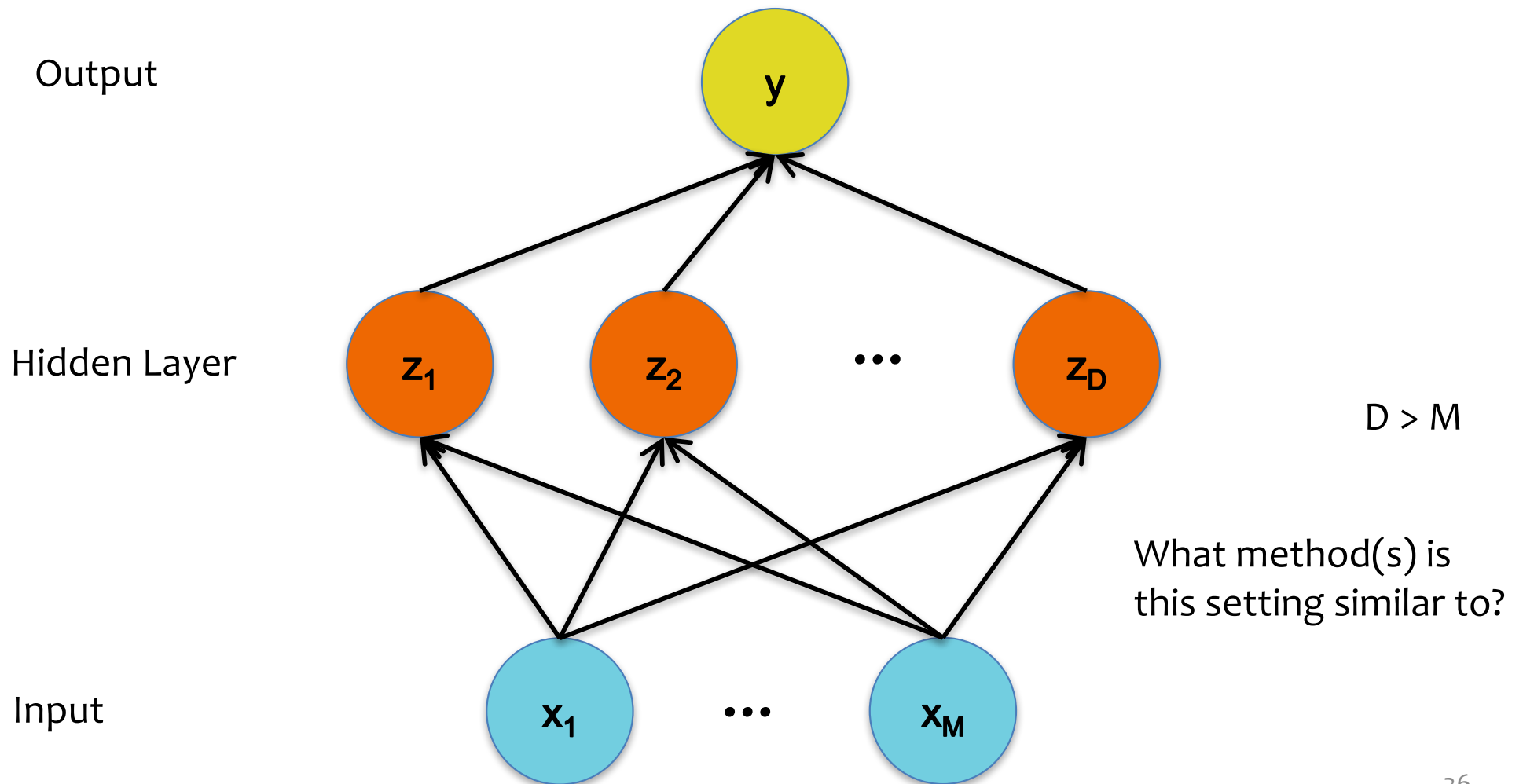
# Building a Neural Net

*Q: How many hidden units,  $D$ , should we use?*

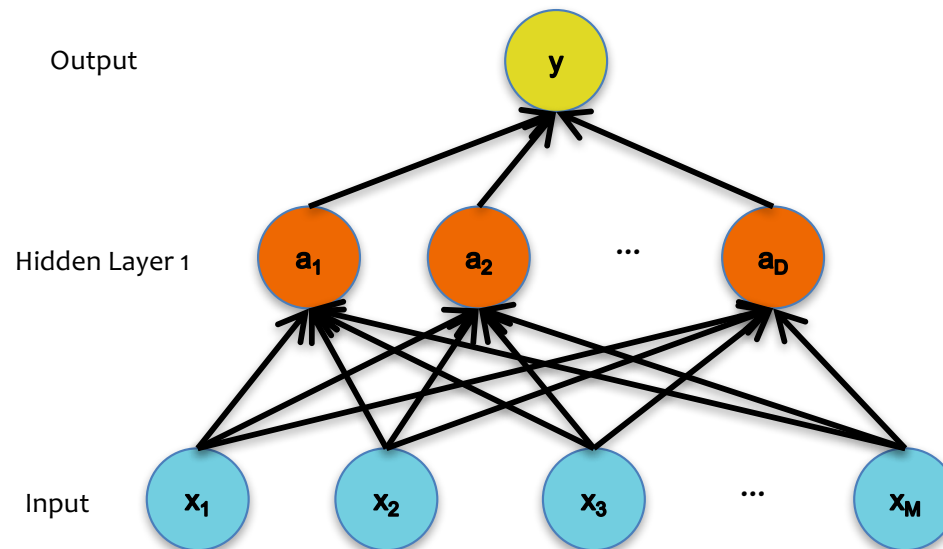


# Building a Neural Net

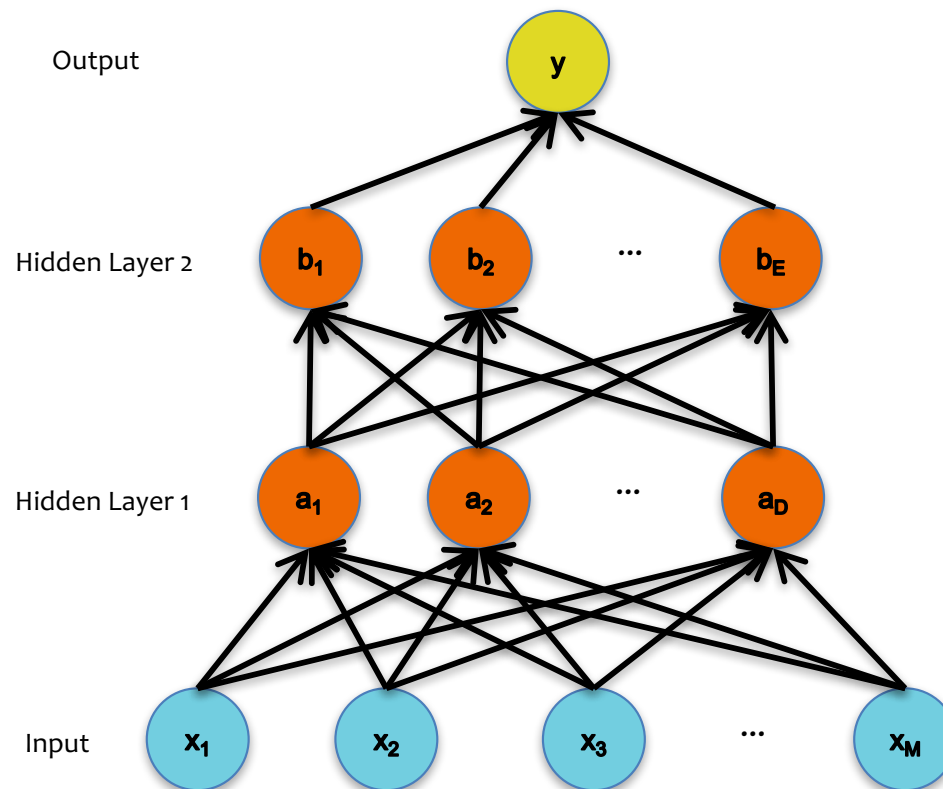
*Q: How many hidden units,  $D$ , should we use?*



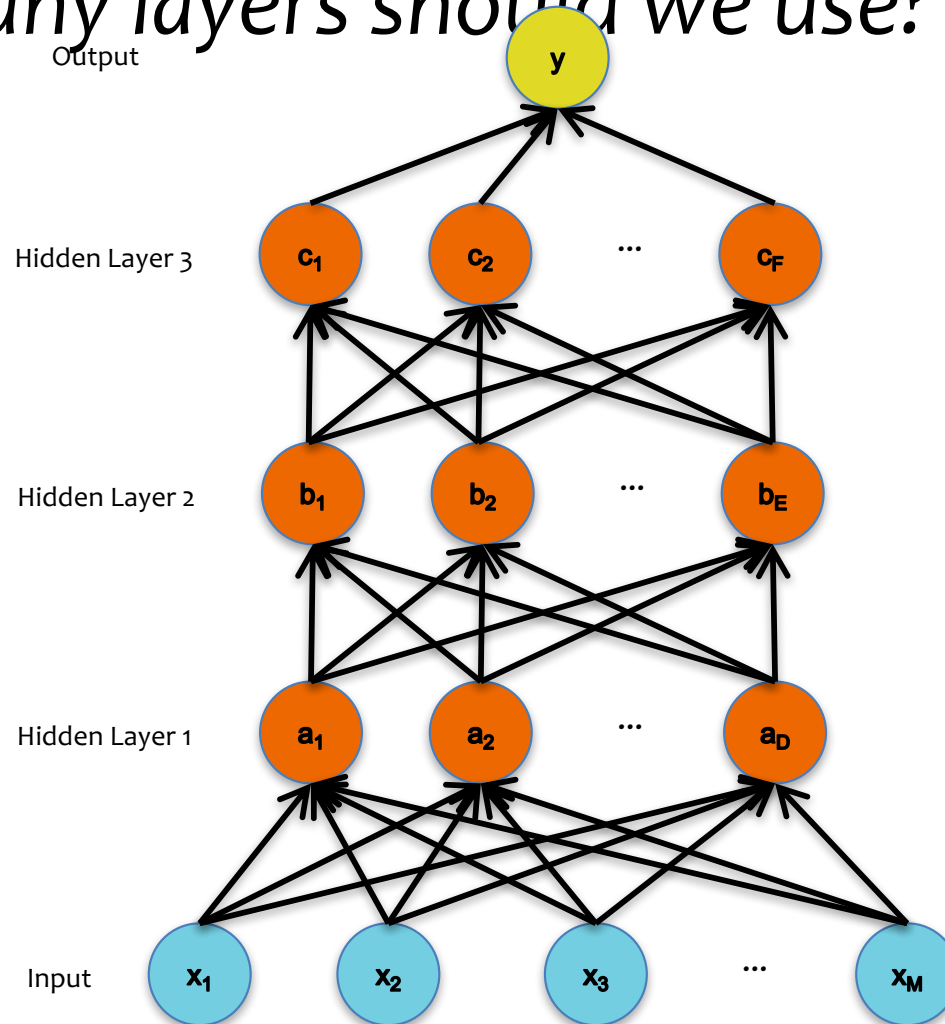
*Q: How many layers should we use?*



*Q: How many layers should we use?*



Q: *How many layers should we use?*



*Q: How many layers should we use?*

- **Theoretical answer:**

- A neural network with 1 hidden layer is a **universal function approximator**
- Cybenko (1989): For any continuous function  $g(\mathbf{x})$ , there exists a 1-hidden-layer neural net  $h_{\theta}(\mathbf{x})$  s.t.  $|h_{\theta}(\mathbf{x}) - g(\mathbf{x})| < \epsilon$  for all  $\mathbf{x}$ , assuming sigmoid activation functions

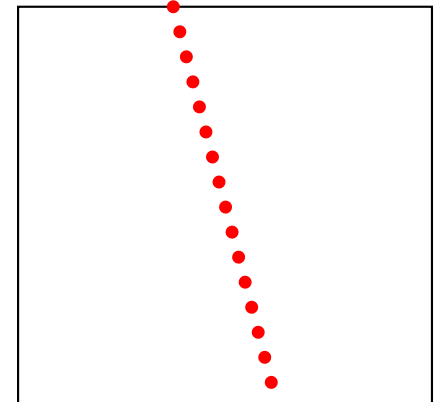
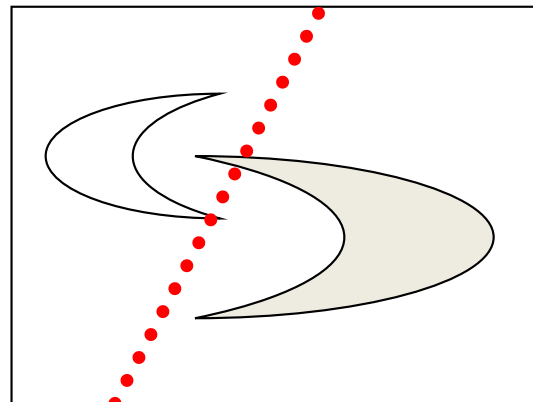
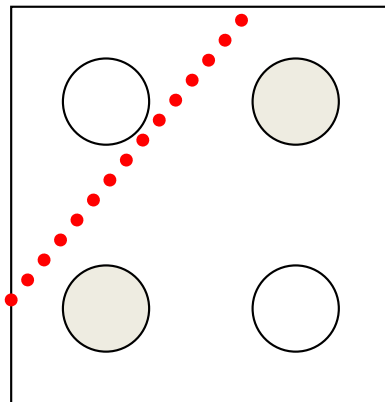
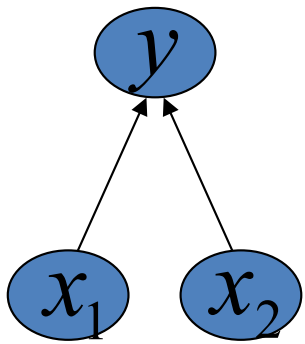
- **Empirical answer:**

- Before 2006: “Deep networks (e.g. 3 or more hidden layers) are too hard to train”
- After 2006: “Deep networks are easier to train than shallow networks (e.g. 2 or fewer layers) for many problems”

Big caveat: You need to know and use the right tricks.

# Decision Boundary

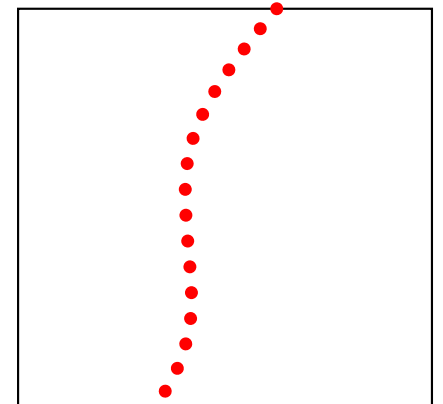
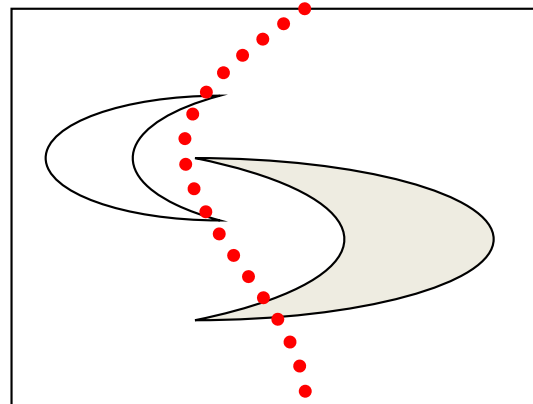
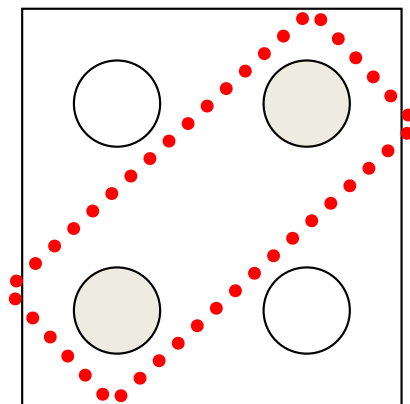
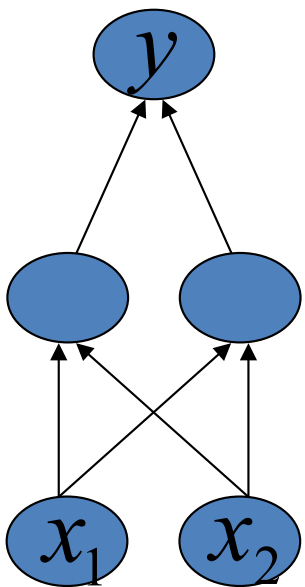
- 0 hidden layers: linear classifier
  - Hyperplanes



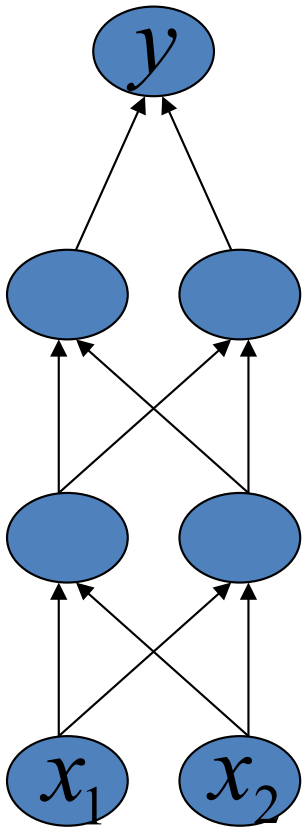


# Decision Boundary

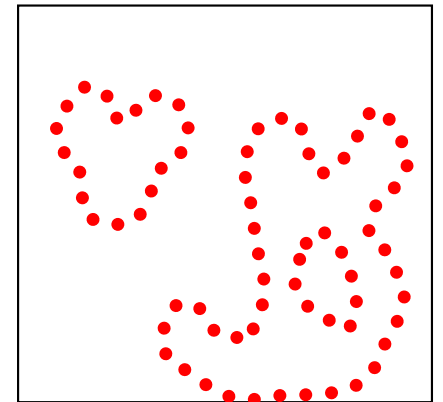
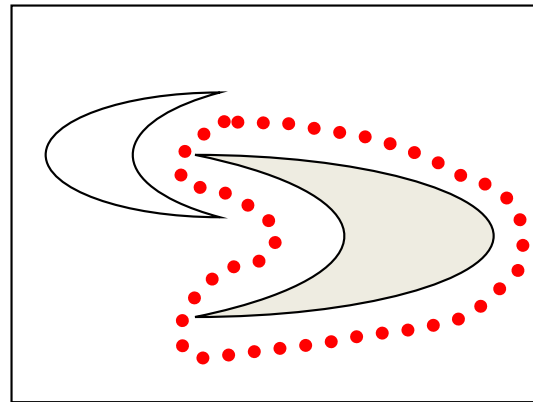
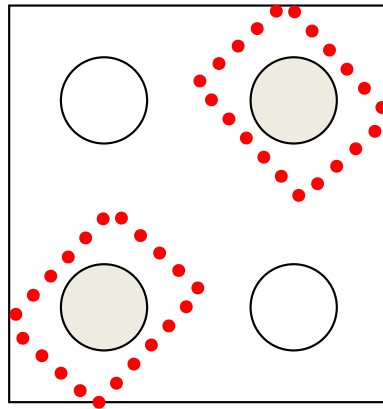
- 1 hidden layer
  - Boundary of convex region (open or closed)



# Decision Boundary



- 2 hidden layers
  - Combinations of convex regions



- We don't know the “right” levels of abstraction
- So let the model figure it out!

Feature representation



3rd layer  
“Objects”



2nd layer  
“Object parts”



1st layer  
“Edges”



Pixels

## Face Recognition:

- Deep Network can build up increasingly higher levels of abstraction
- Lines, parts, regions

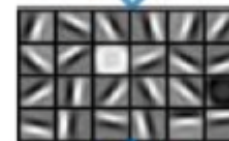
Feature representation



3rd layer  
"Objects"



2nd layer  
"Object parts"



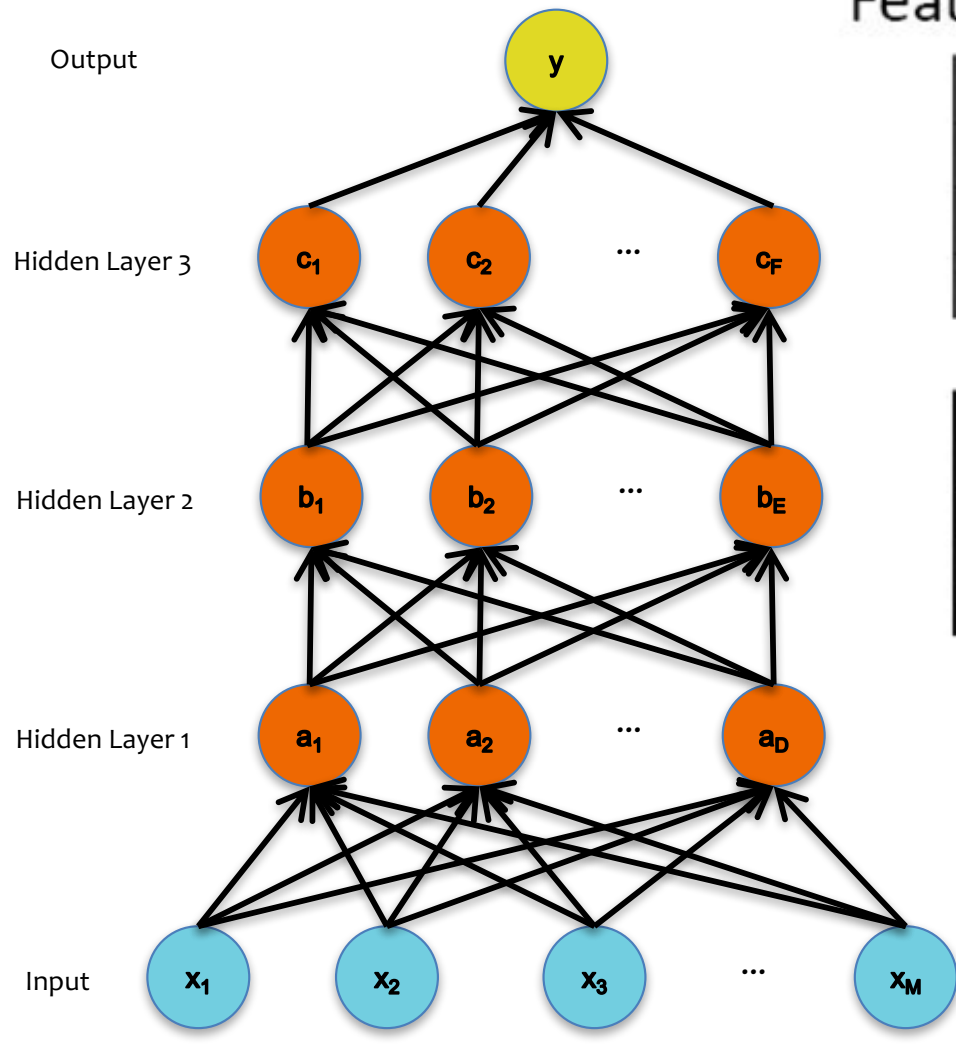
1st layer  
"Edges"



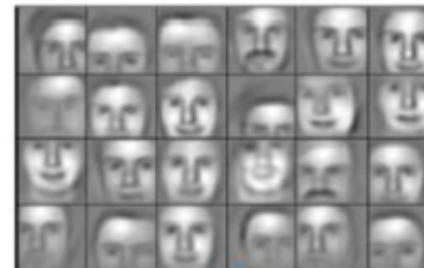
Pixels

# Decision Functions

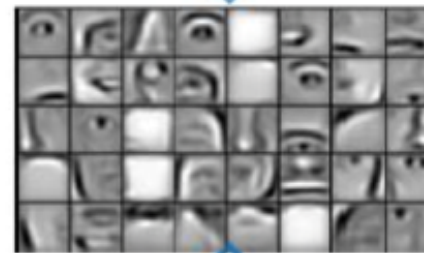
# Different Levels of Abstraction



Feature representation



3rd layer  
"Objects"



2nd layer  
"Object parts"



1st layer  
"Edges"



Pixels

Example from Honglak Lee (NIPS 2010)