



10-601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Neural Networks

Matt Gormley
Lecture 11
Oct. 8, 2018

Reminders

- **Homework 4: Logistic Regression**
 - Out: Sun, Sep 30
 - Due: Tue, Oct 9 at 11:59pm
- **Homework 5: Neural Networks**
 - Out: Tue, Oct 9
 - Due: Sat, Oct 20 at 11:59pm

Q&A

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)

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:

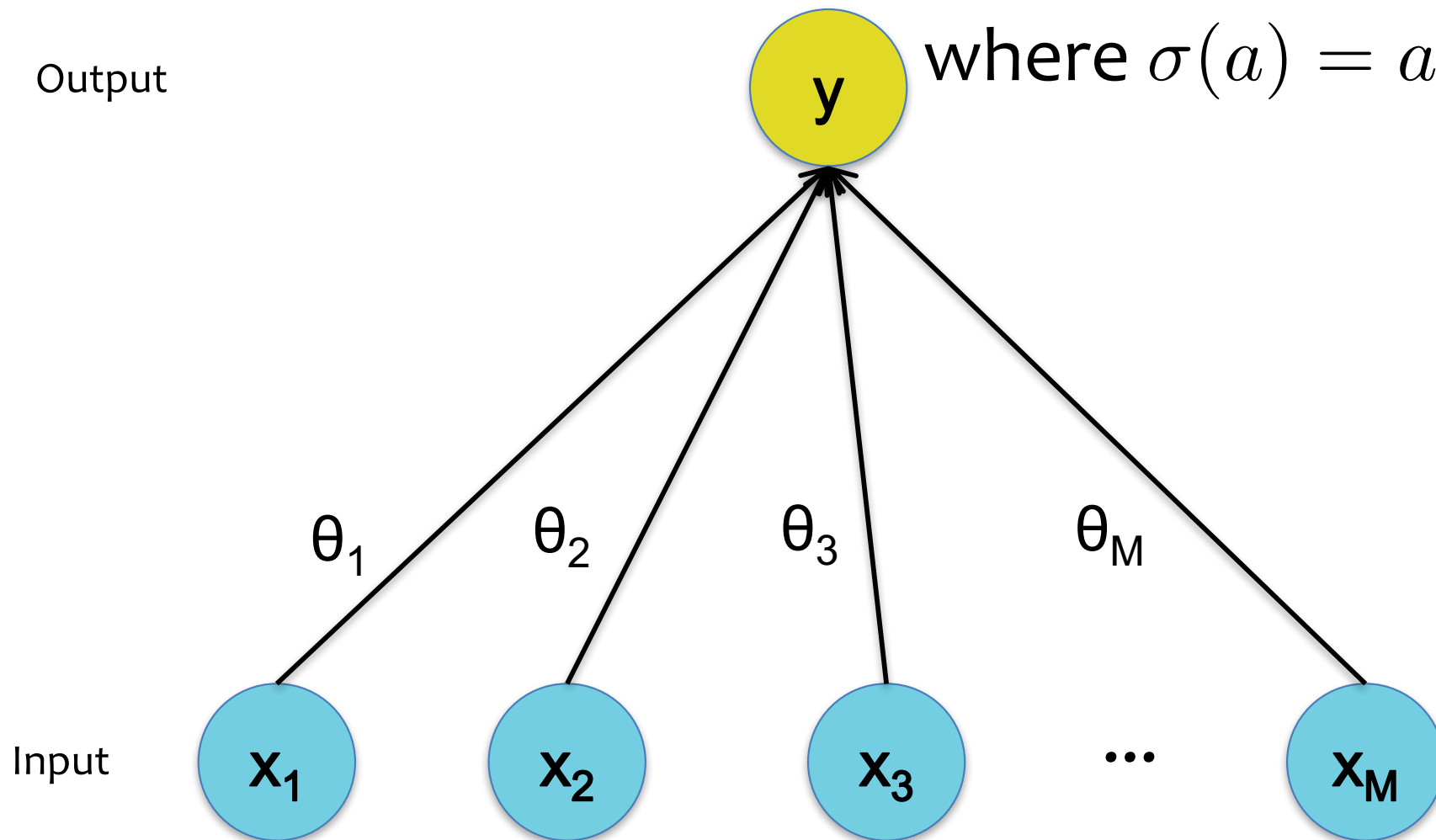
– Take small steps opposite the gradient)

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

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

where $\sigma(a) = a$

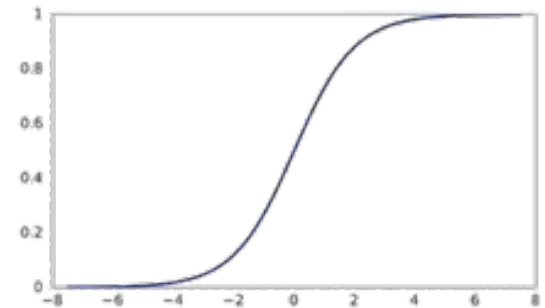
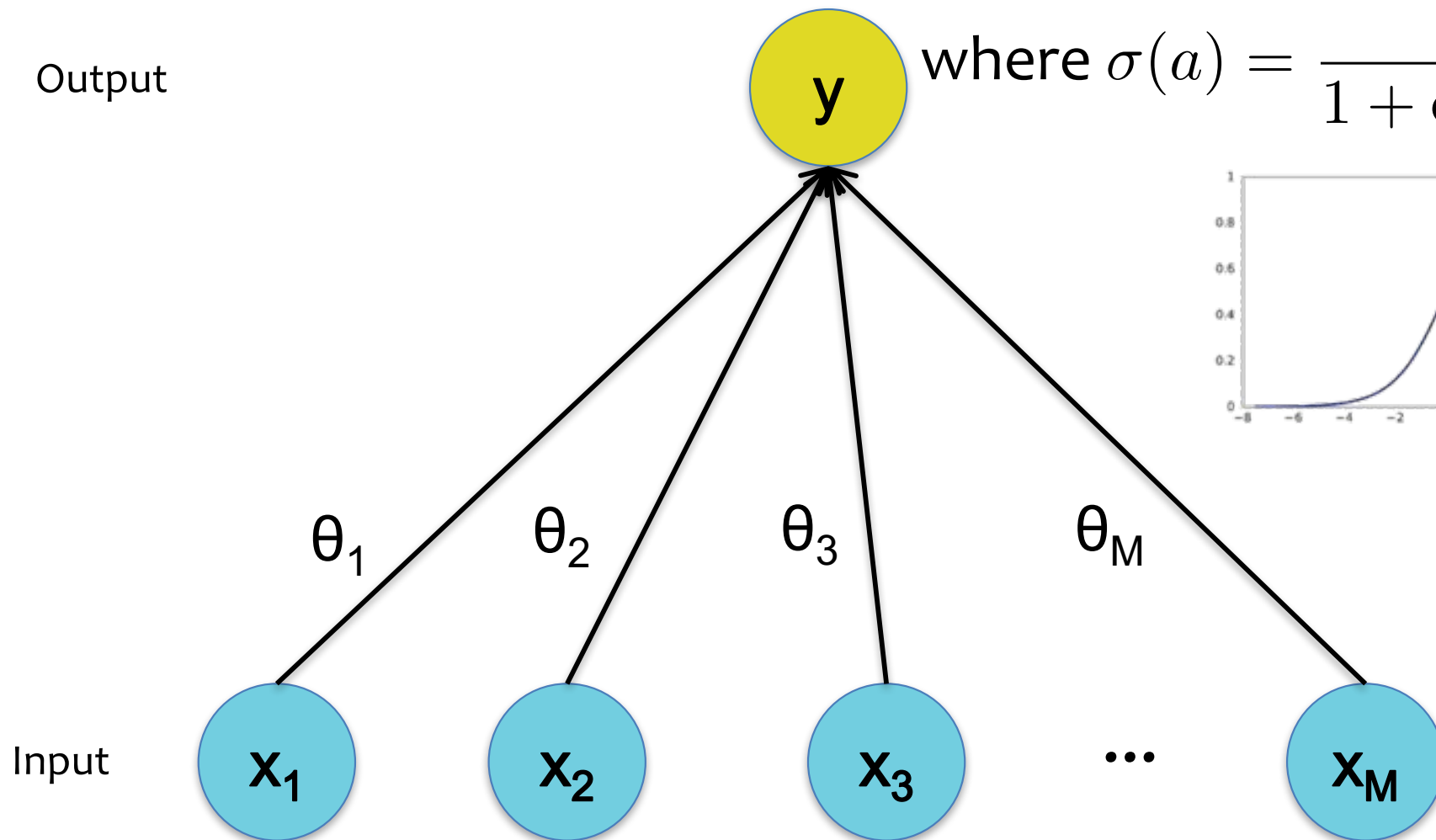
Output



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

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

Output



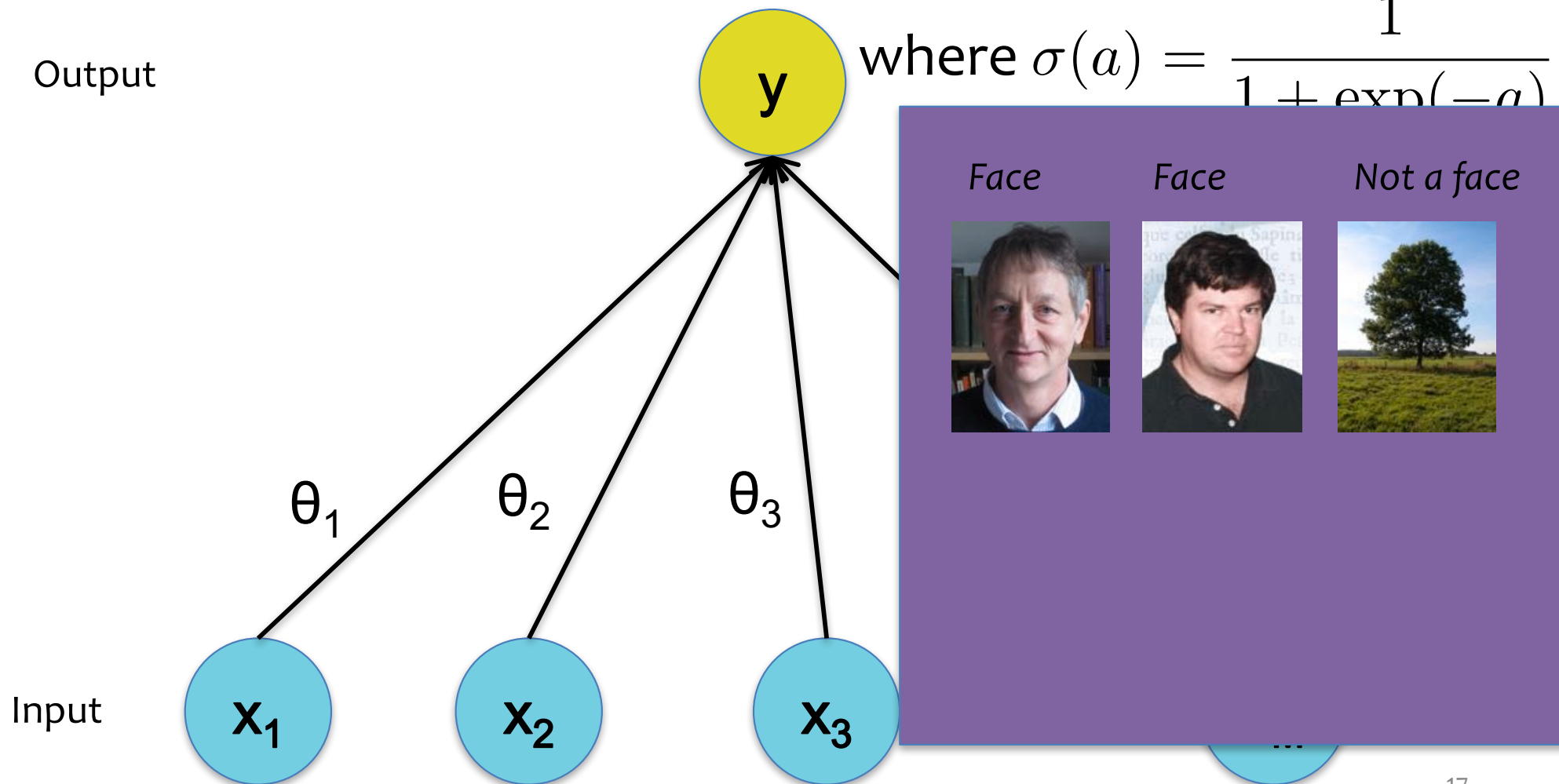
Decision Functions

Logistic Regression

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

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

Output

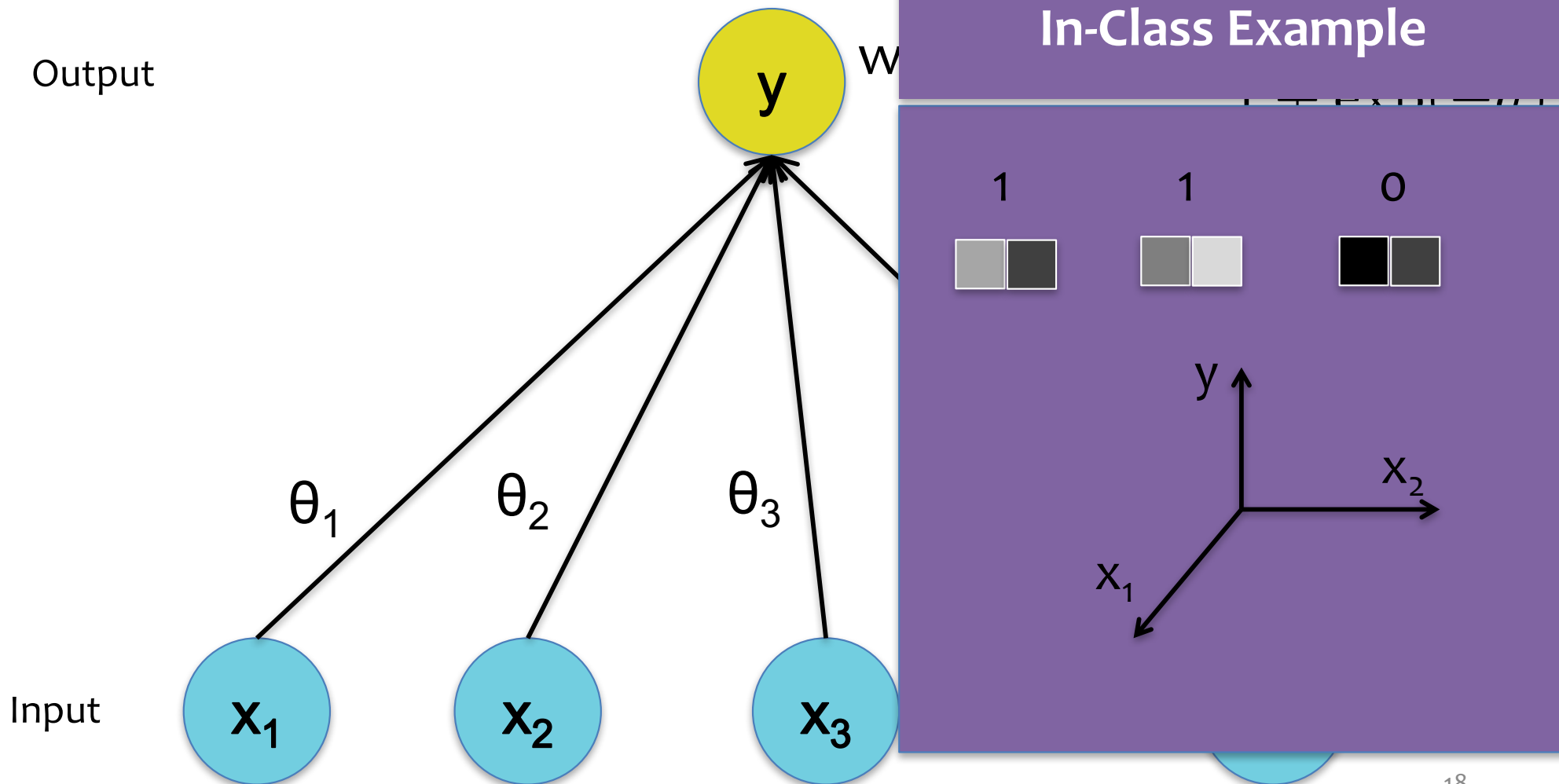


Input

Decision Functions

Logistic Regression

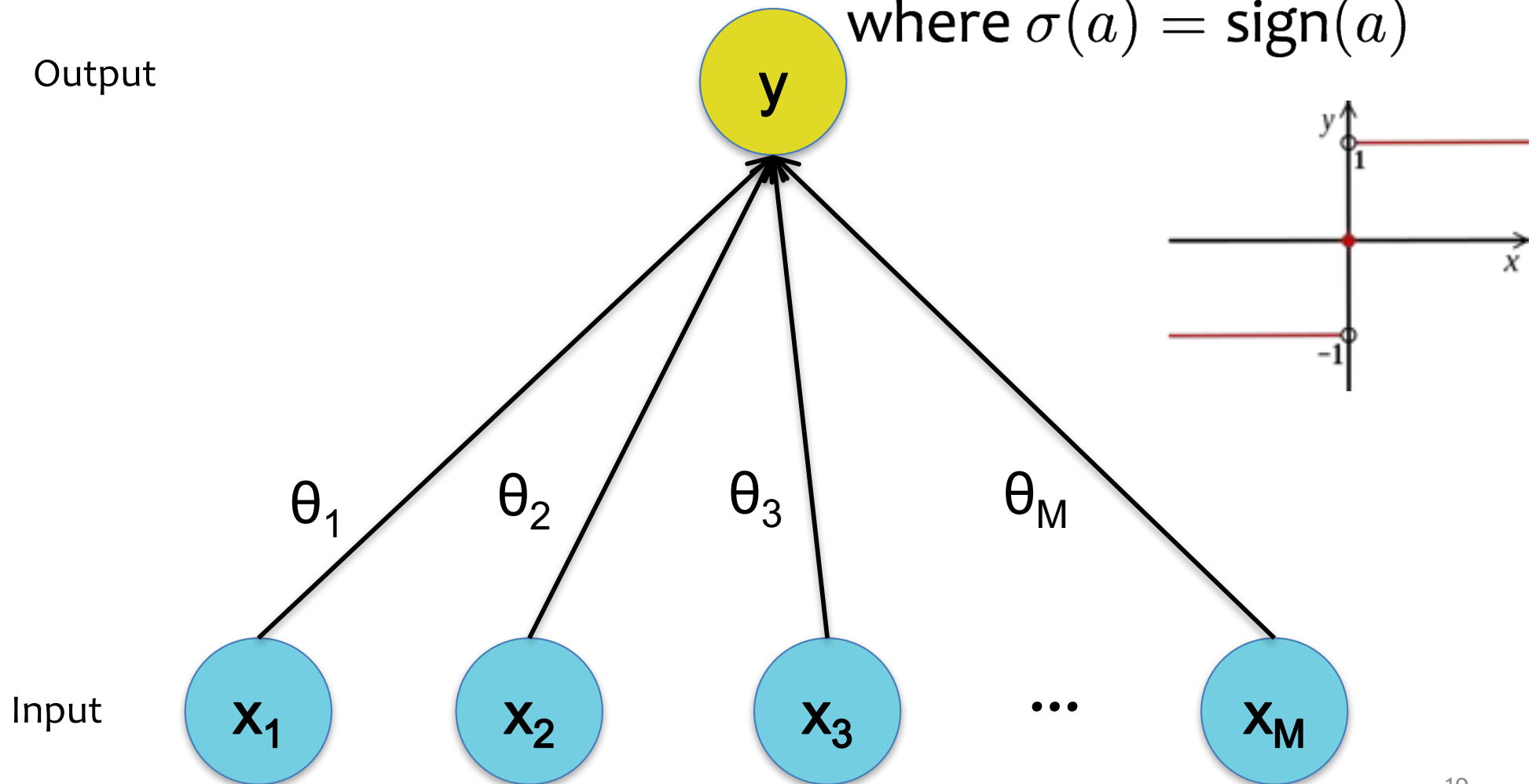
$$y = h_{\theta}(x) = \sigma(\theta^T x)$$



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

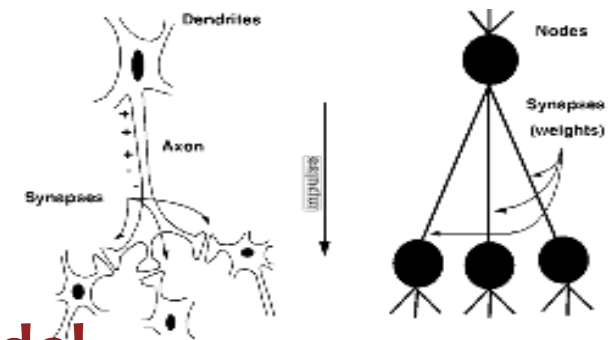
$$\text{where } \sigma(a) = \text{sign}(a)$$

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

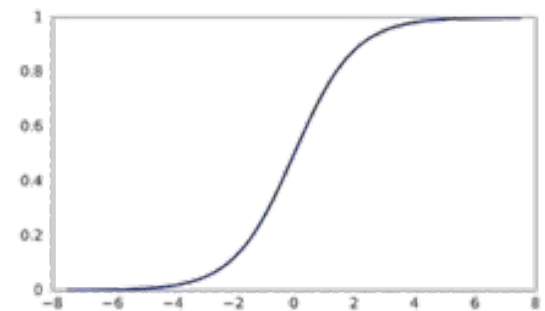
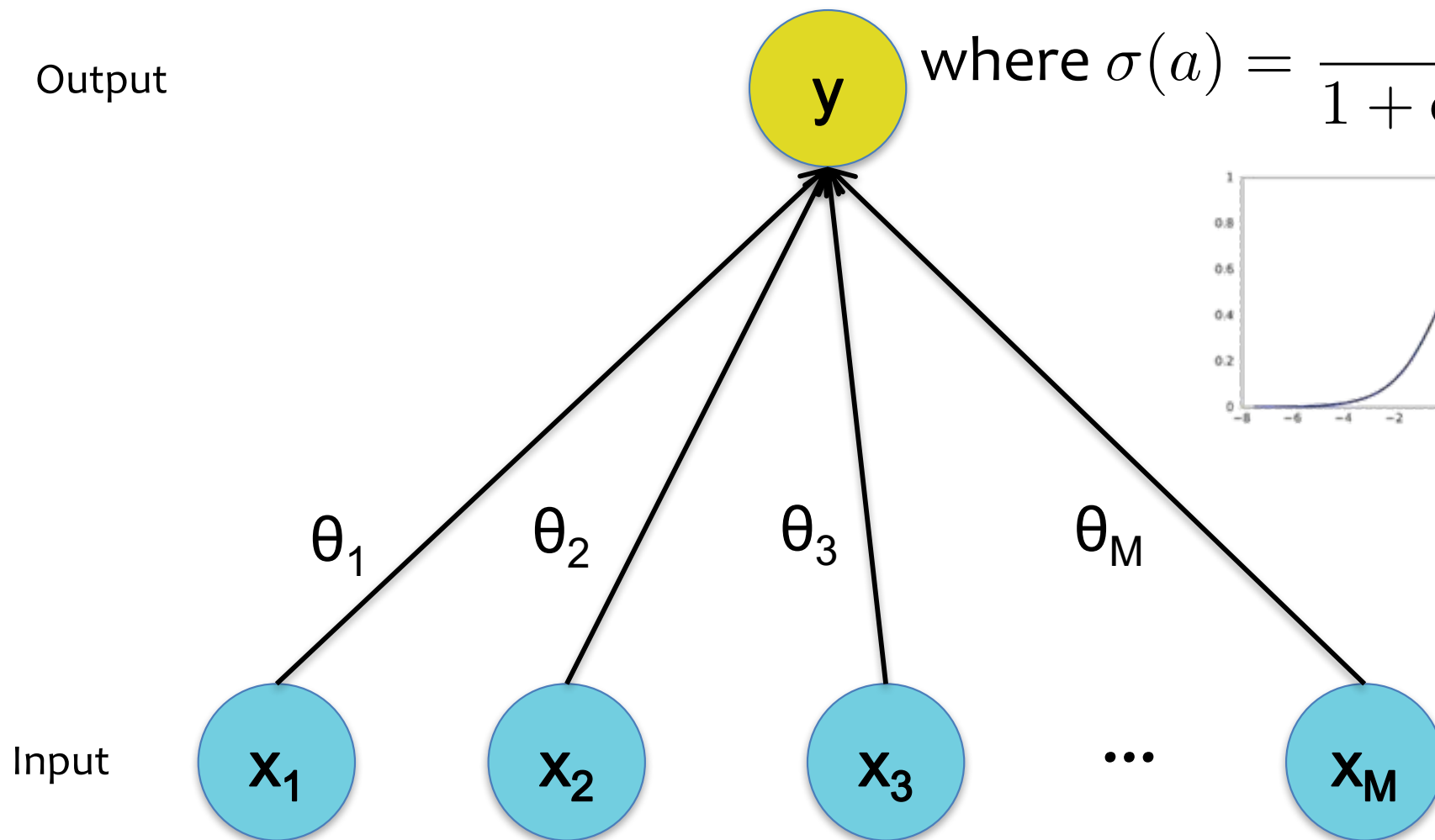
Artificial Computation

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

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

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

Output



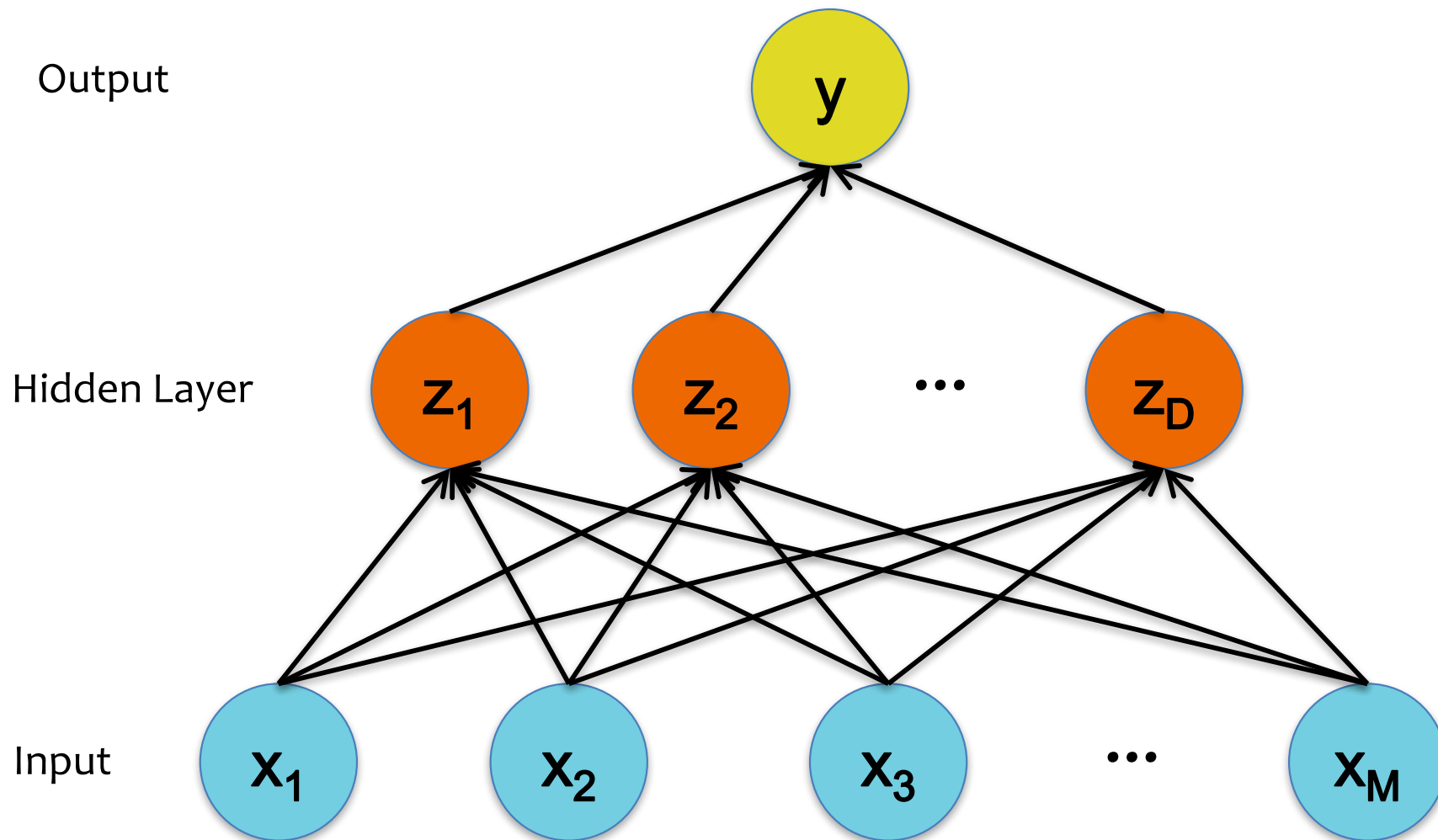
Neural Networks

Chalkboard

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

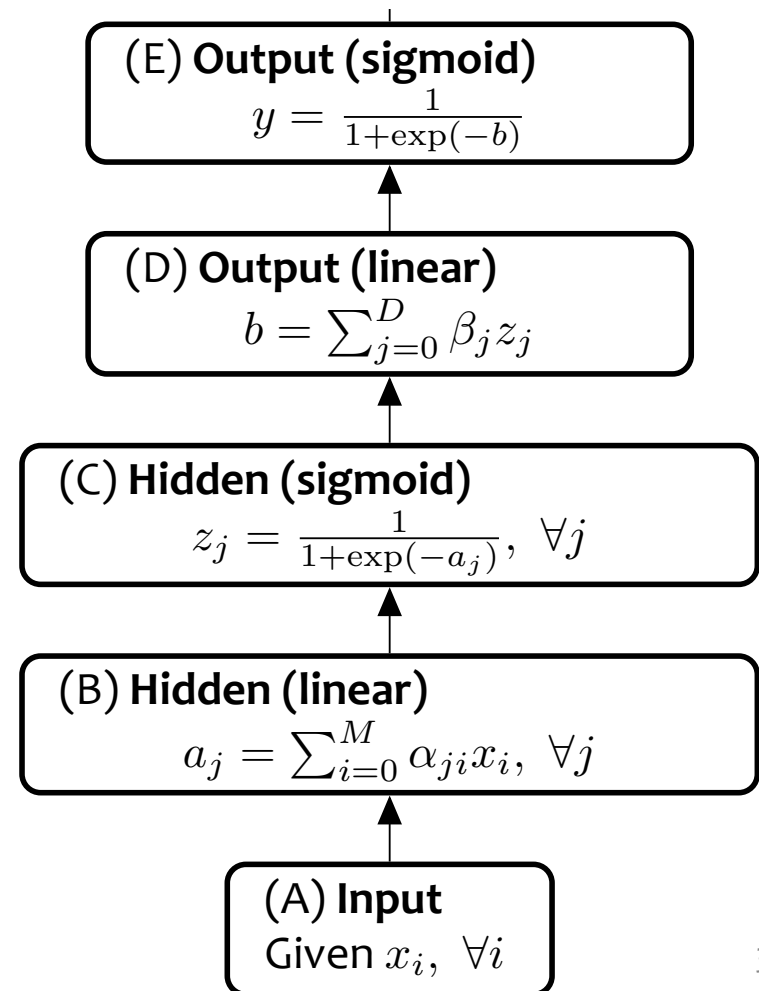
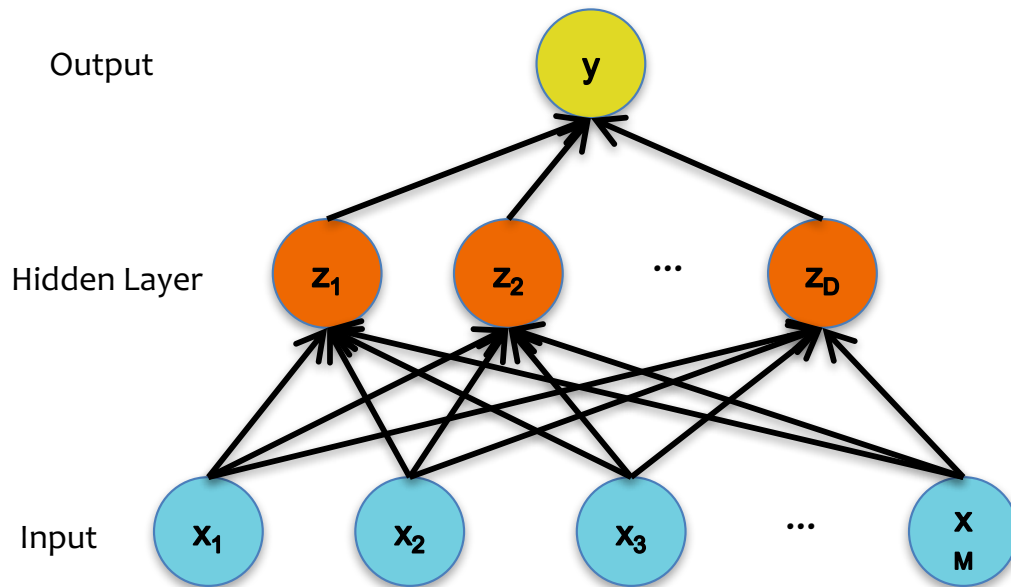
Decision Functions

Neural Network



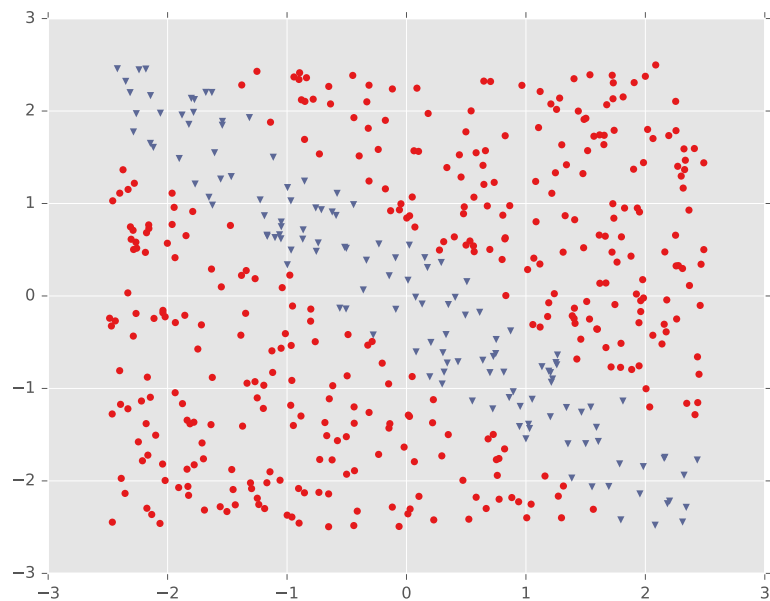
Decision Functions

Neural Network

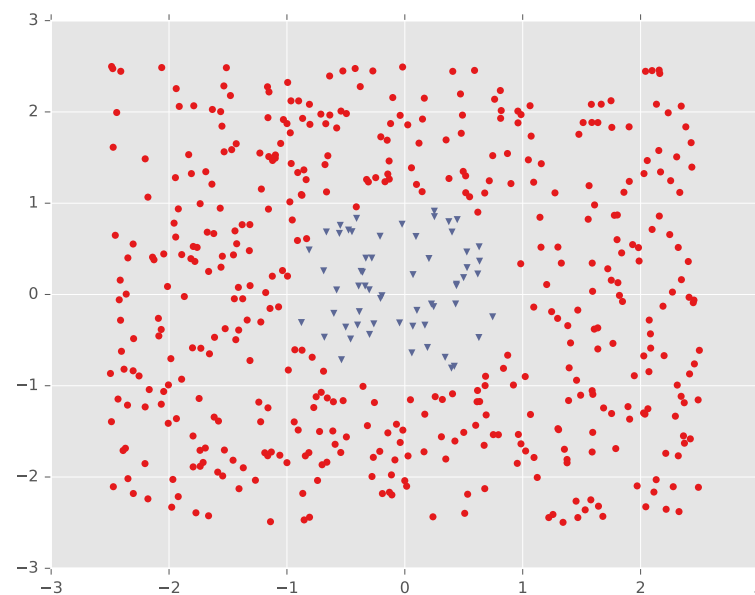


DECISION BOUNDARY EXAMPLES

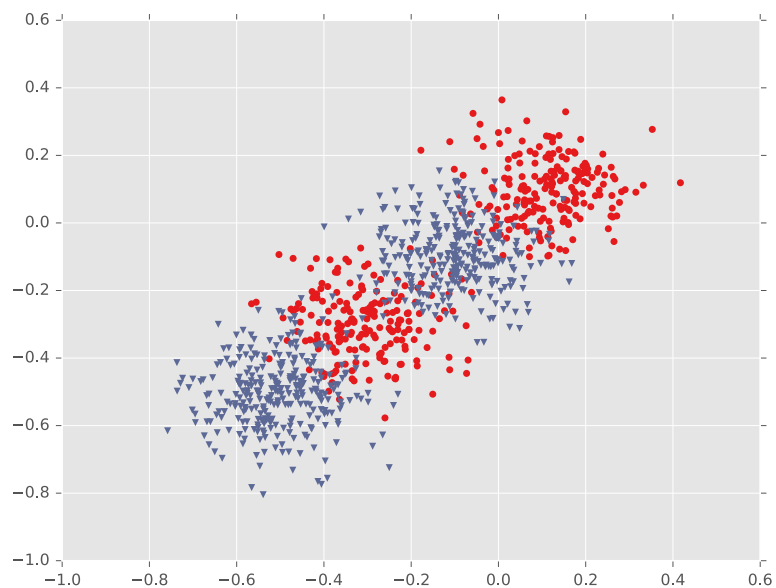
Example #1: Diagonal Band



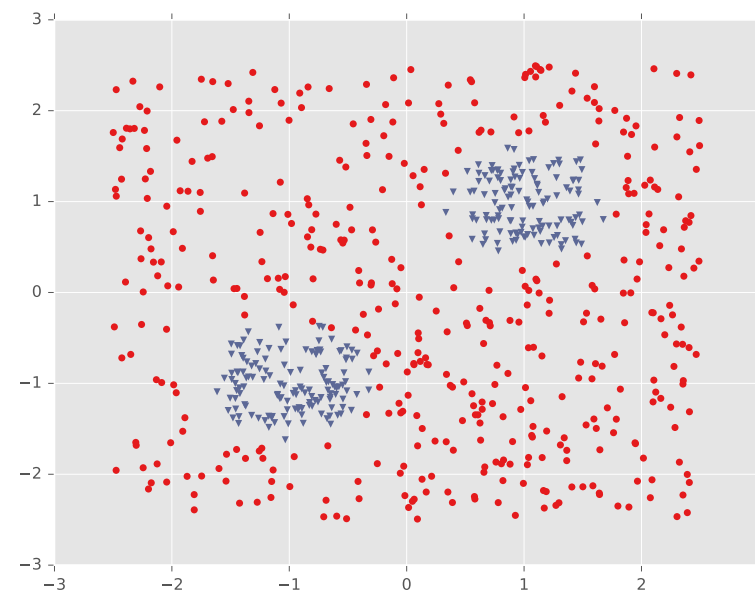
Example #2: One Pocket



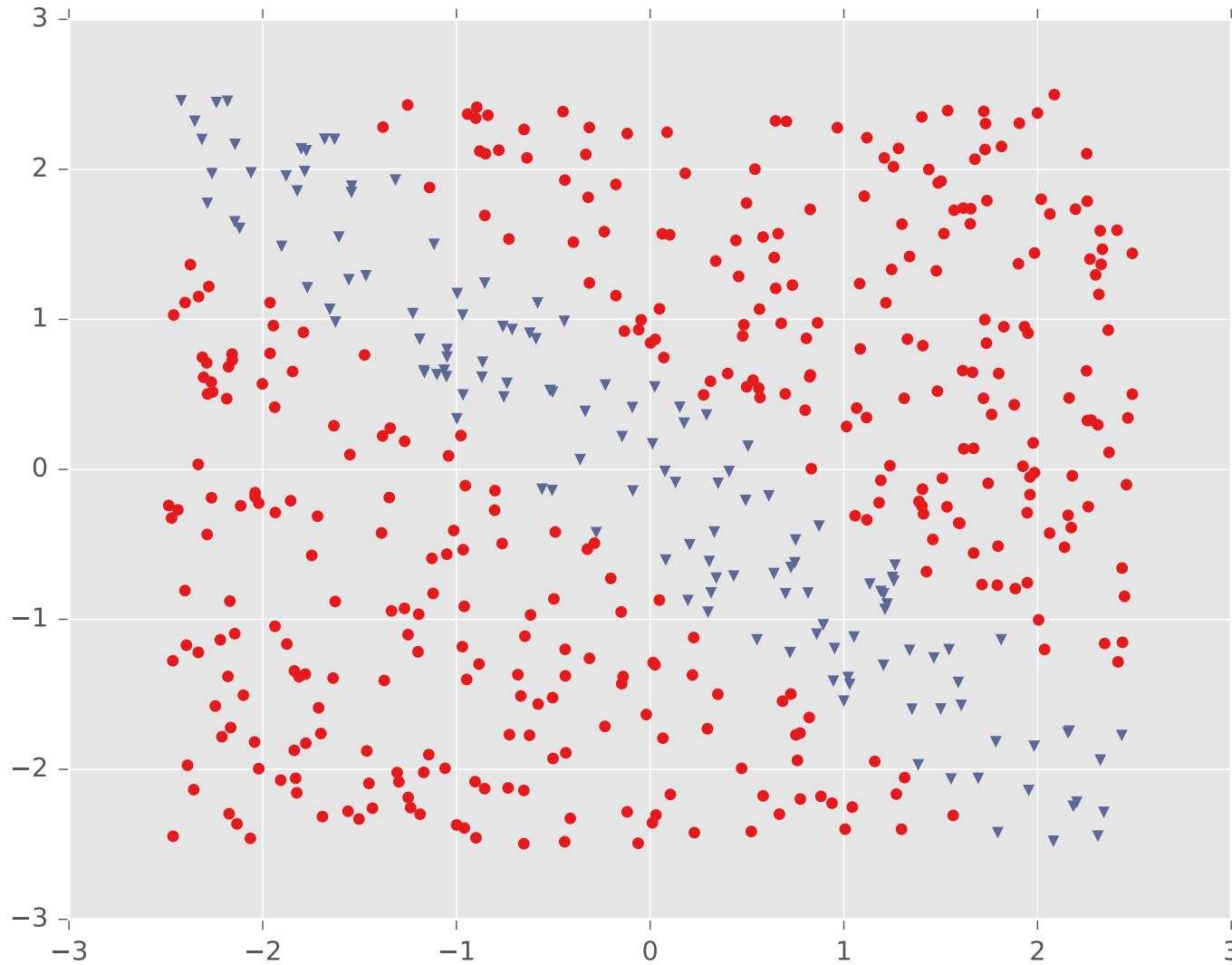
Example #3: Four Gaussians



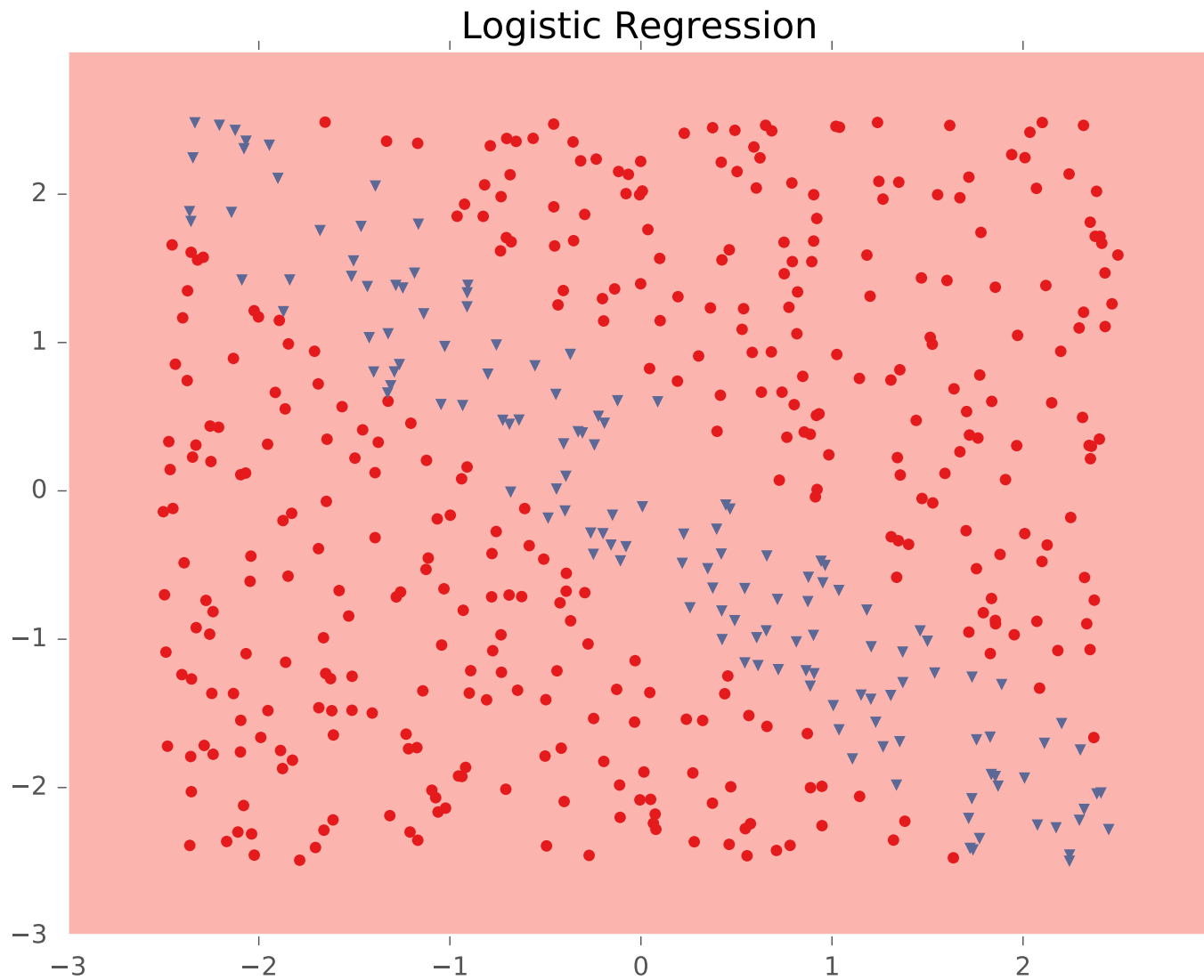
Example #4: Two Pockets



Example #1: Diagonal Band

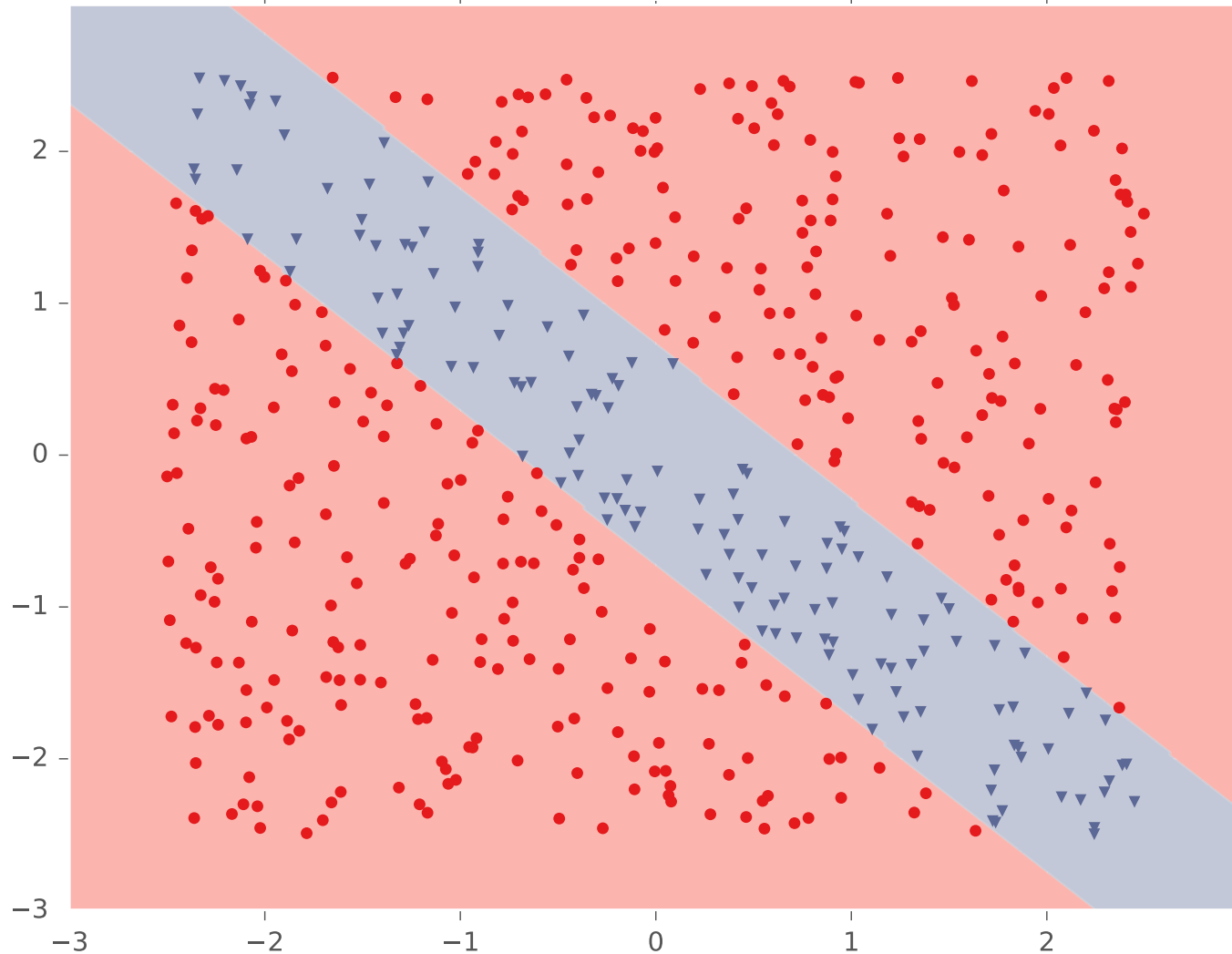


Example #1: Diagonal Band



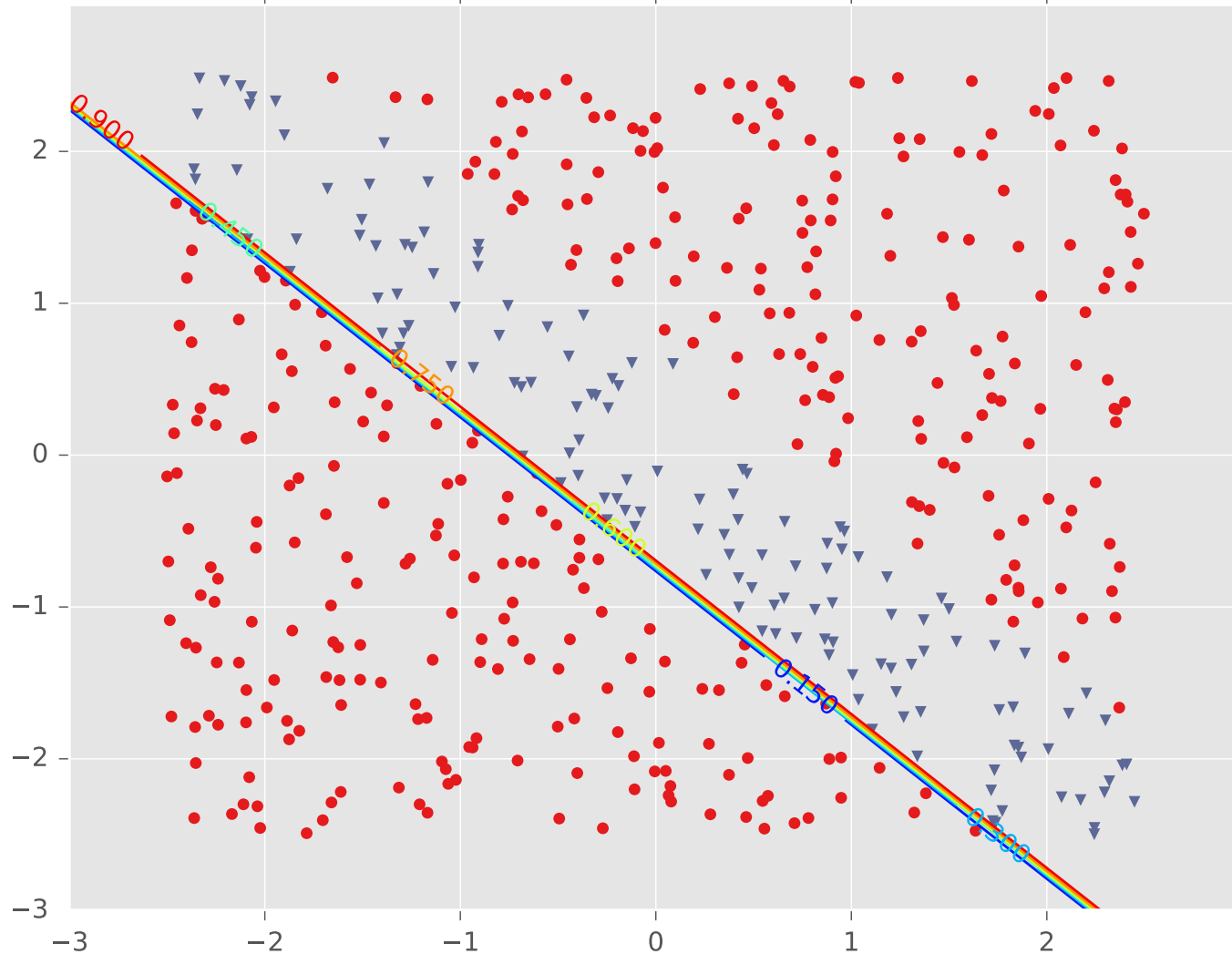
Example #1: Diagonal Band

Tuned Neural Network (hidden=2, activation=logistic)



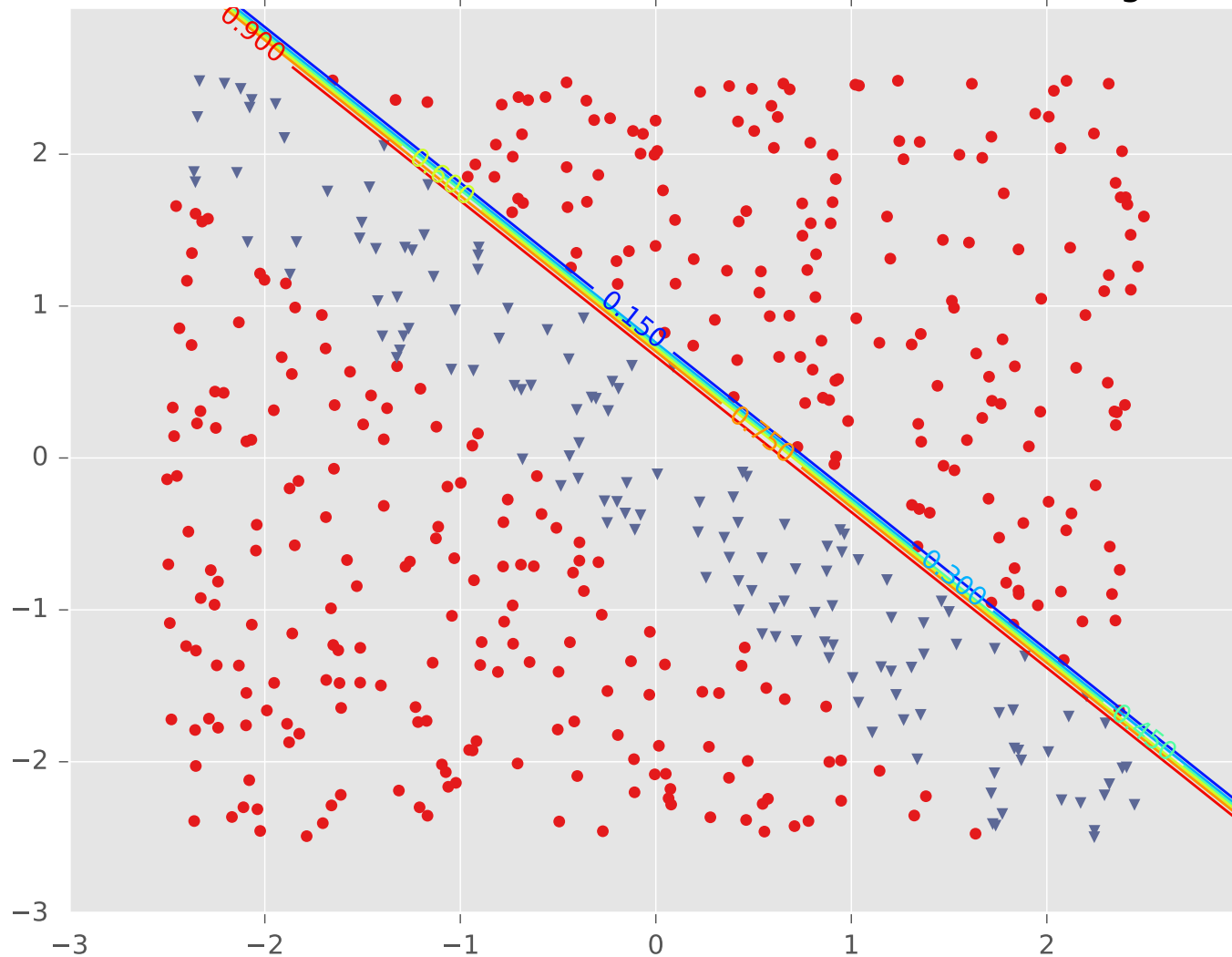
Example #1: Diagonal Band

LR1 for Tuned Neural Network (hidden=2, activation=logistic)

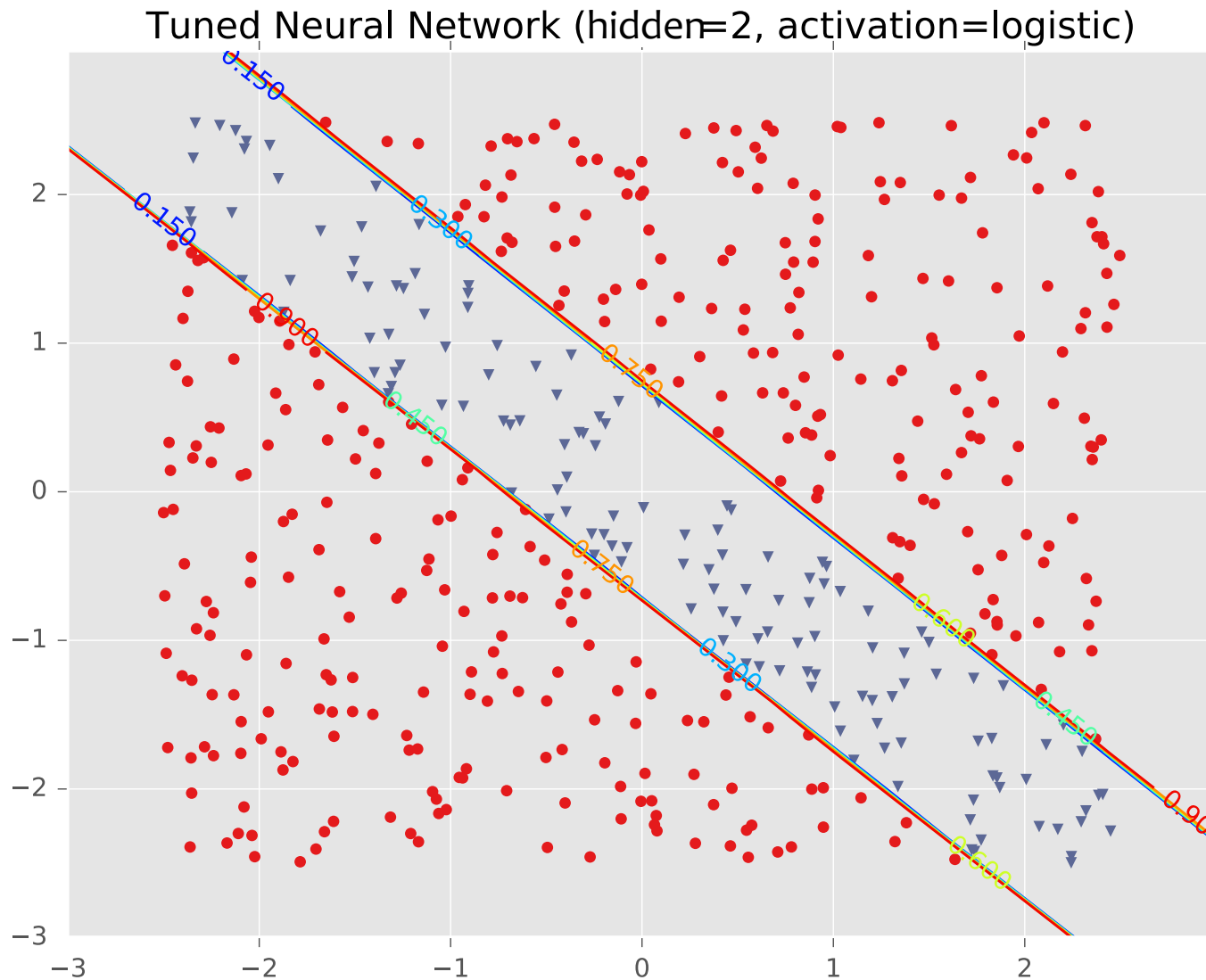


Example #1: Diagonal Band

LR2 for Tuned Neural Network (hidden=2, activation=logistic)

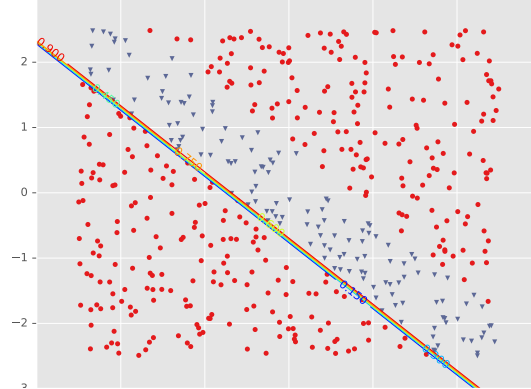


Example #1: Diagonal Band

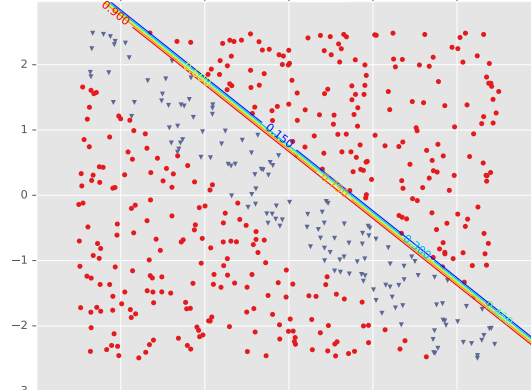


Example #1: Diagonal Band

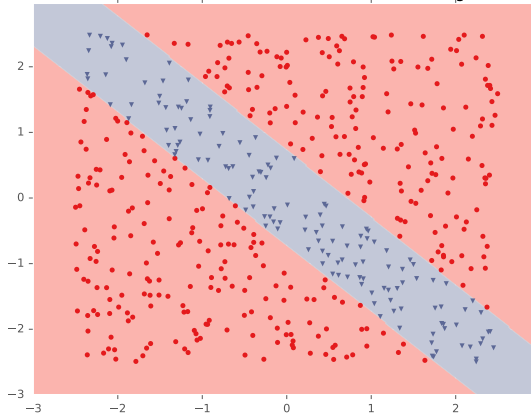
LR1 for Tuned Neural Network (hidden=2, activation=logistic)



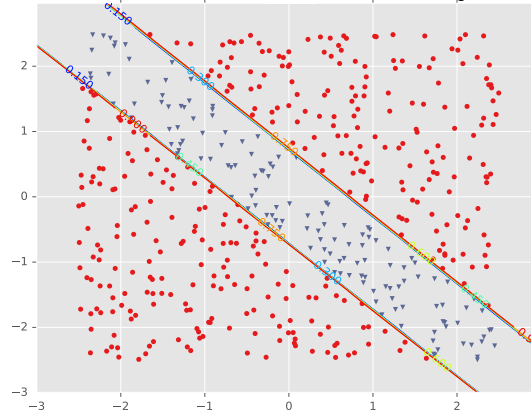
LR2 for Tuned Neural Network (hidden=2, activation=logistic)



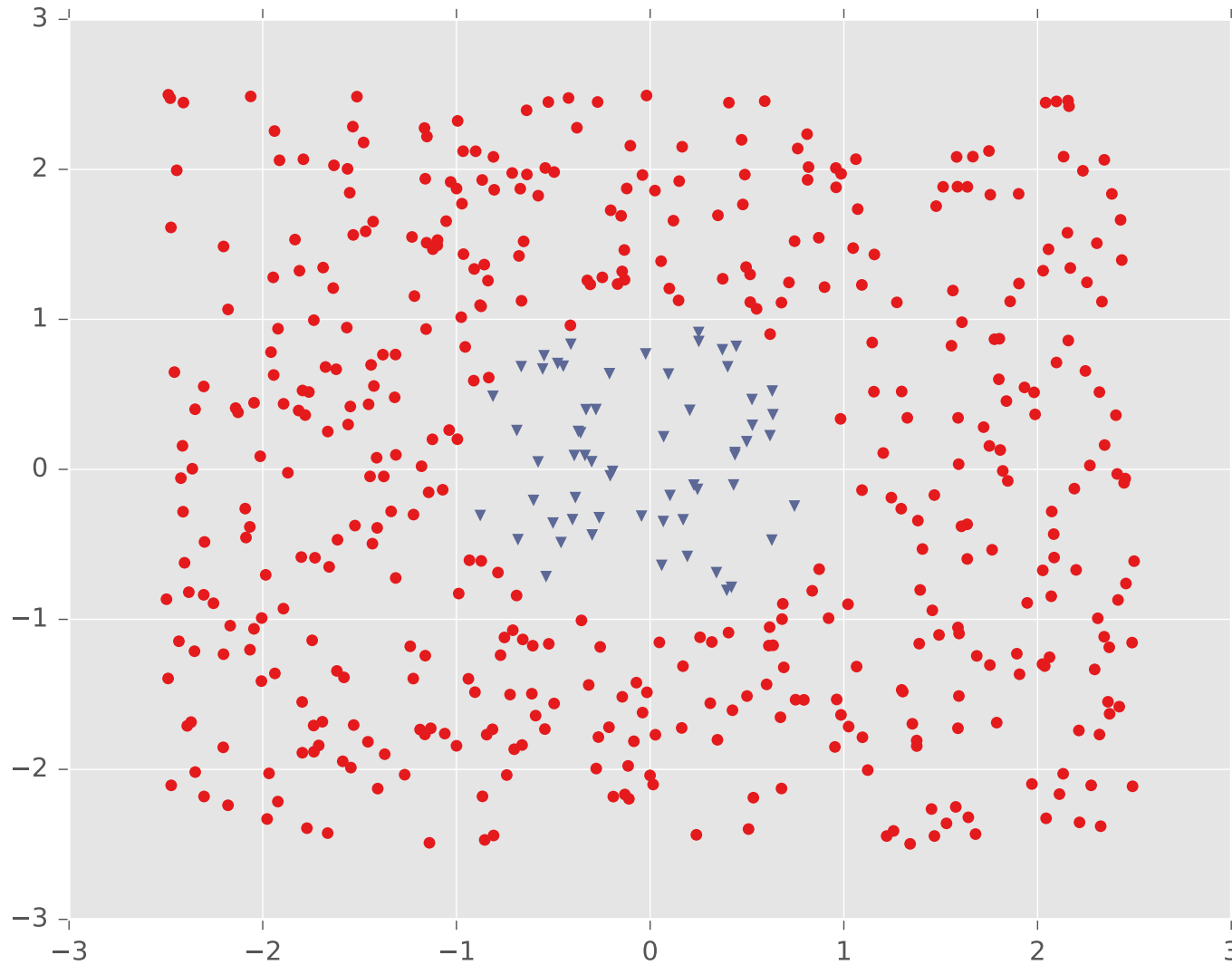
Tuned Neural Network (hidden=2, activation=logistic)



Tuned Neural Network (hidden=2, activation=logistic)



Example #2: One Pocket

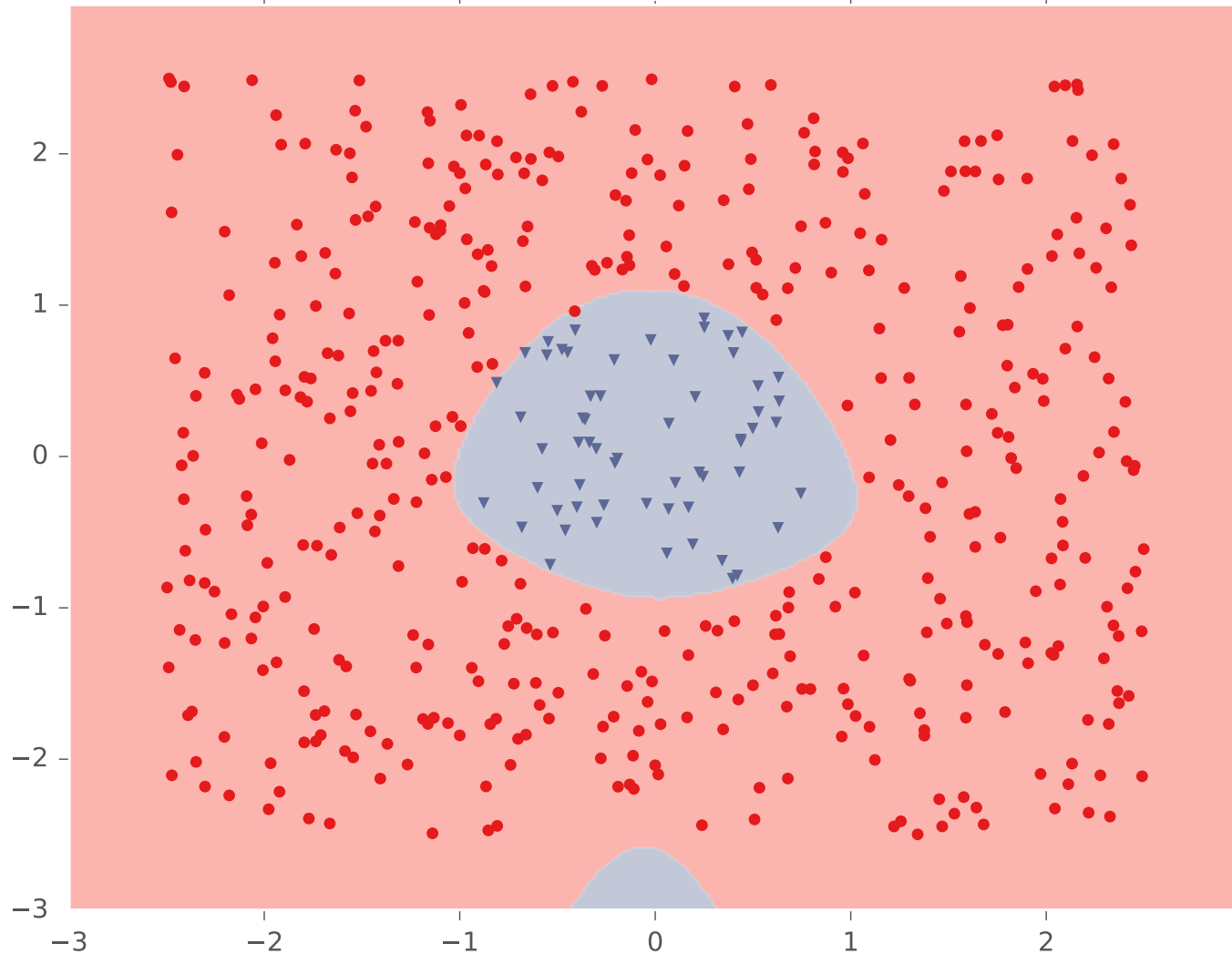


Example #2: One Pocket



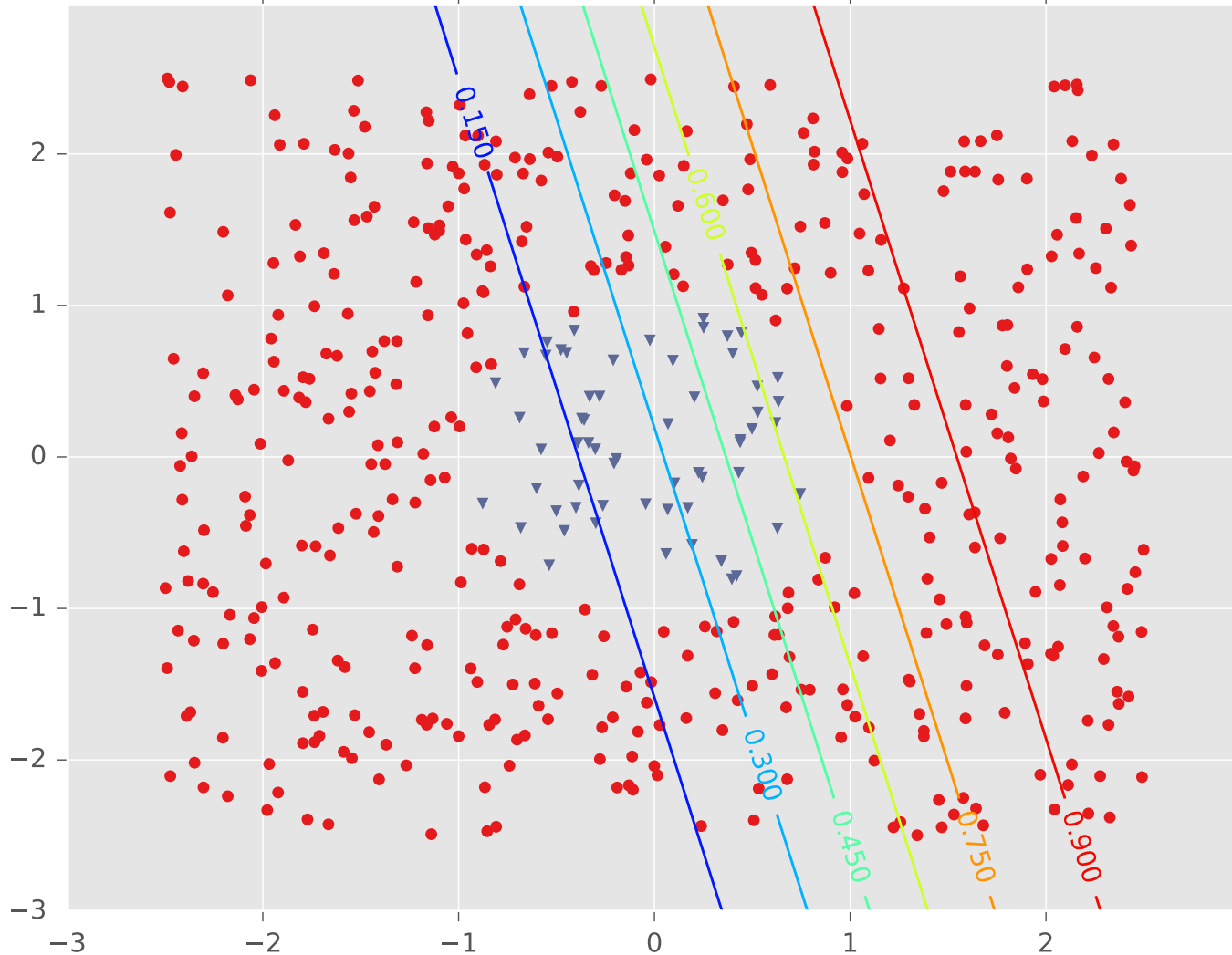
Example #2: One Pocket

Tuned Neural Network (hidden=3, activation=logistic)



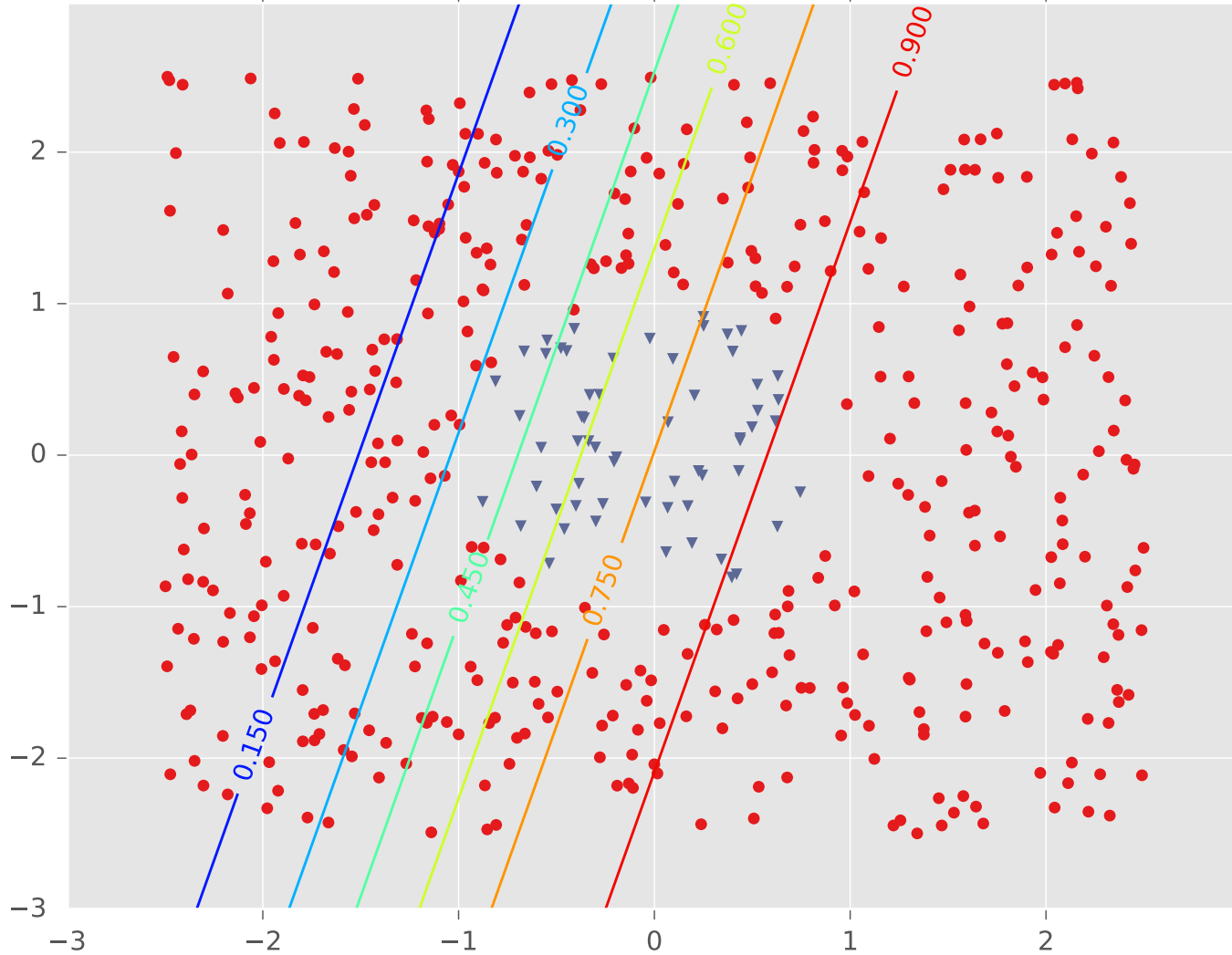
Example #2: One Pocket

LR1 for Tuned Neural Network (hidden=3, activation=logistic)



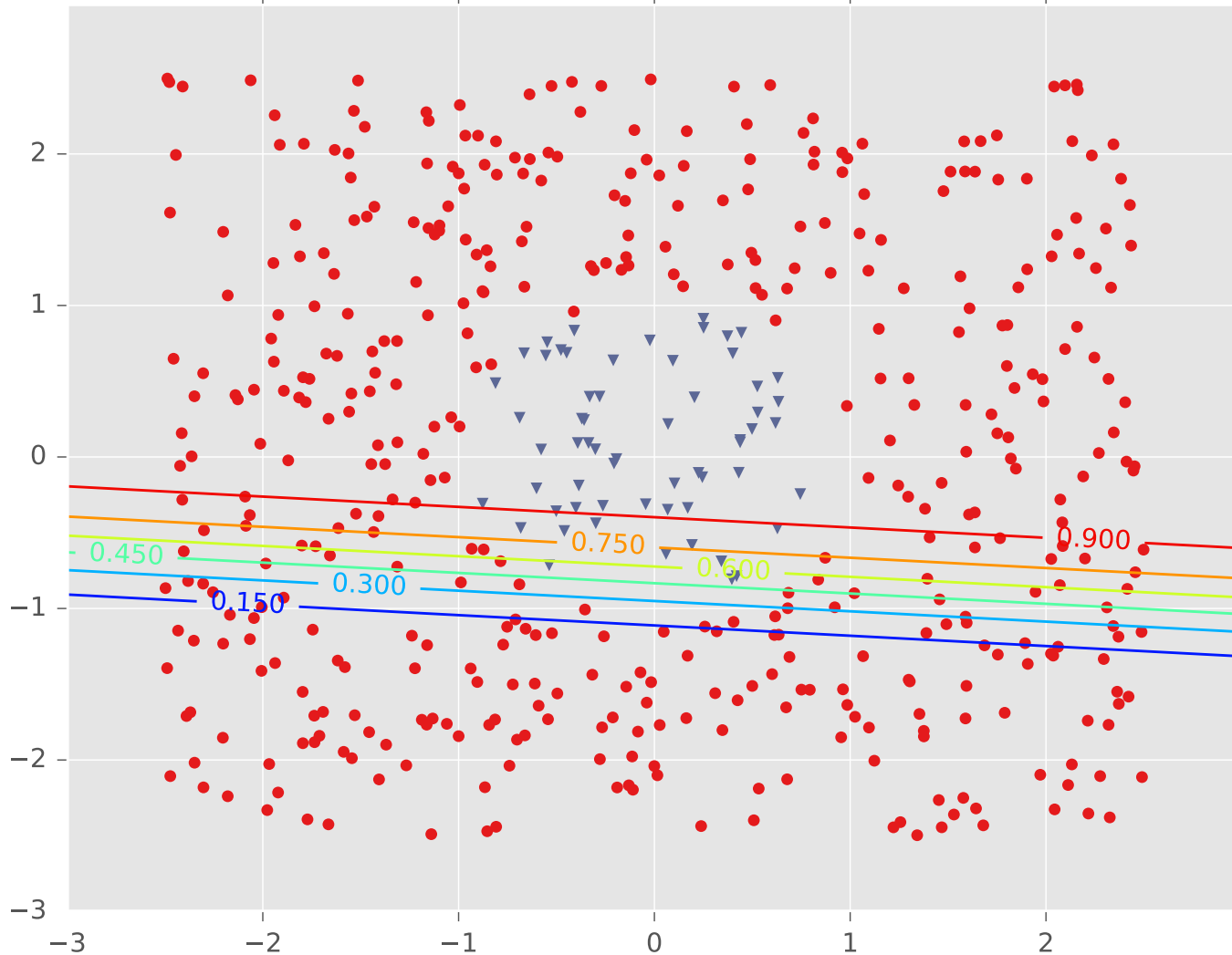
Example #2: One Pocket

LR2 for Tuned Neural Network (hidden=3, activation=logistic)



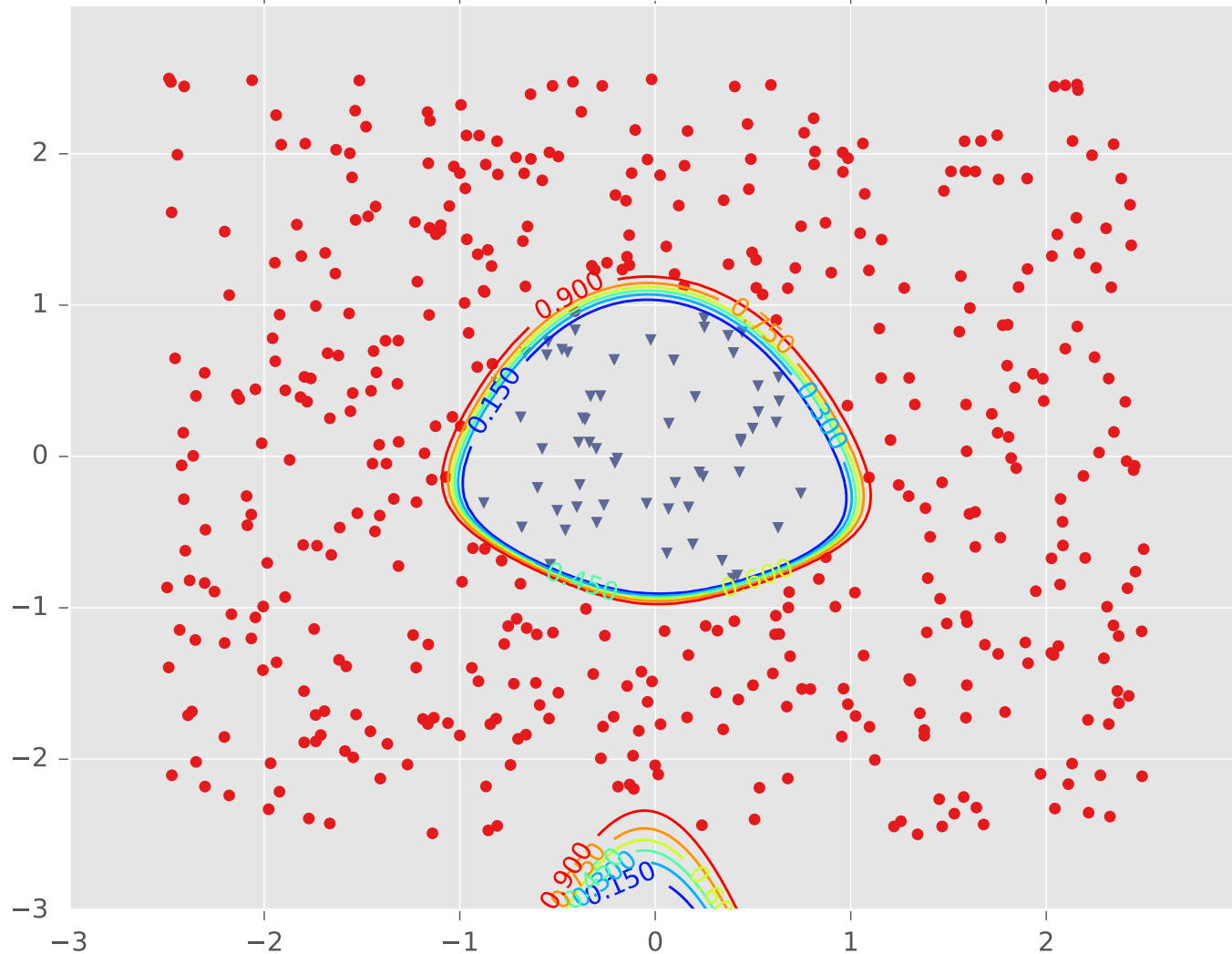
Example #2: One Pocket

LR3 for Tuned Neural Network (hidden=3, activation=logistic)



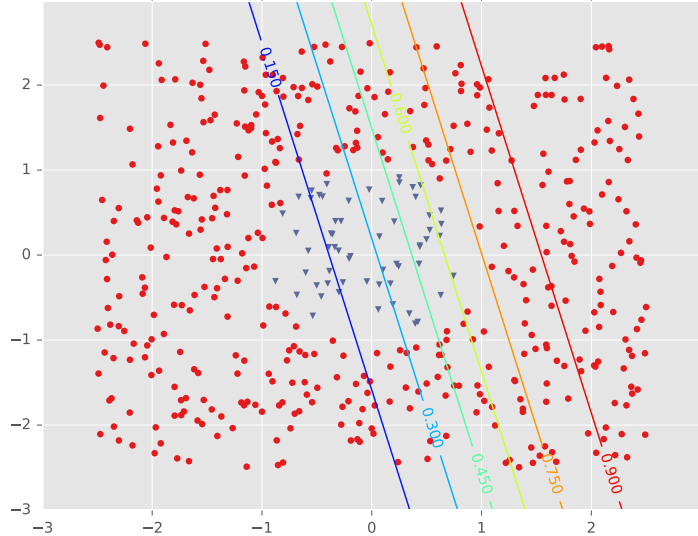
Example #2: One Pocket

Tuned Neural Network (hidden=3, activation=logistic)

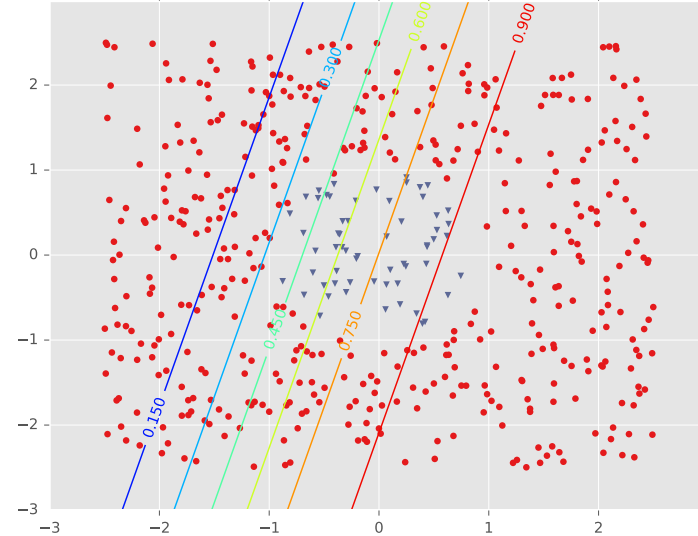


Example #2: One Pocket

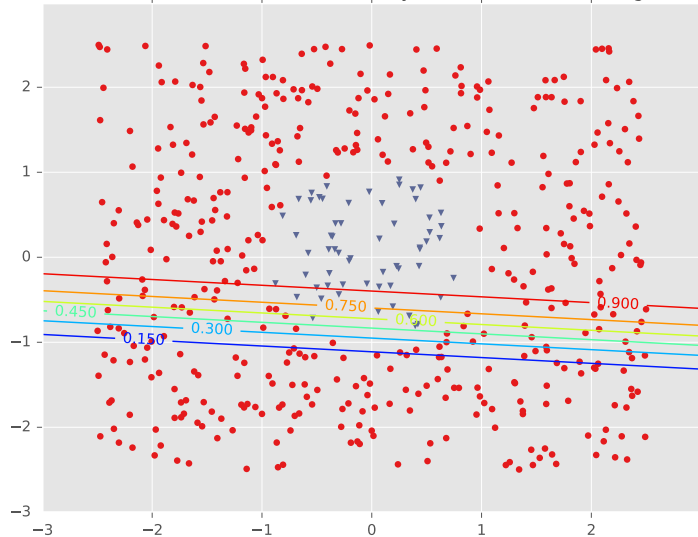
LR1 for Tuned Neural Network (hidden=3, activation=logistic)



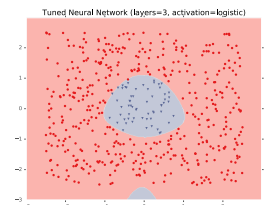
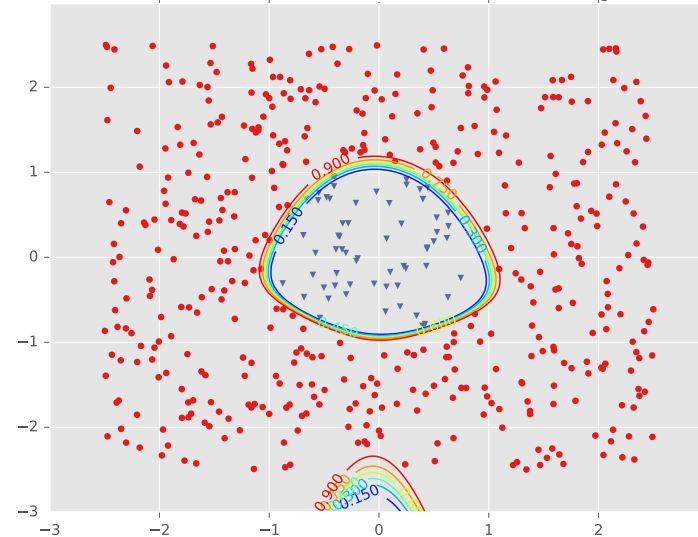
LR2 for Tuned Neural Network (hidden=3, activation=logistic)



LR3 for Tuned Neural Network (hidden=3, activation=logistic)



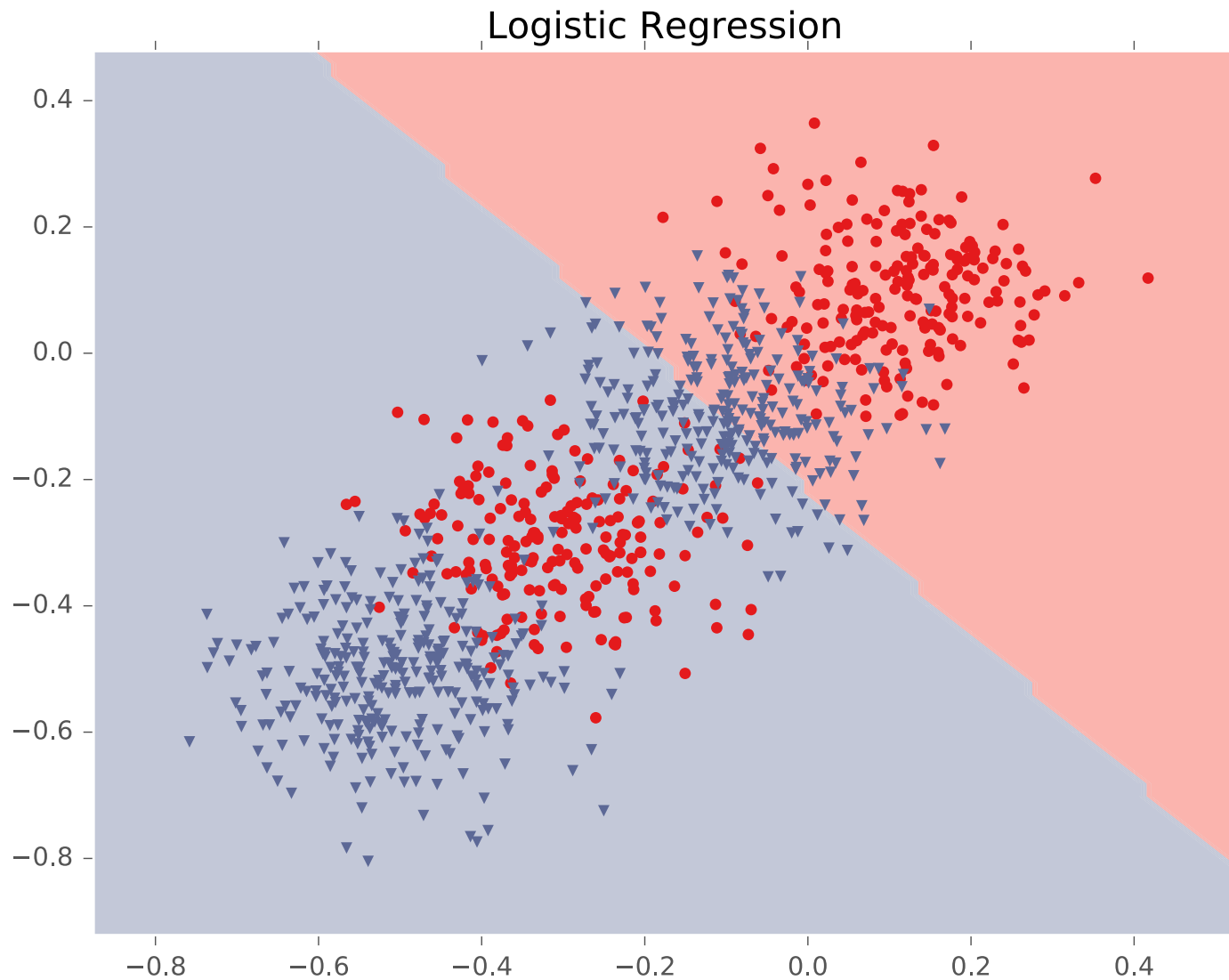
Tuned Neural Network (hidden=3, activation=logistic)



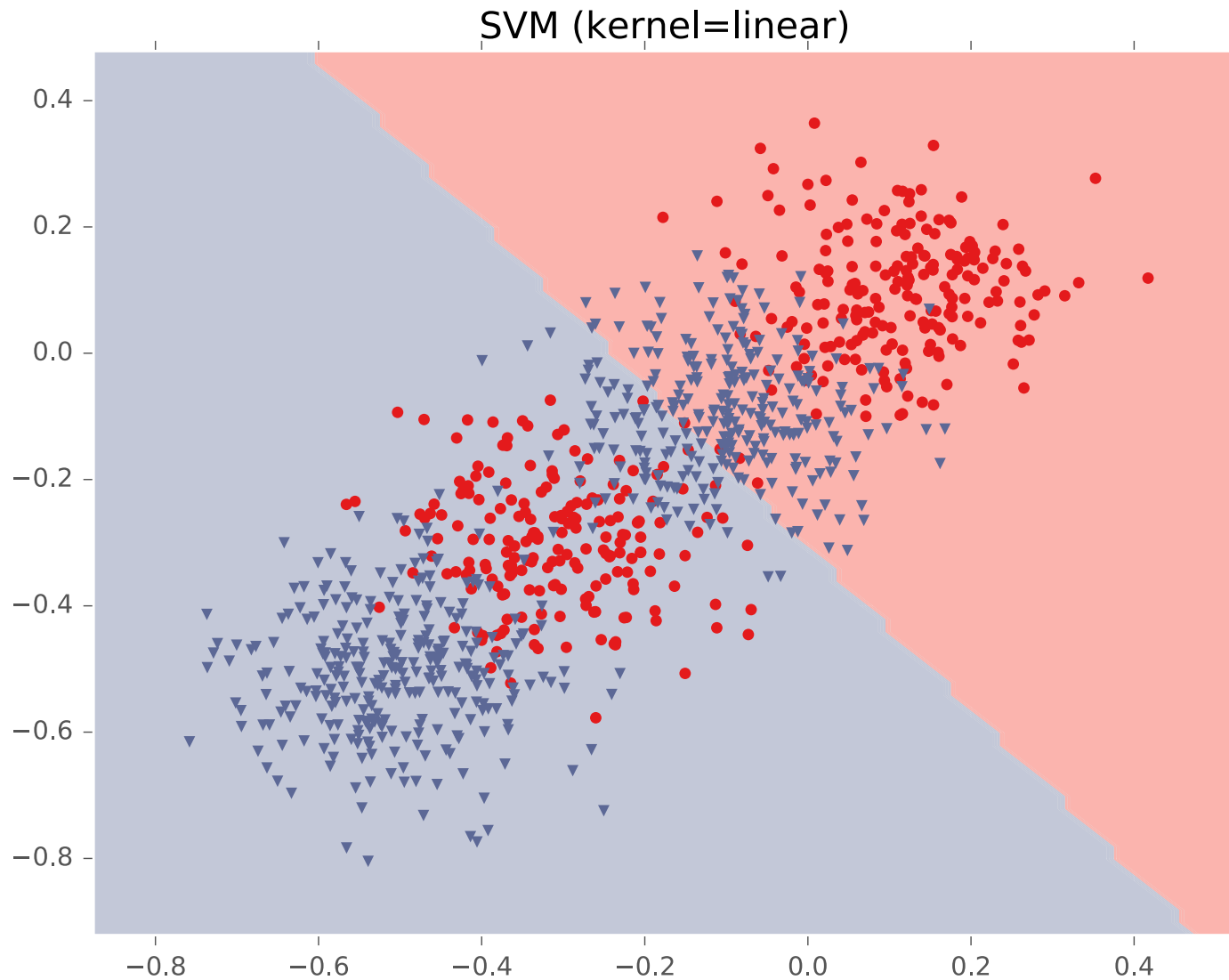
Example #3: Four Gaussians



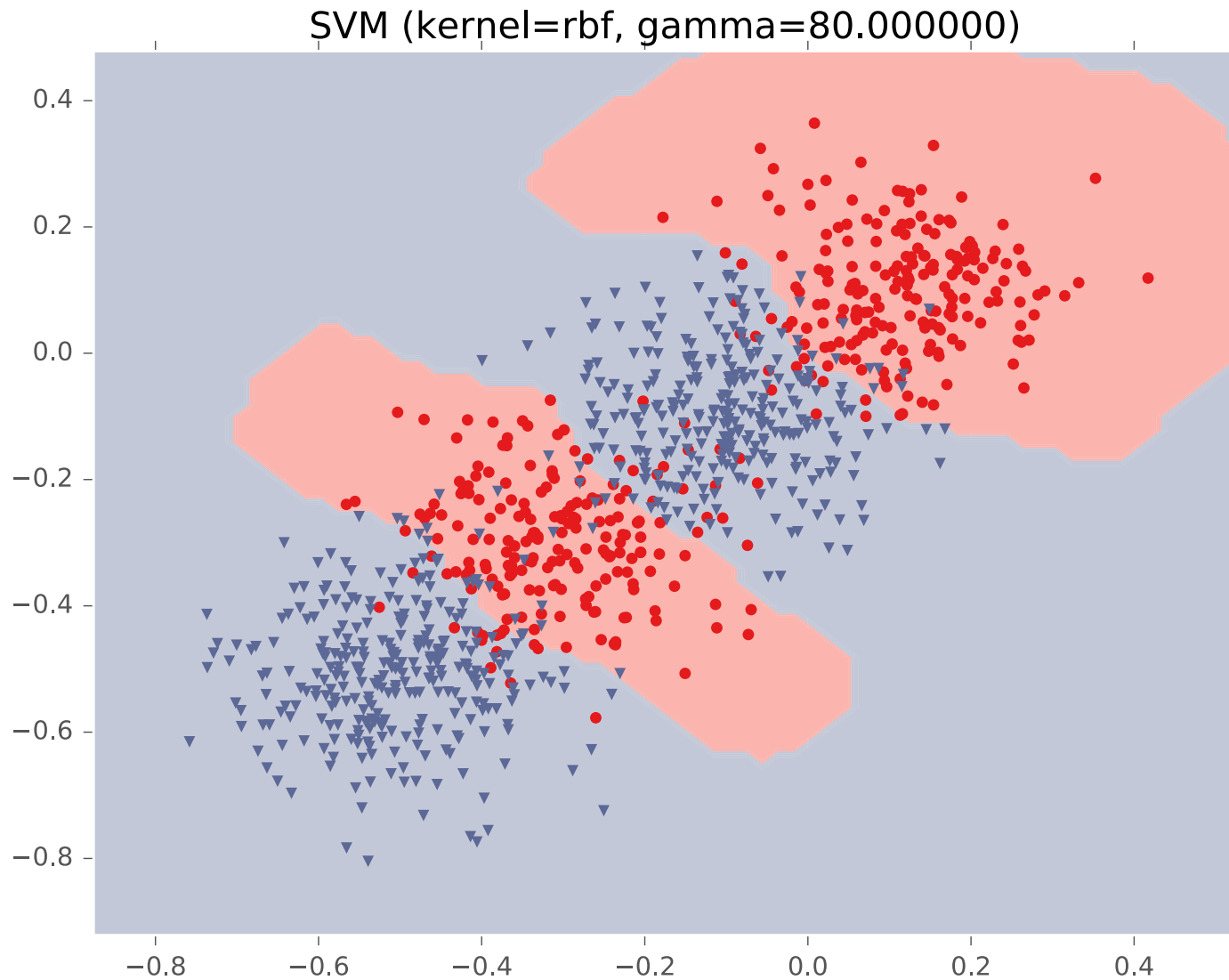
Example #3: Four Gaussians



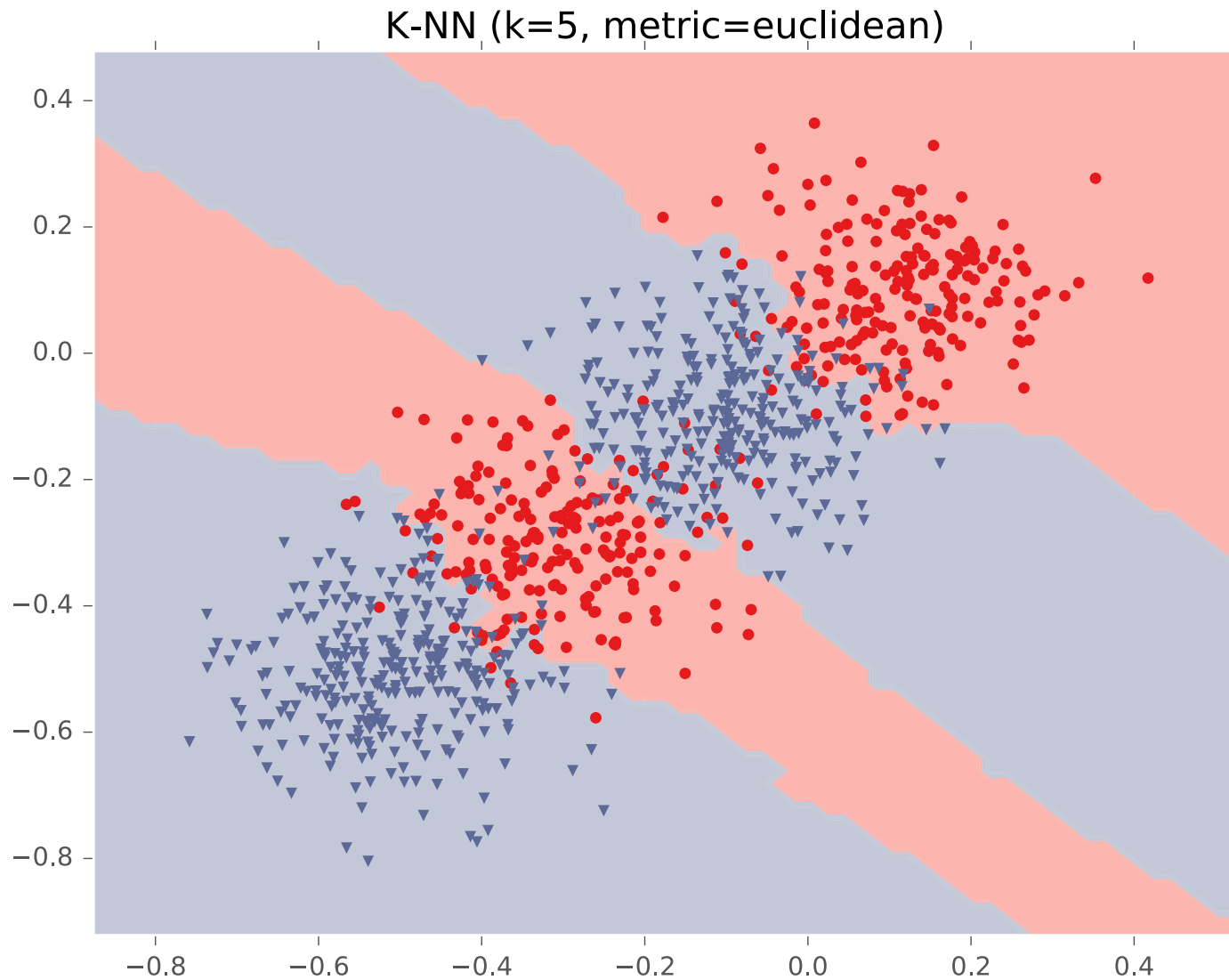
Example #3: Four Gaussians



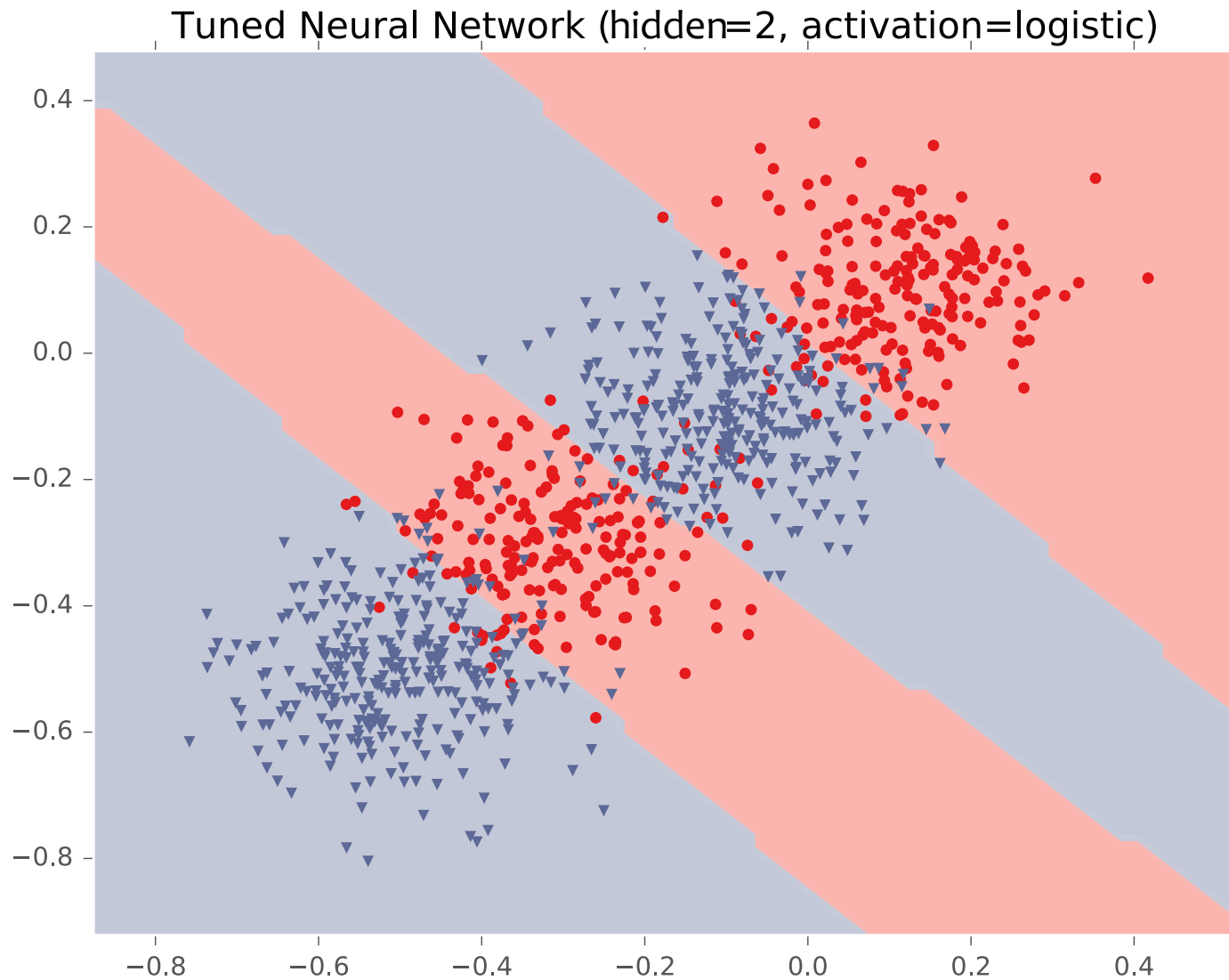
Example #3: Four Gaussians



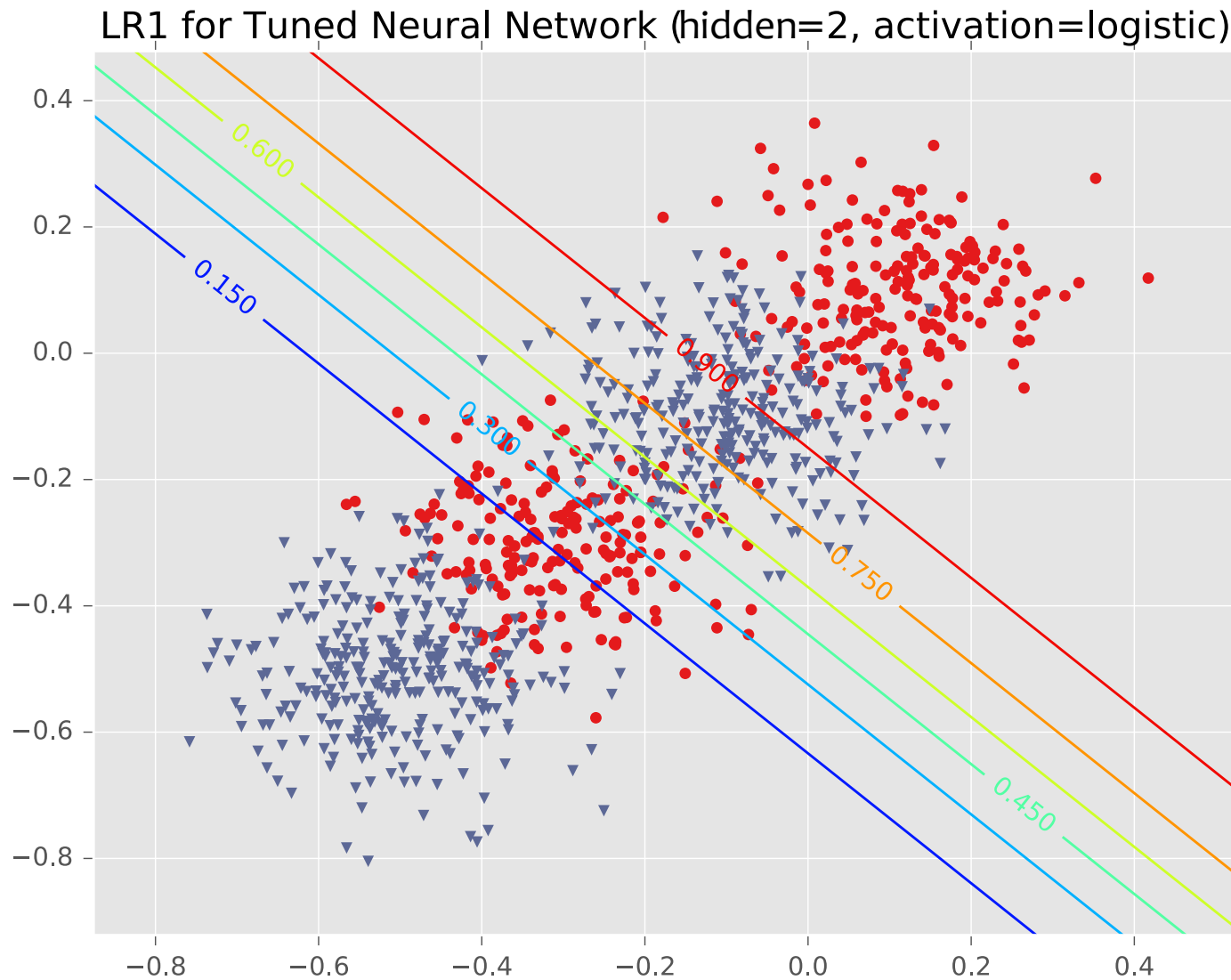
Example #3: Four Gaussians



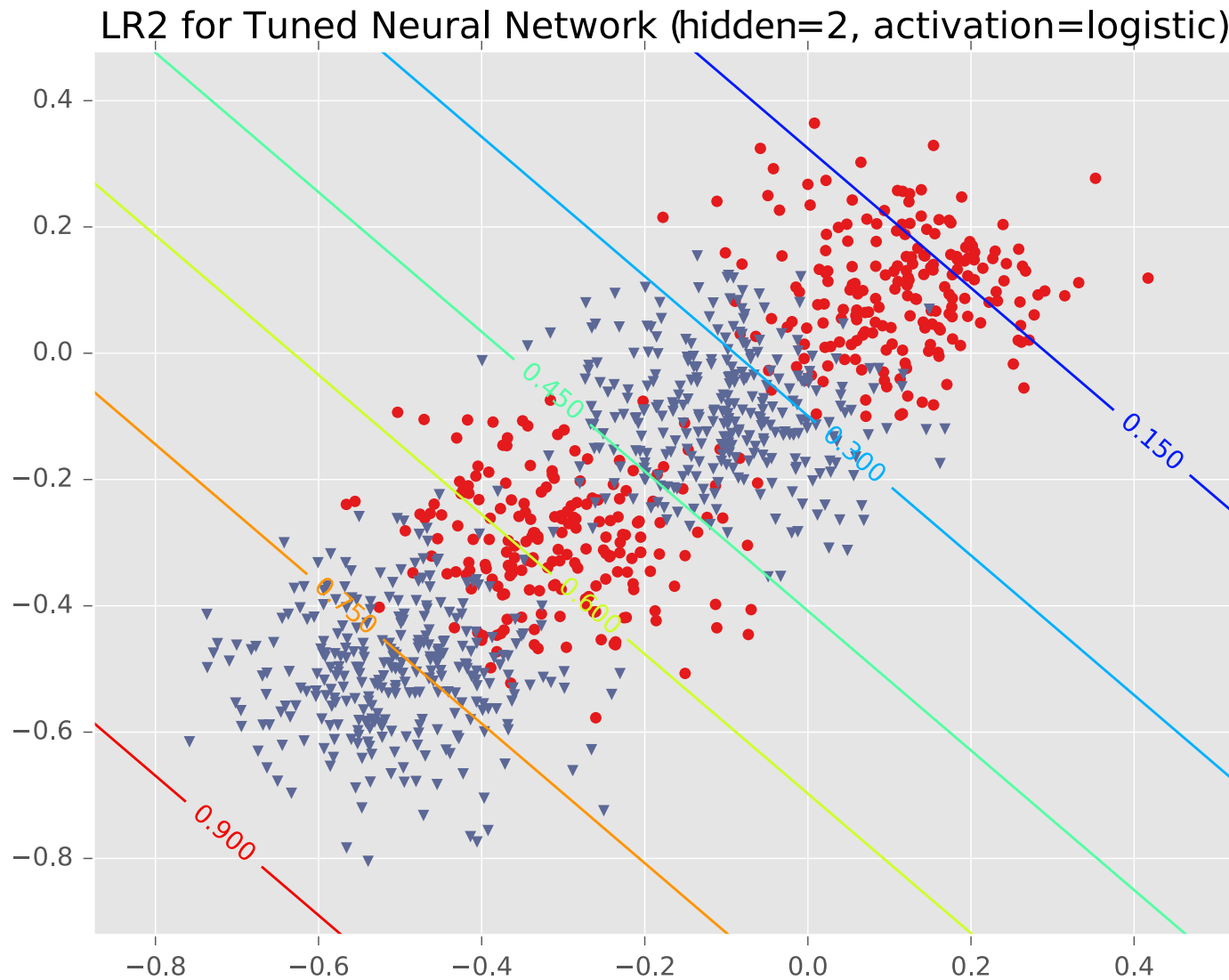
Example #3: Four Gaussians



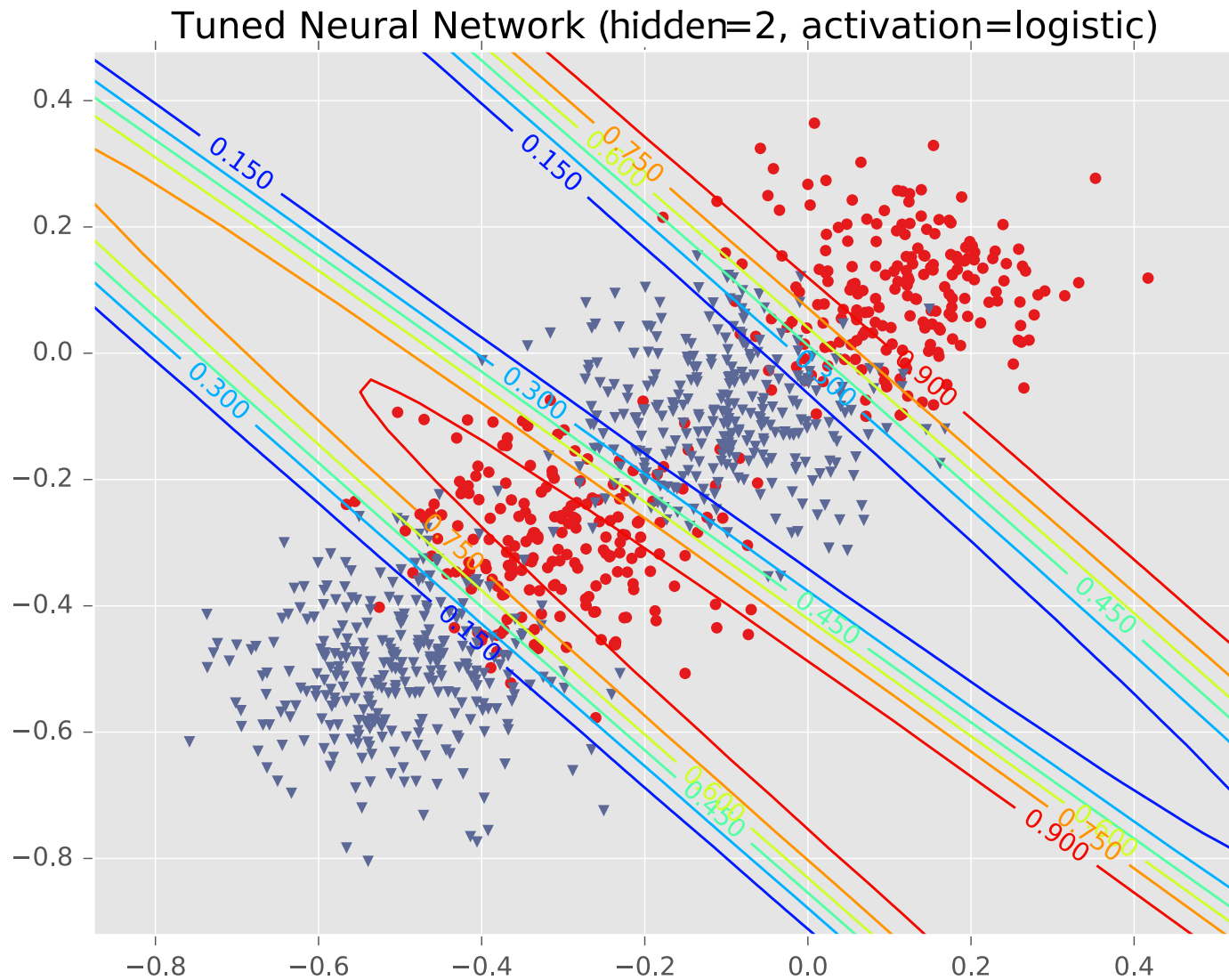
Example #3: Four Gaussians



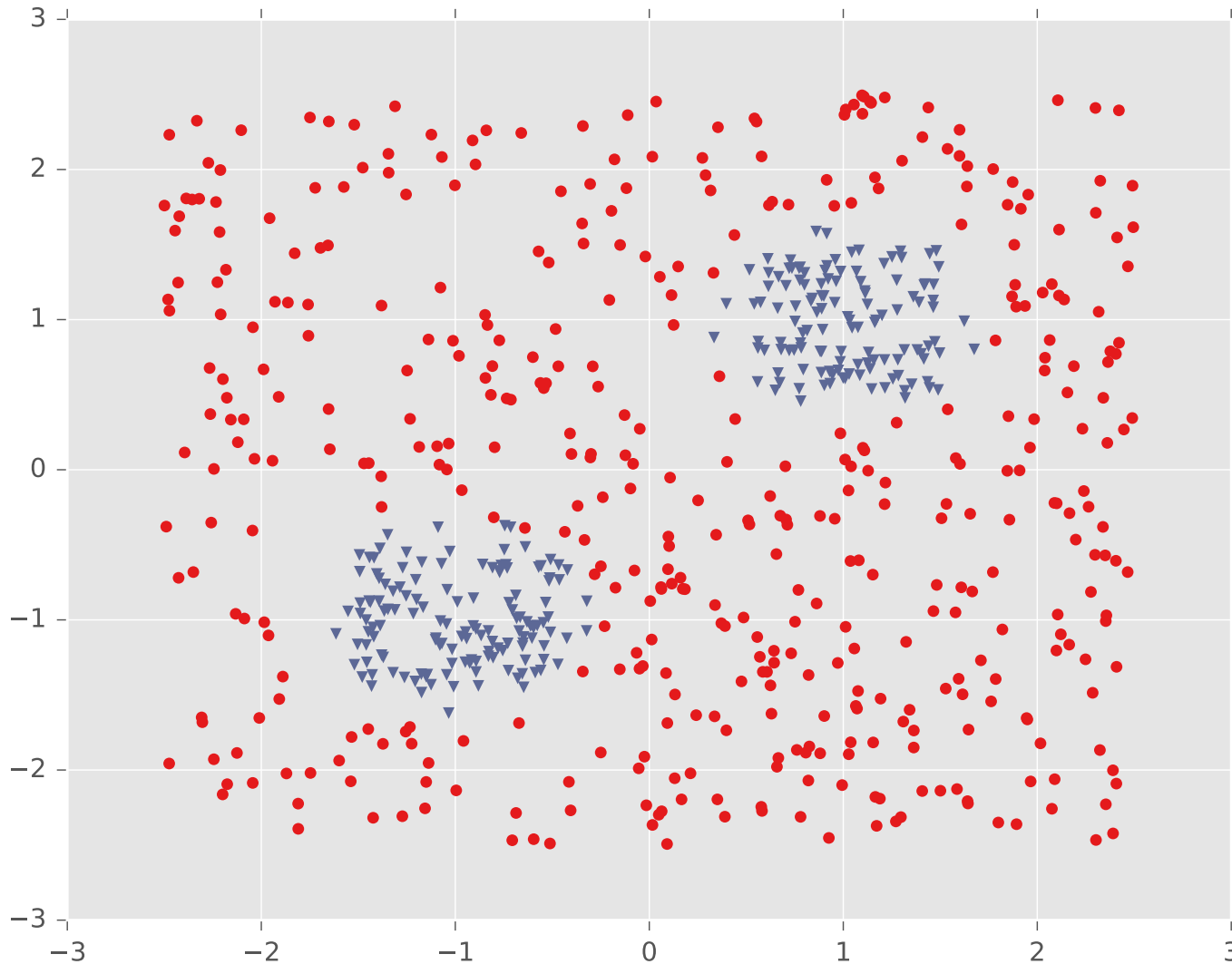
Example #3: Four Gaussians



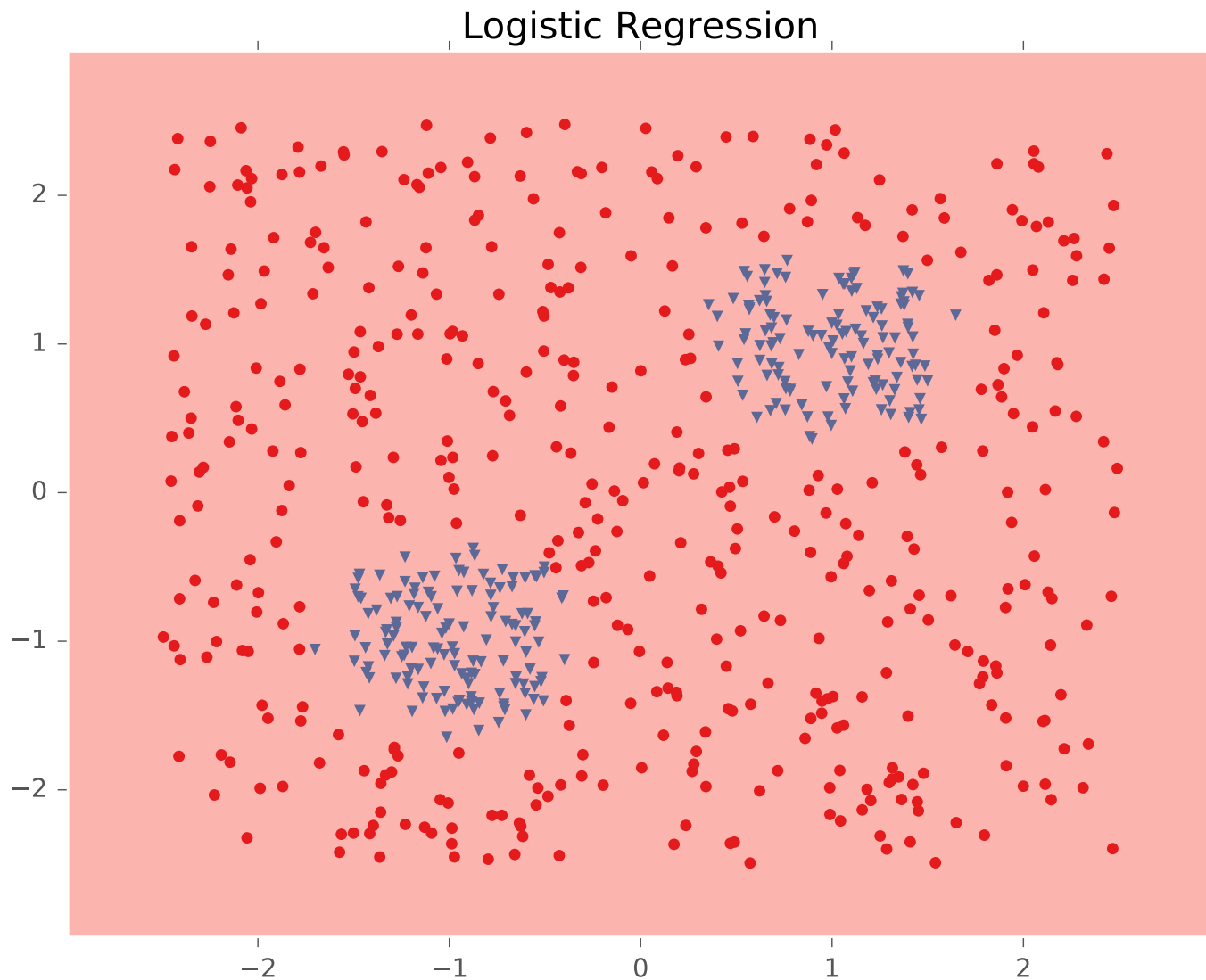
Example #3: Four Gaussians



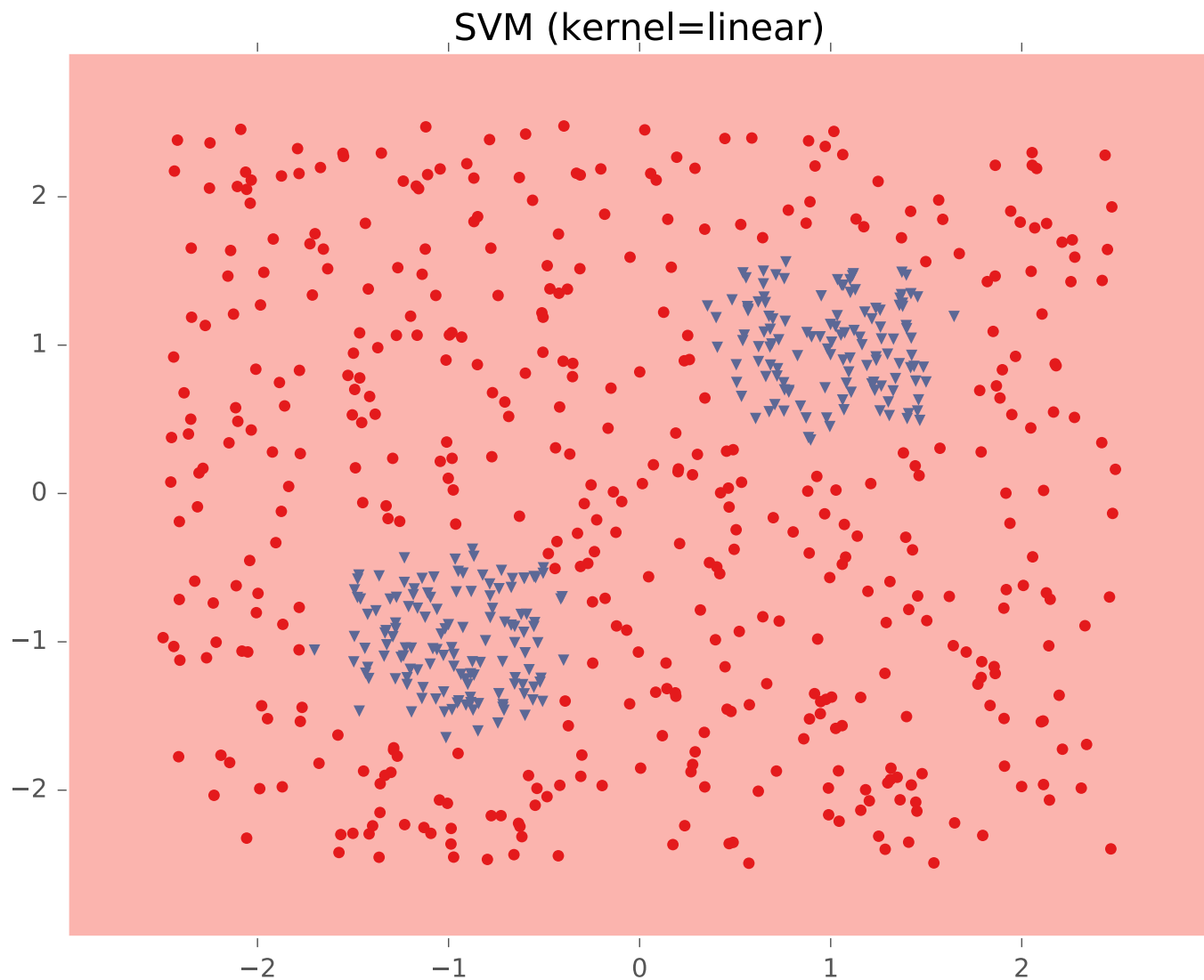
Example #4: Two Pockets



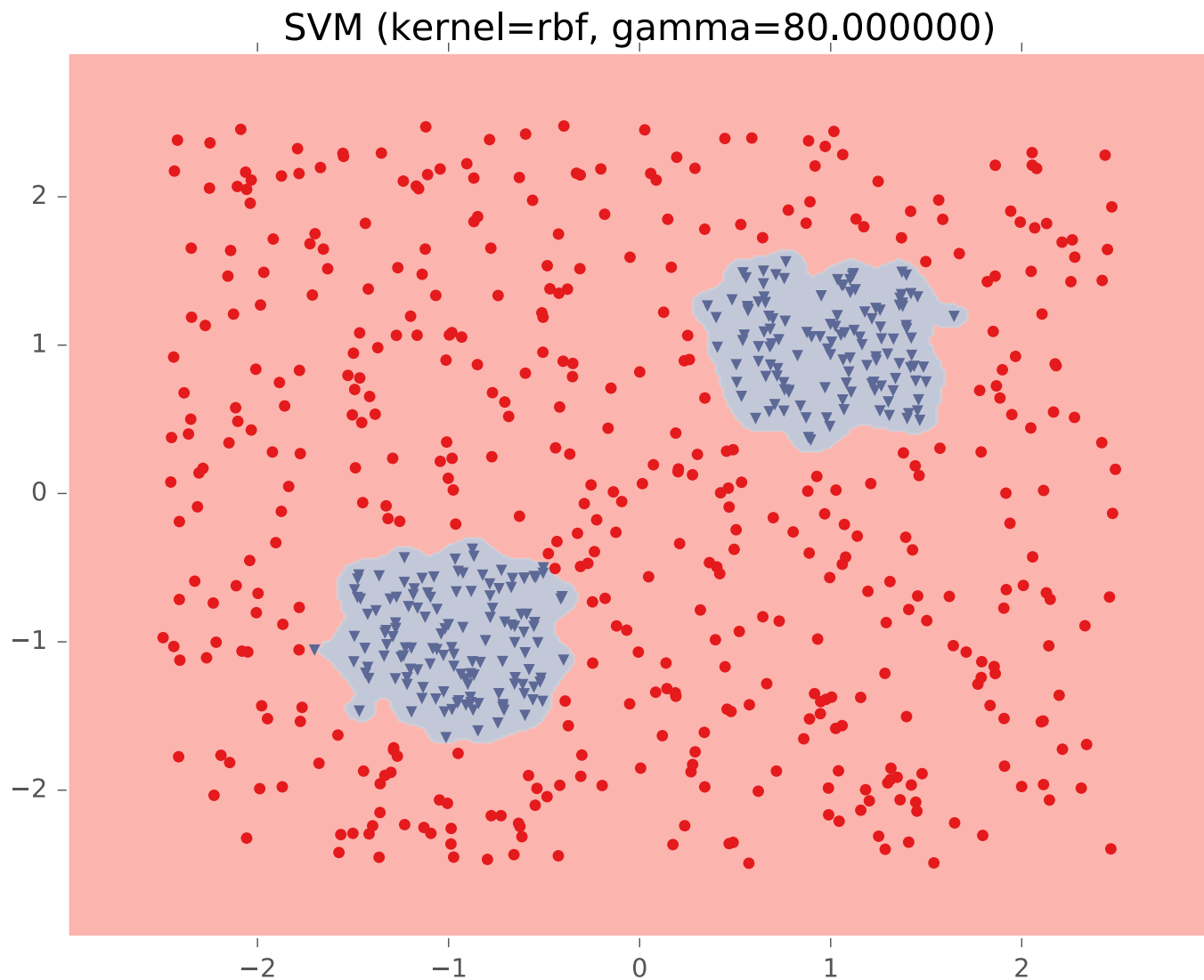
Example #4: Two Pockets



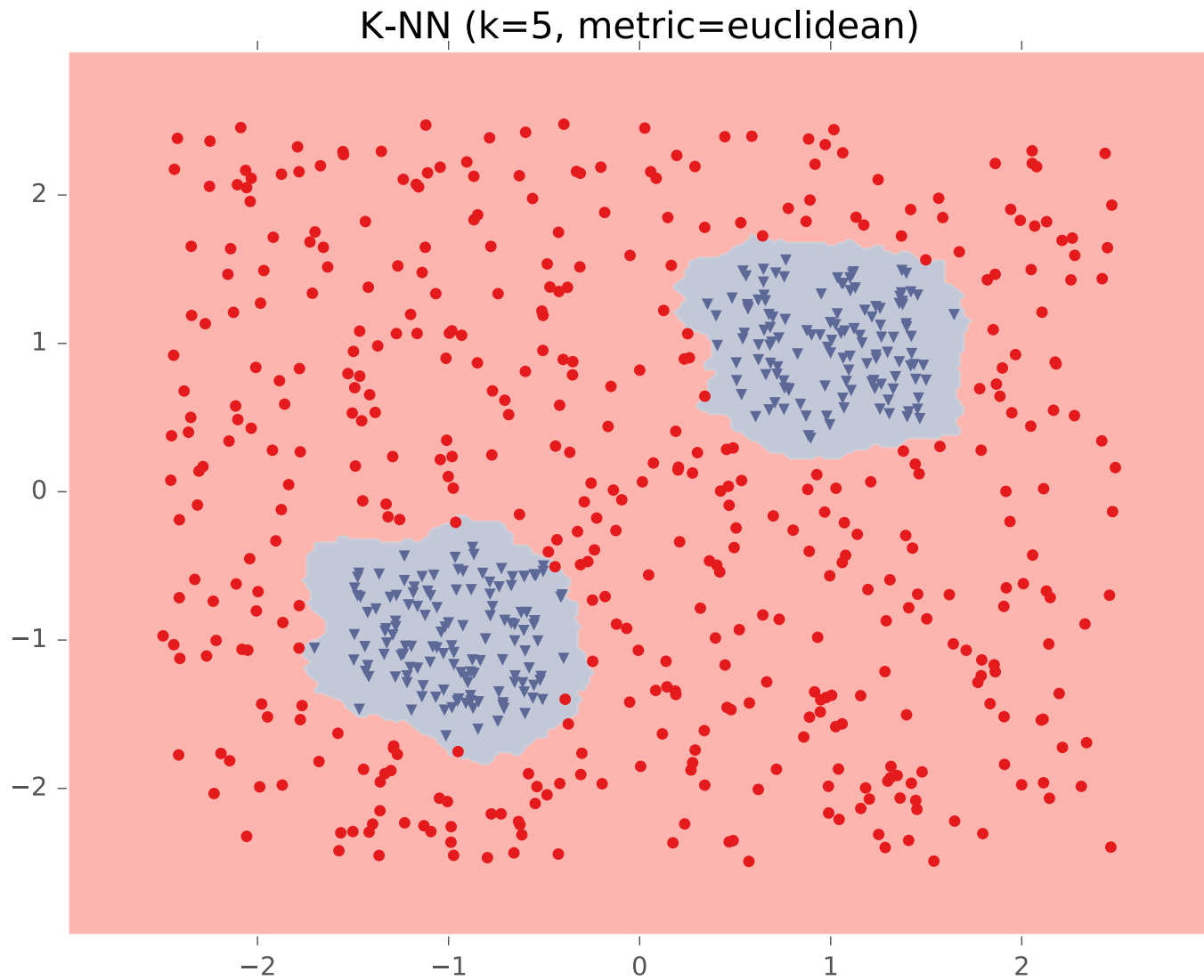
Example #4: Two Pockets



Example #4: Two Pockets

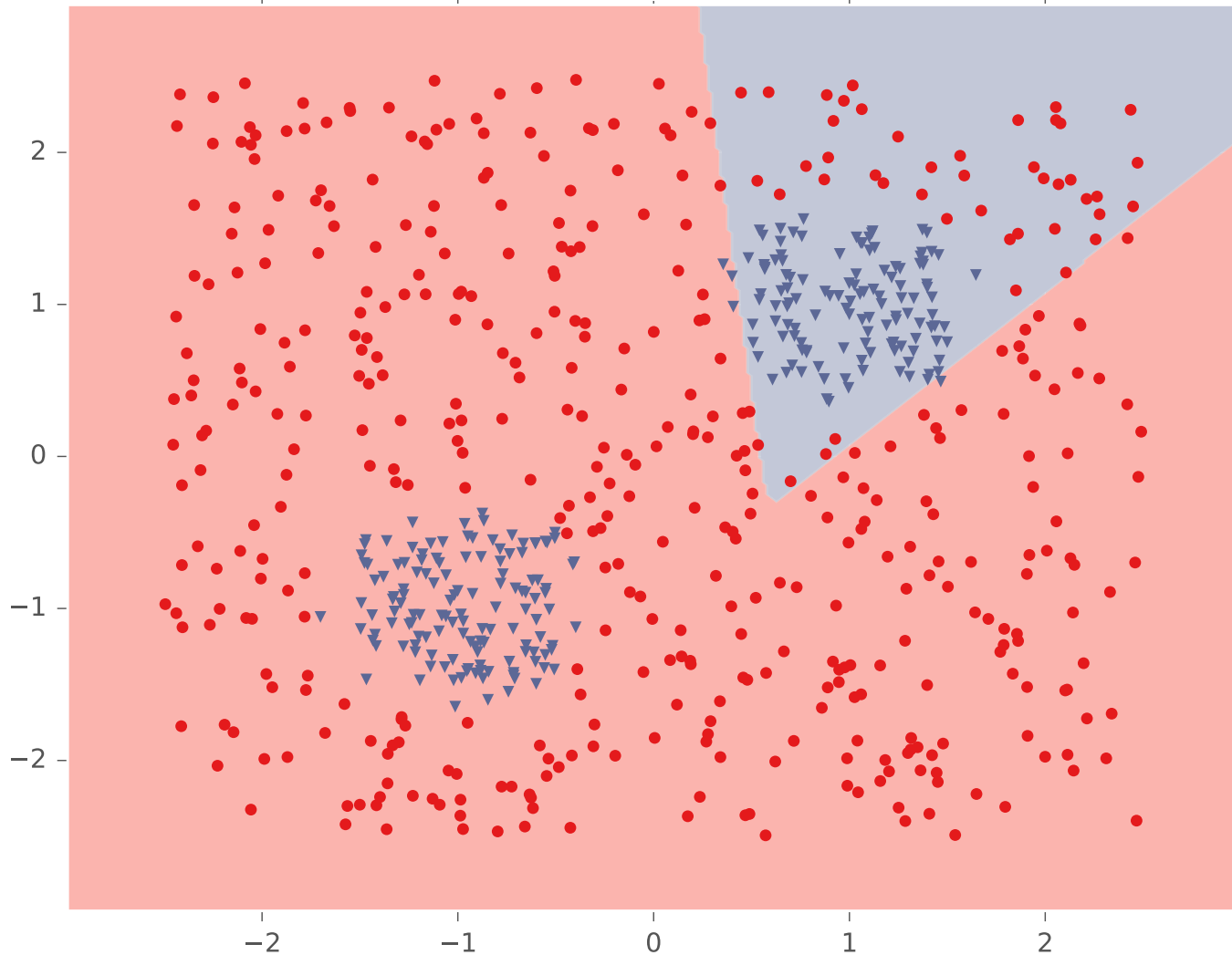


Example #4: Two Pockets



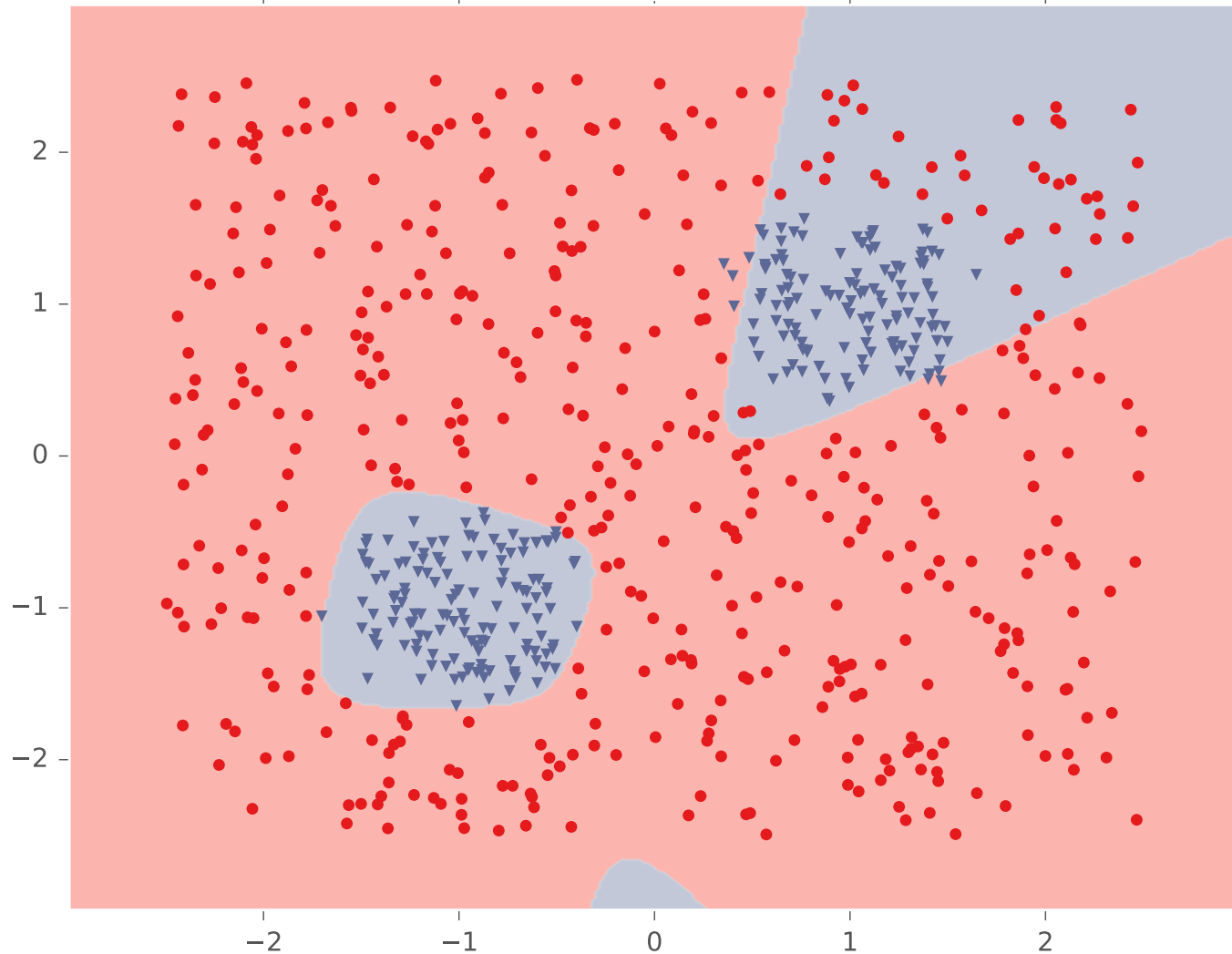
Example #4: Two Pockets

Tuned Neural Network (hidden=2, activation=logistic)



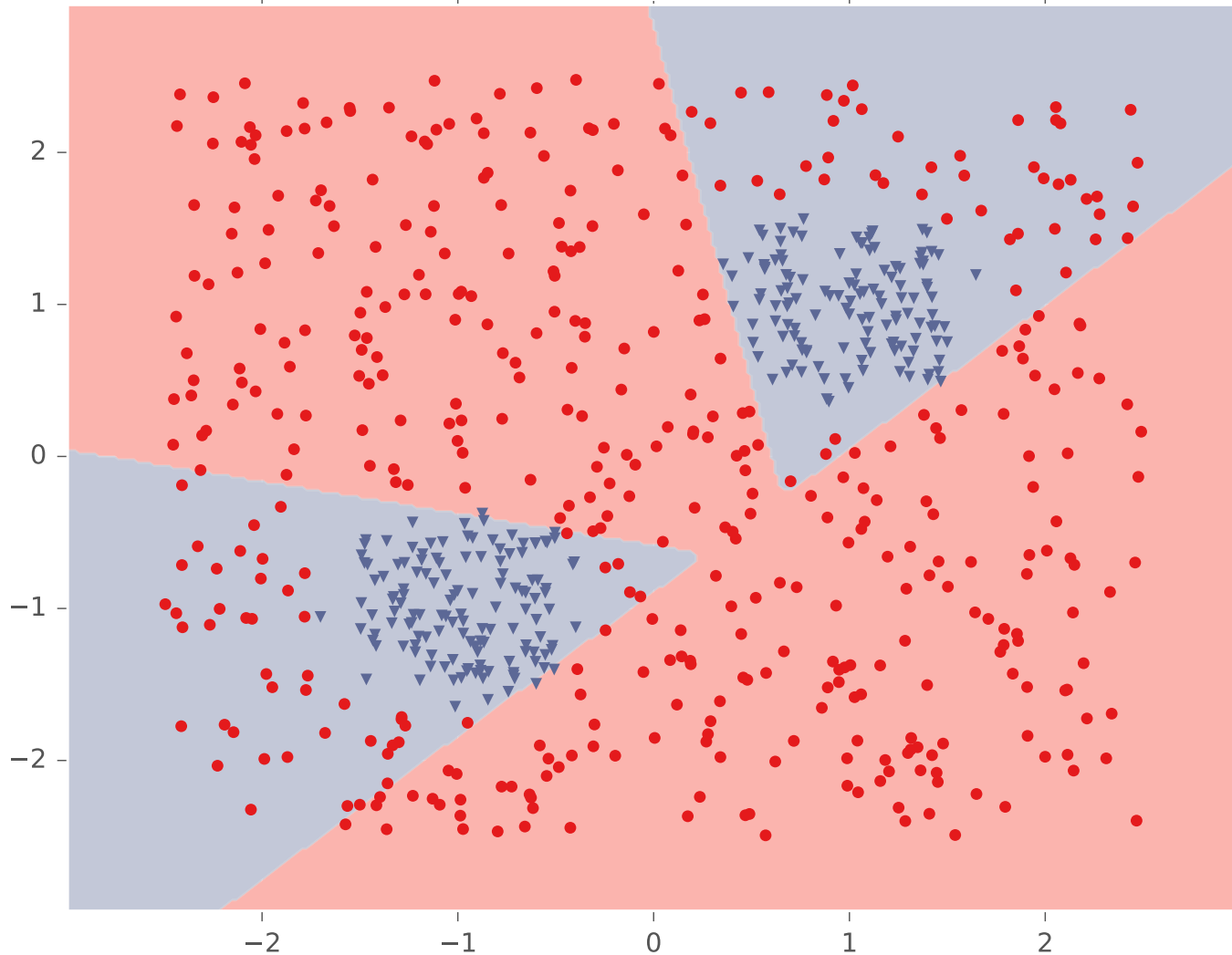
Example #4: Two Pockets

Tuned Neural Network (hidden=3, activation=logistic)



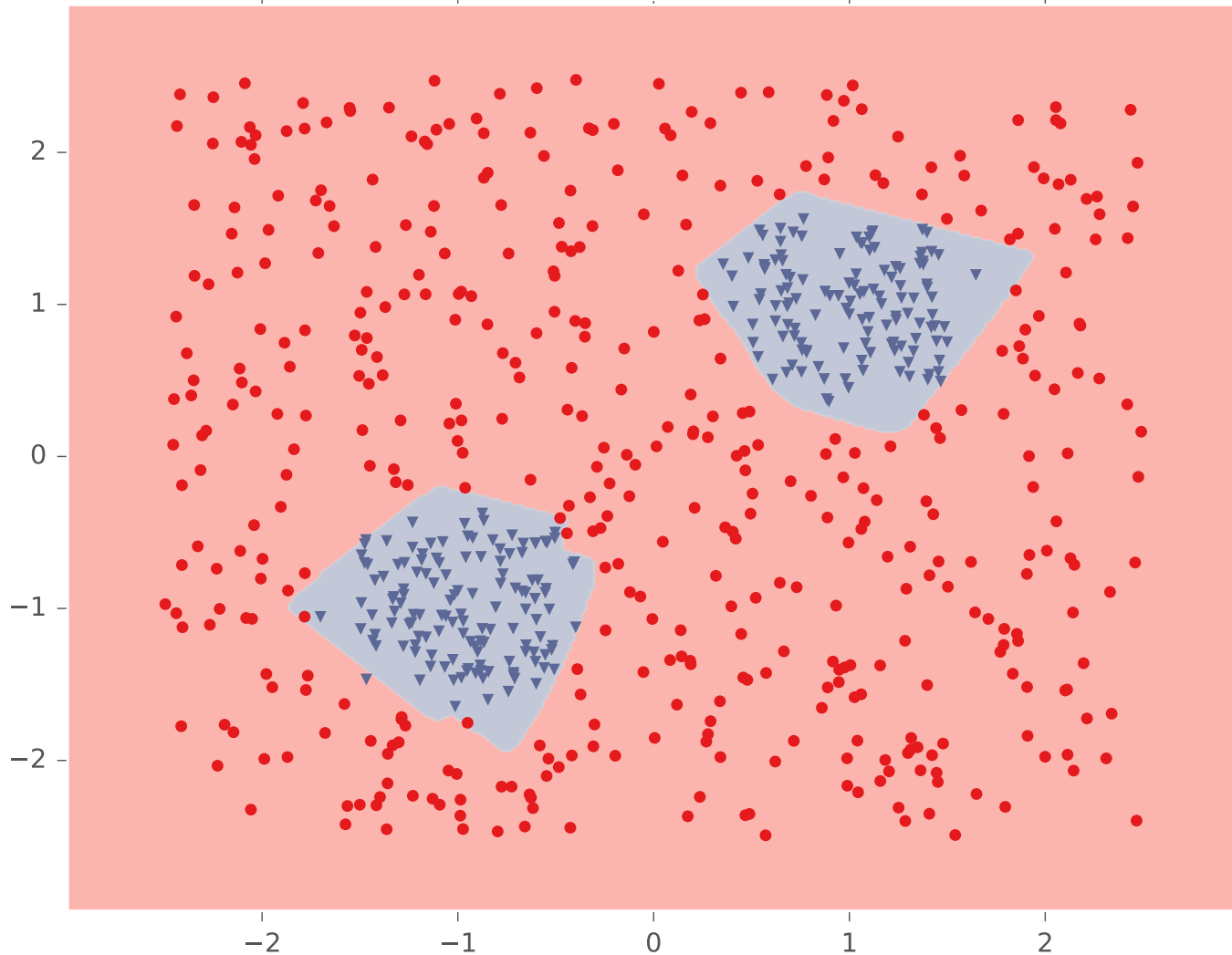
Example #4: Two Pockets

Tuned Neural Network (hidden=4, activation=logistic)



Example #4: Two Pockets

Tuned Neural Network (hidden=10, activation=logistic)



ARCHITECTURES

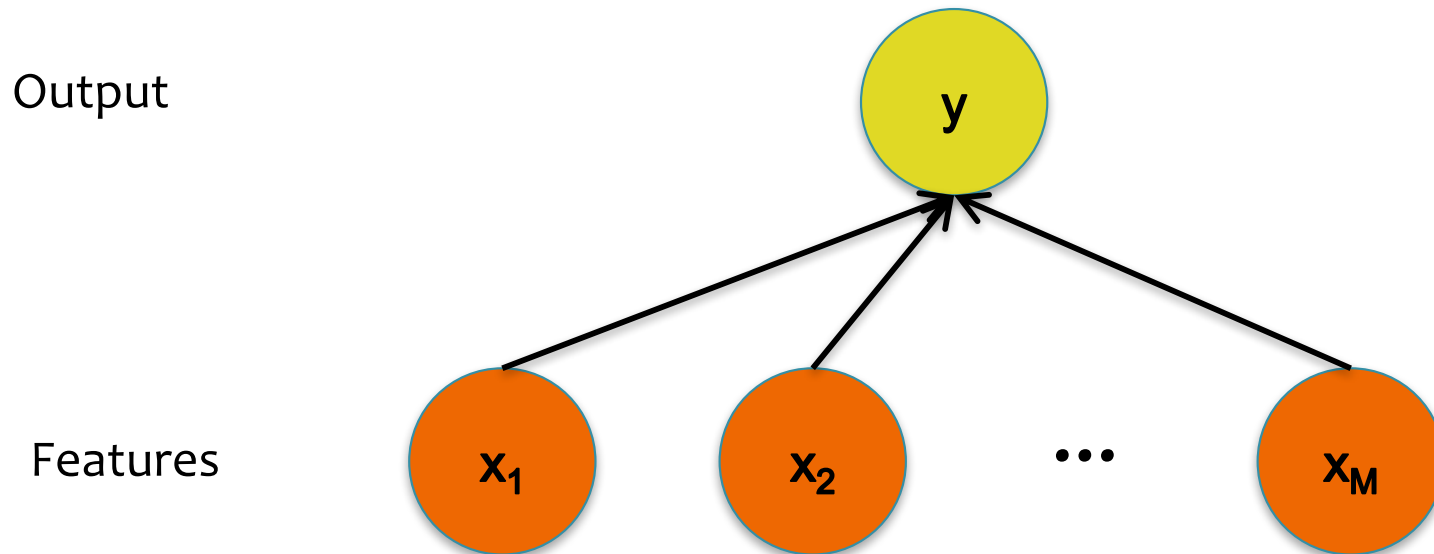
Neural Network Architectures

Even for a basic Neural Network, there are many design decisions to make:

1. # of hidden layers (depth)
2. # of units per hidden layer (width)
3. Type of activation function (nonlinearity)
4. Form of objective function

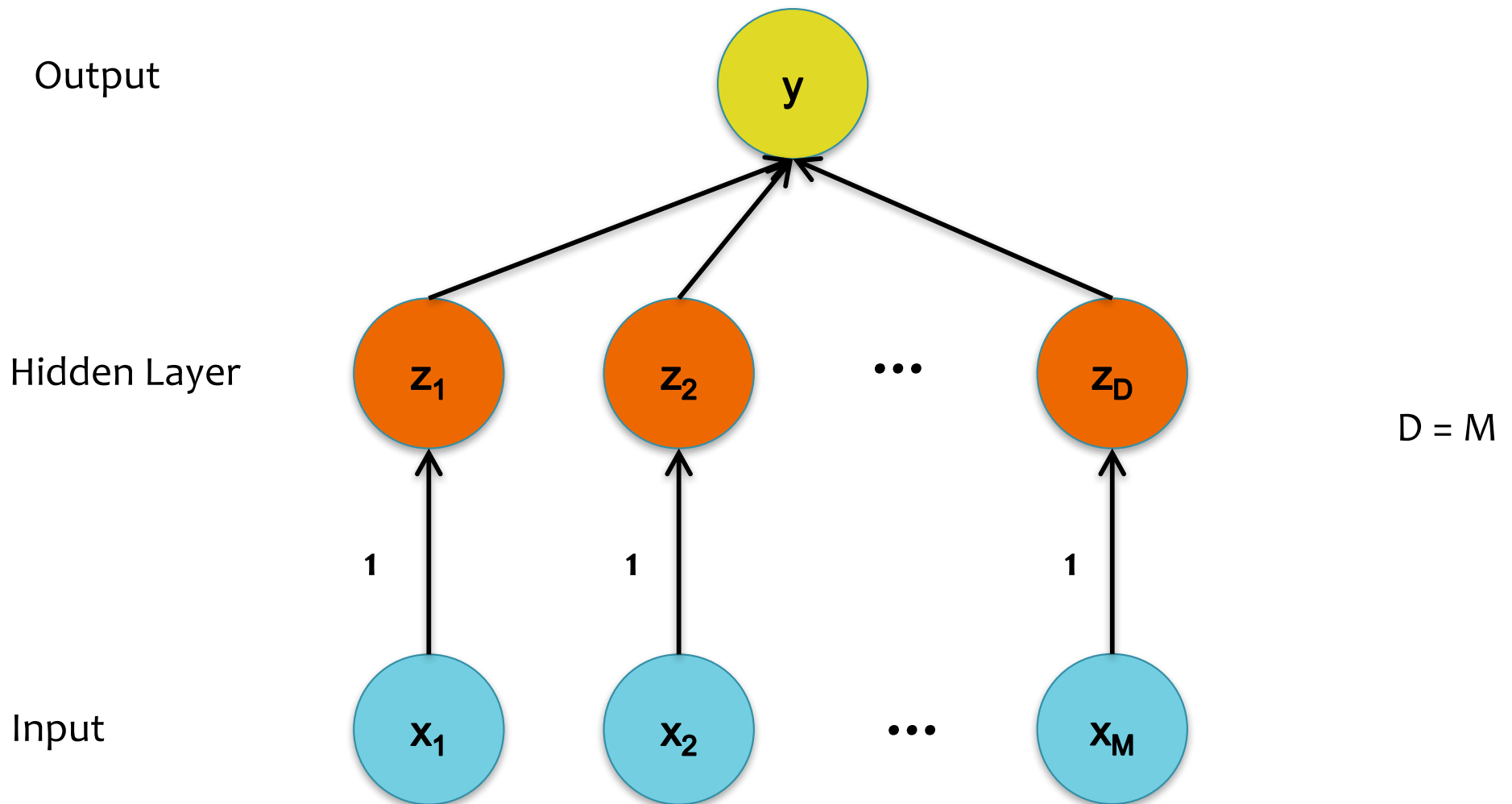
Building a Neural Net

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



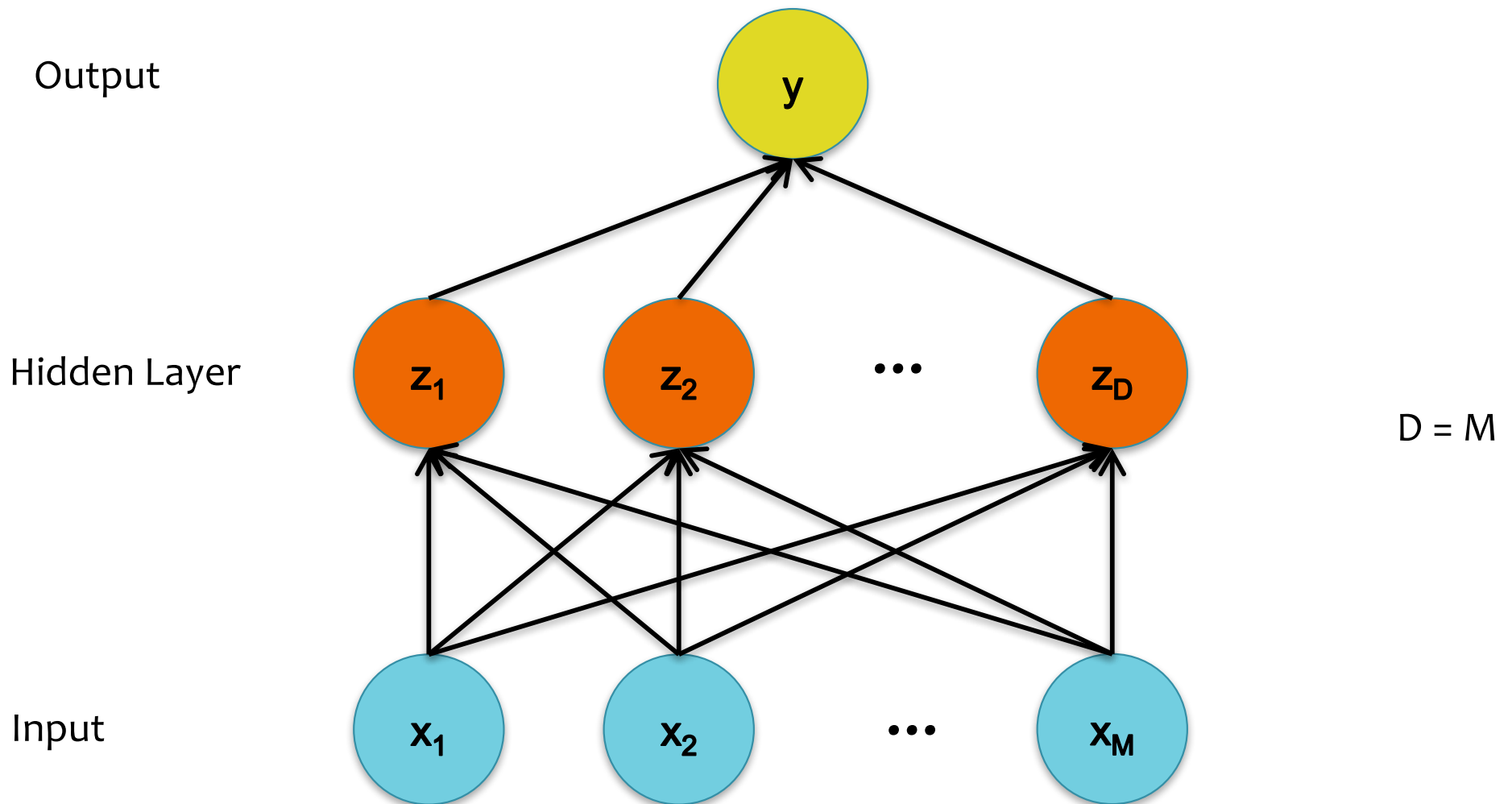
Building a Neural Net

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



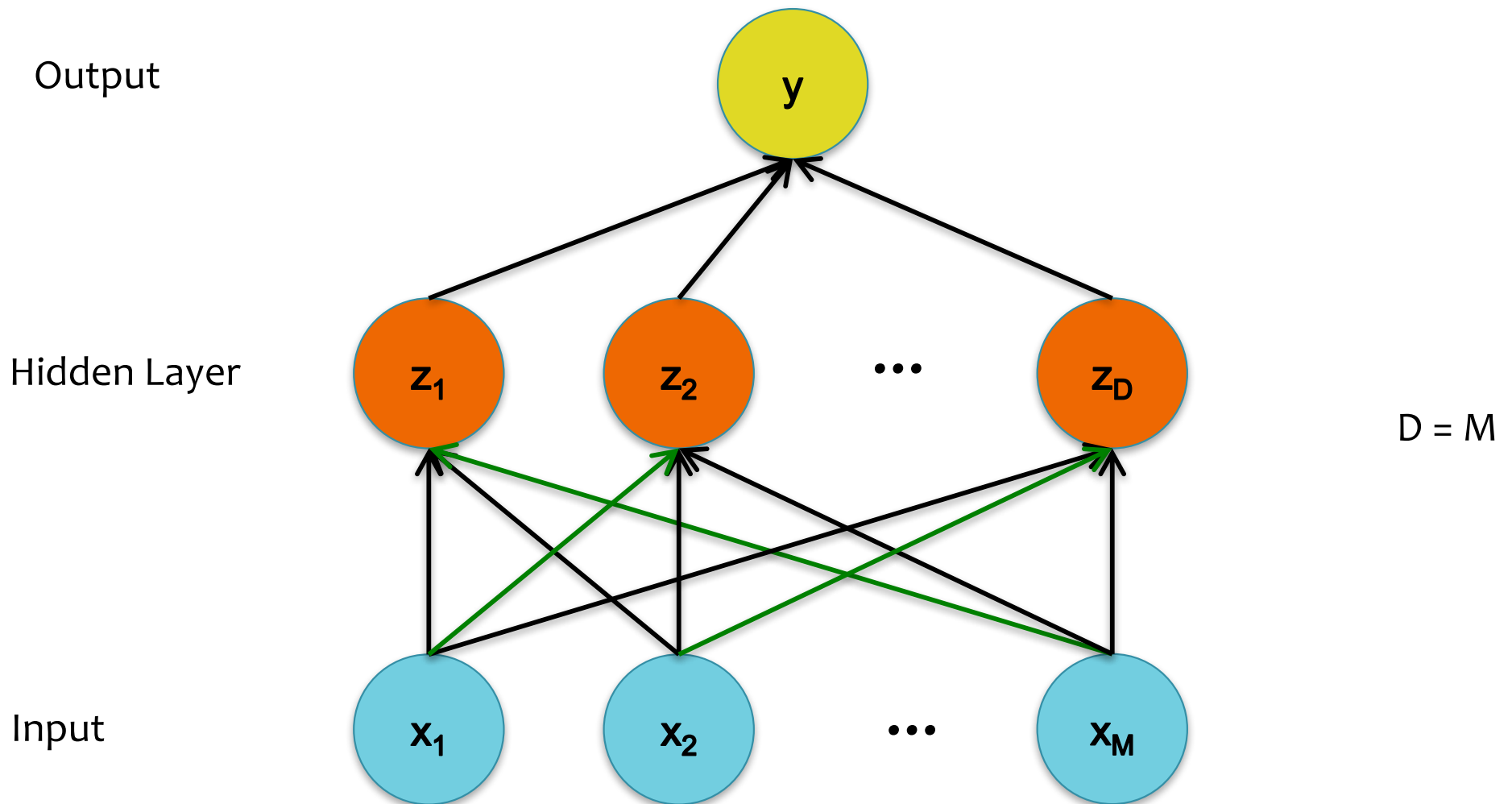
Building a Neural Net

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



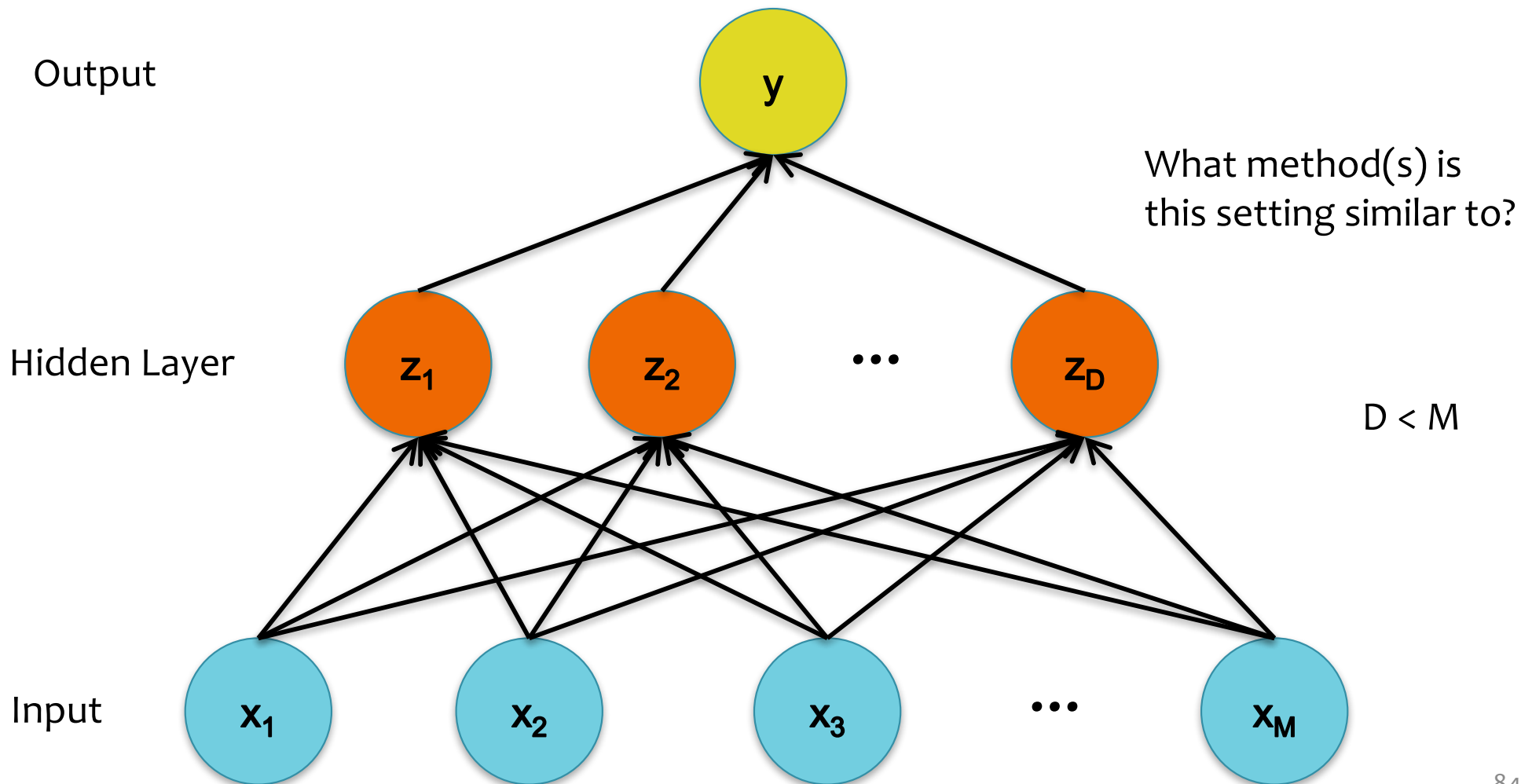
Building a Neural Net

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



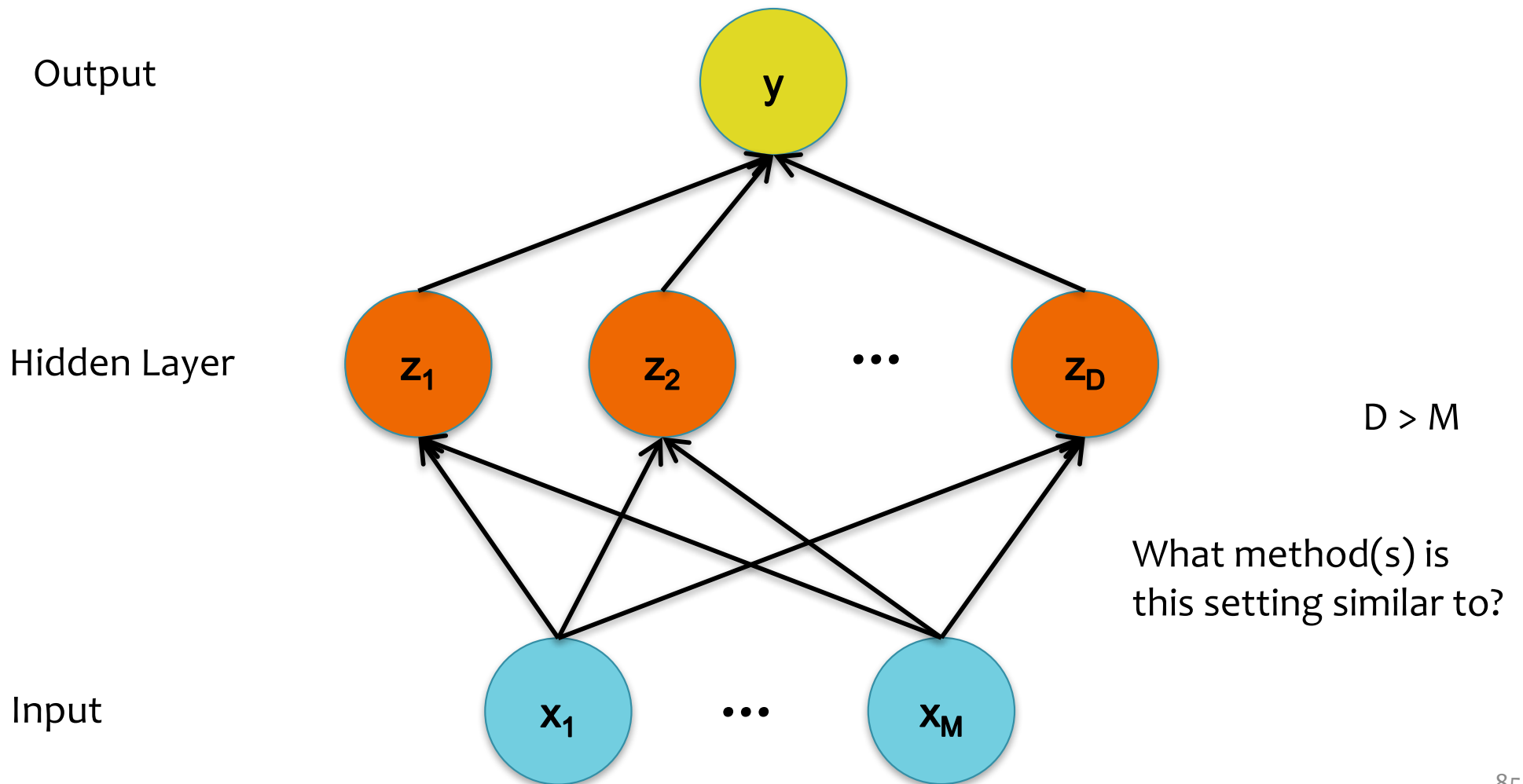
Building a Neural Net

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



Building a Neural Net

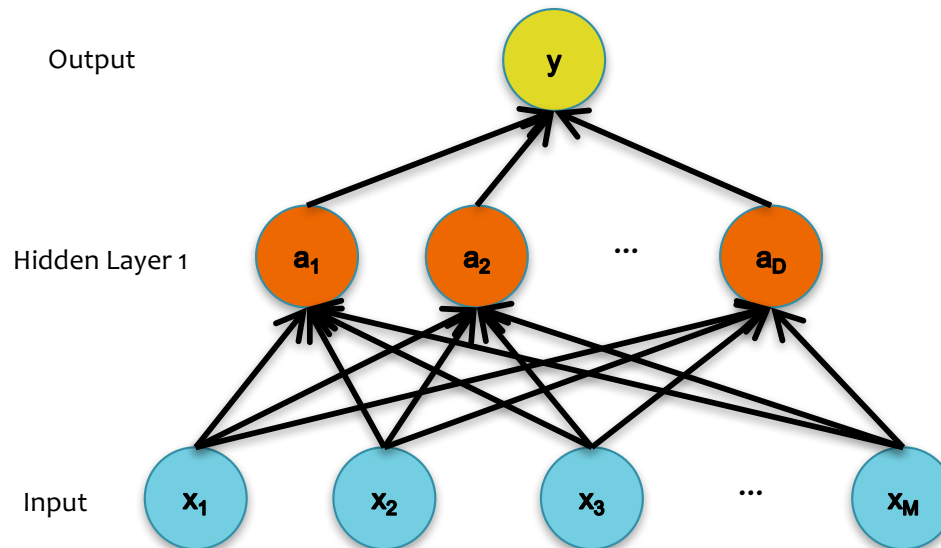
Q: How many hidden units, D , should we use?



Decision Functions

Deeper Networks

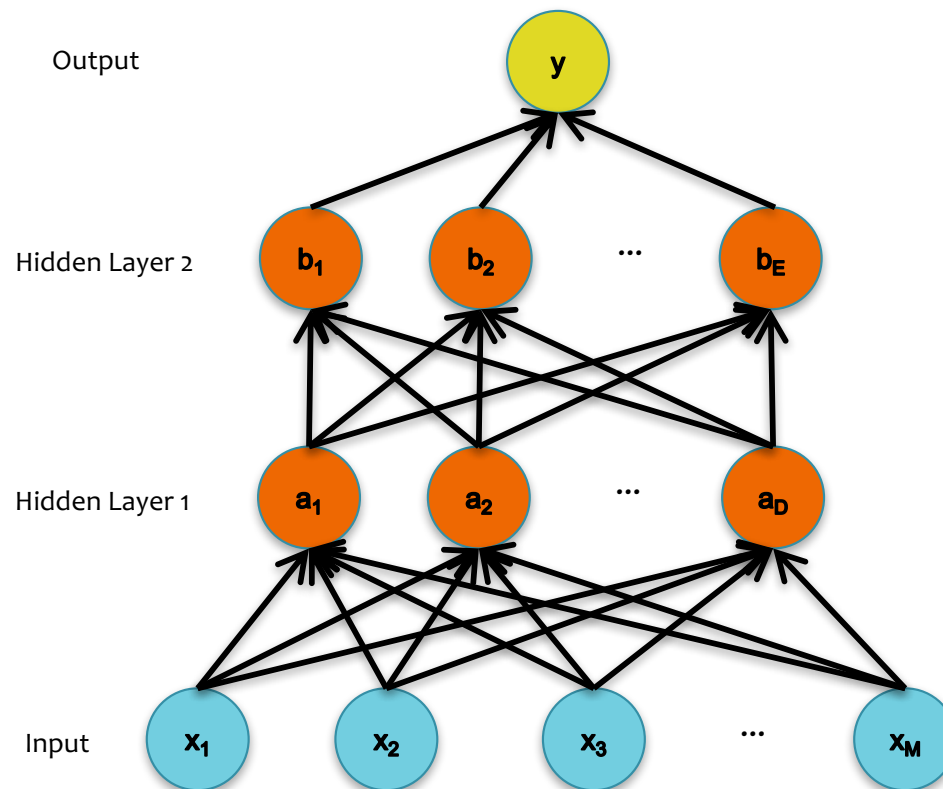
Q: How many layers should we use?



Decision Functions

Deeper Networks

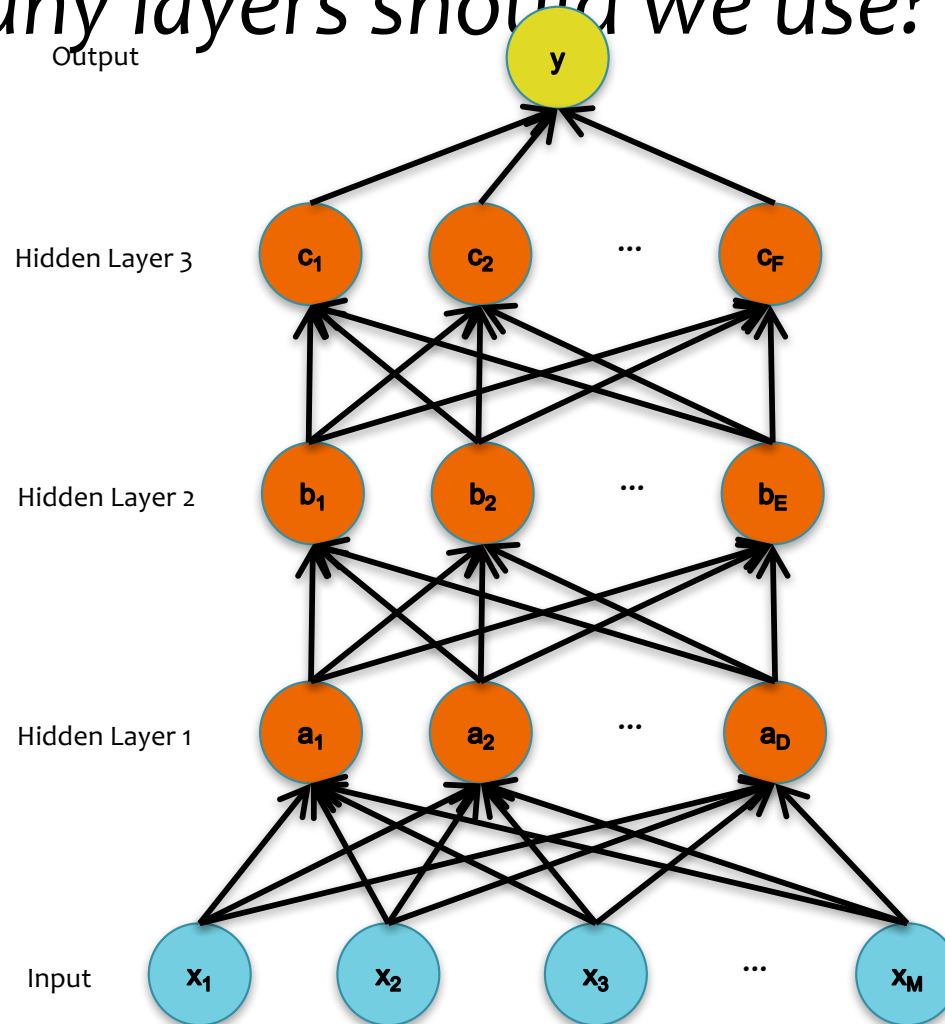
Q: How many layers should we use?



Decision Functions

Deeper Networks

Q: How many layers should we use?



Decision Functions Deeper Networks

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:**

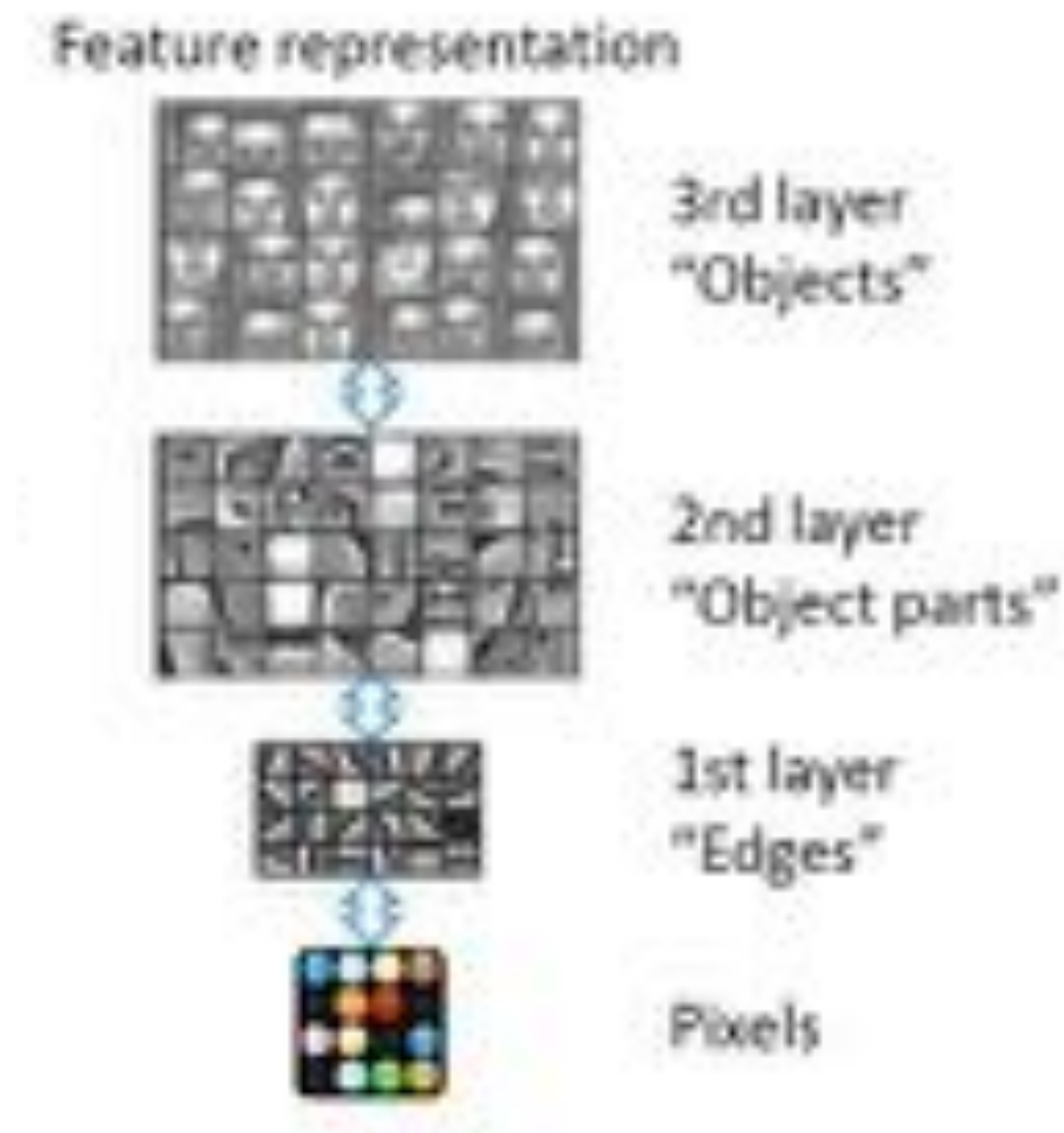
- Before 2006: “Deep networks (e.g. 3 or more hidden layers) are too hard to train”
- After 2006: “Deep networks are easier to train than shallow networks (e.g. 2 or fewer layers) for many problems”

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

Decision Functions

Different Levels of Abstraction

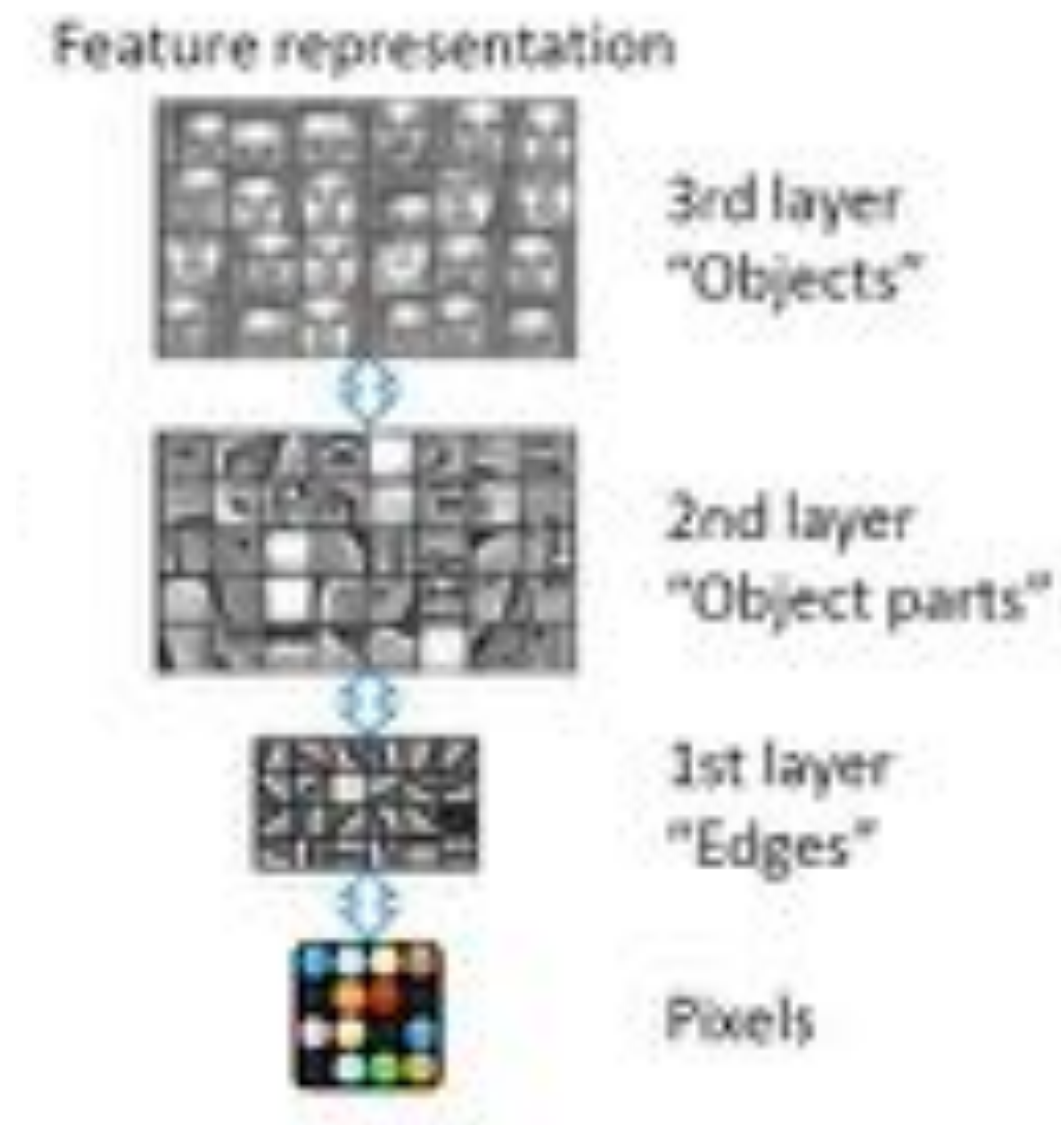
- We don't know the “right” levels of abstraction
- So let the model figure it out!



Decision Functions Different Levels of Abstraction

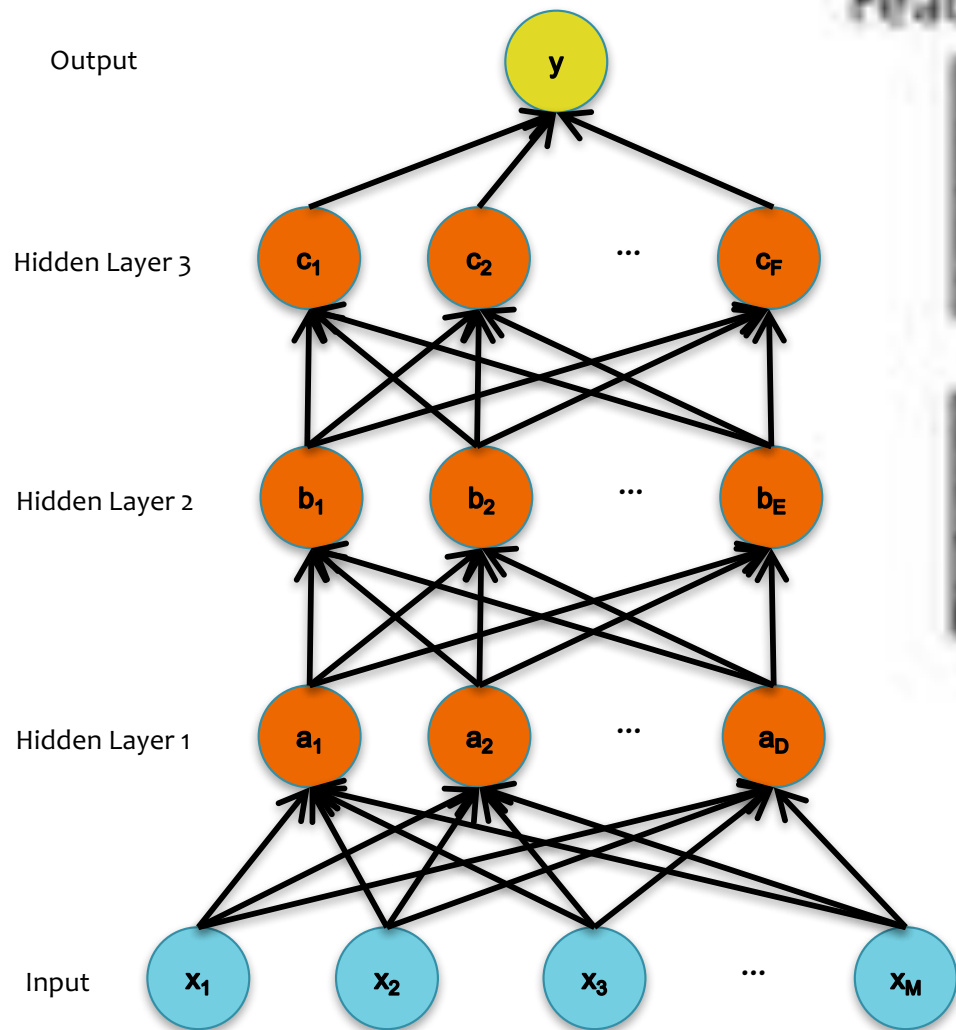
Face Recognition:

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



Decision Functions

Different Levels of Abstraction



Feature representation



3rd layer
"Objects"



2nd layer
"Object parts"



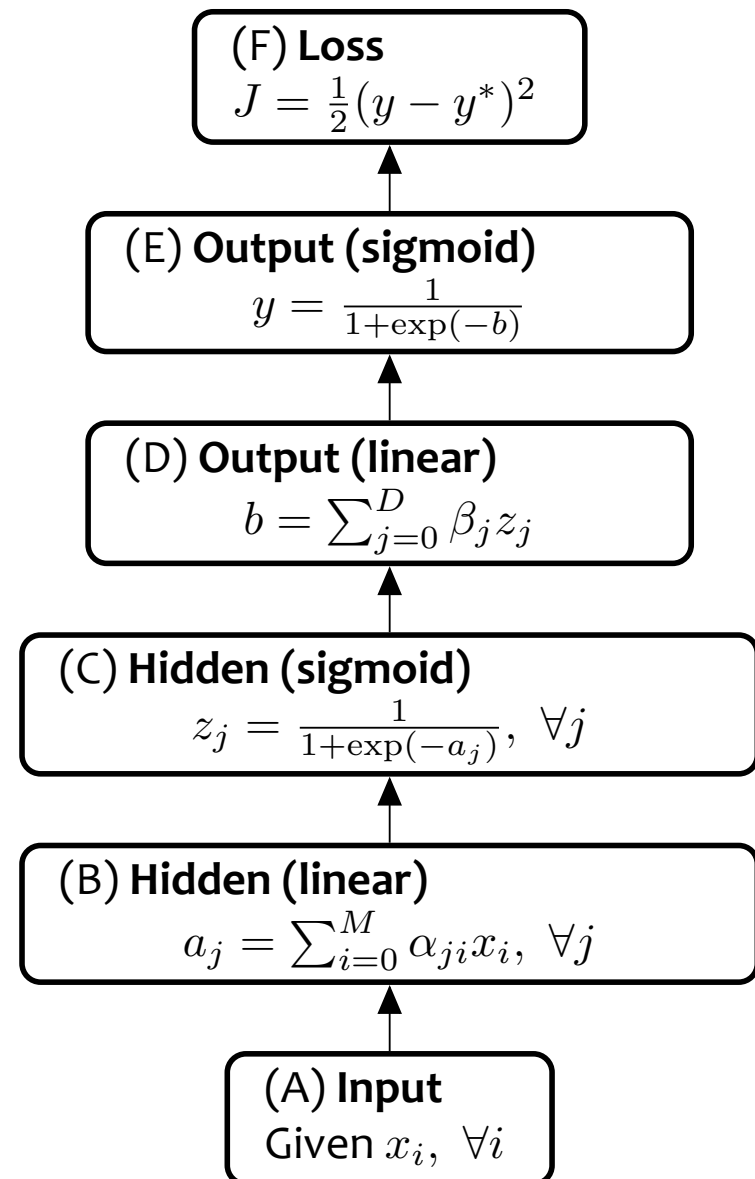
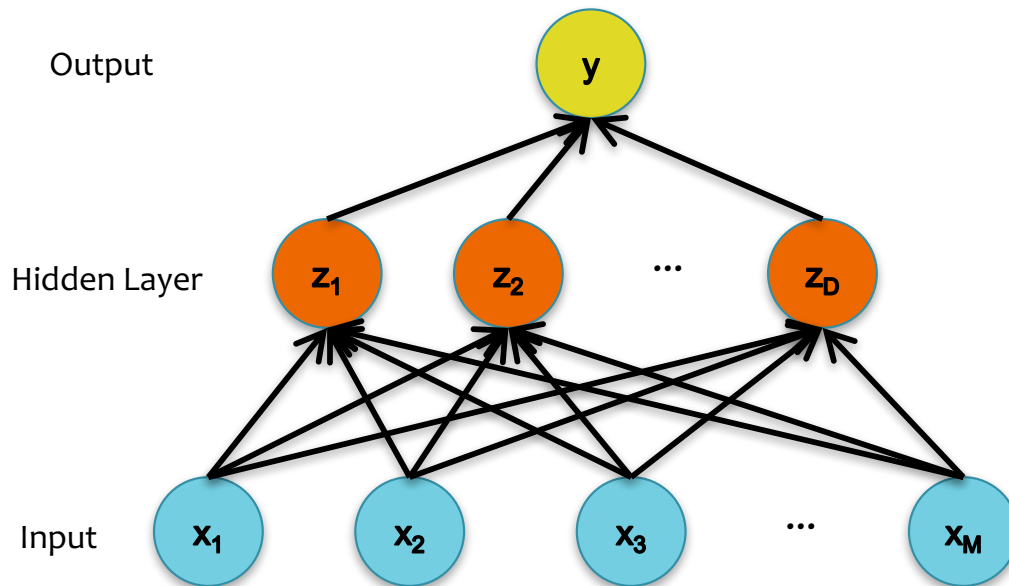
1st layer
"Edges"



Pixels

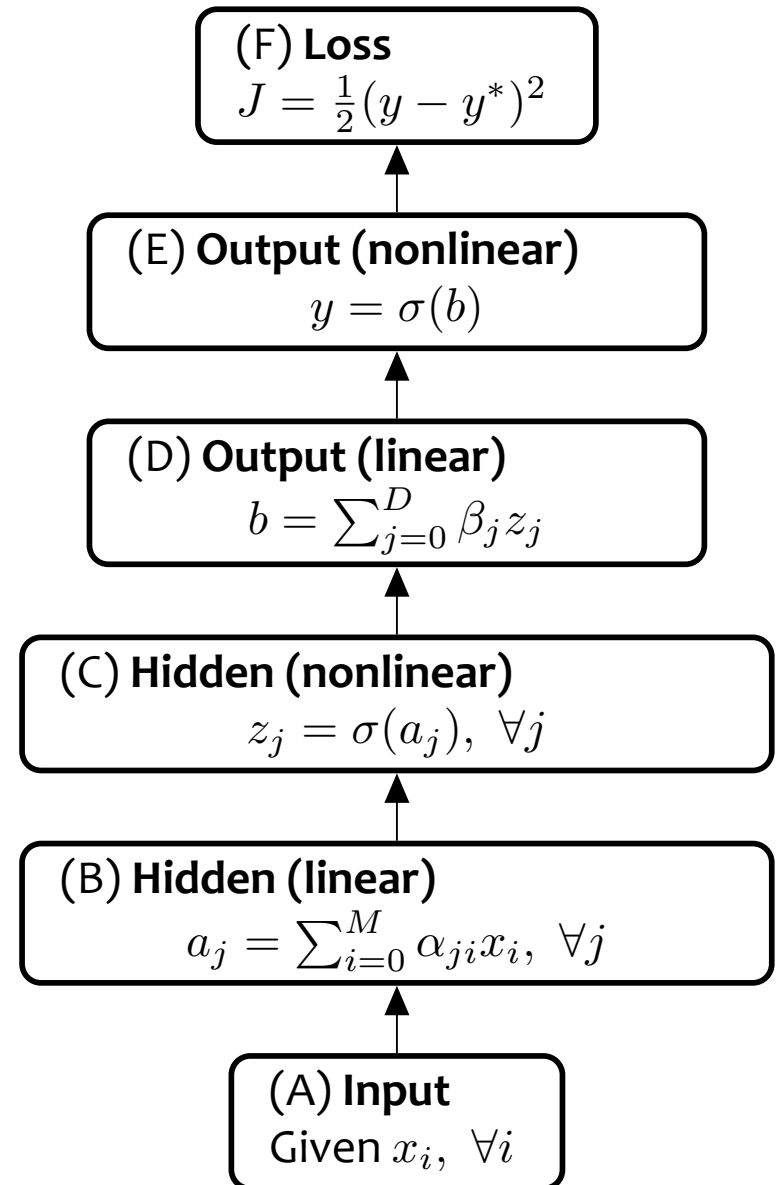
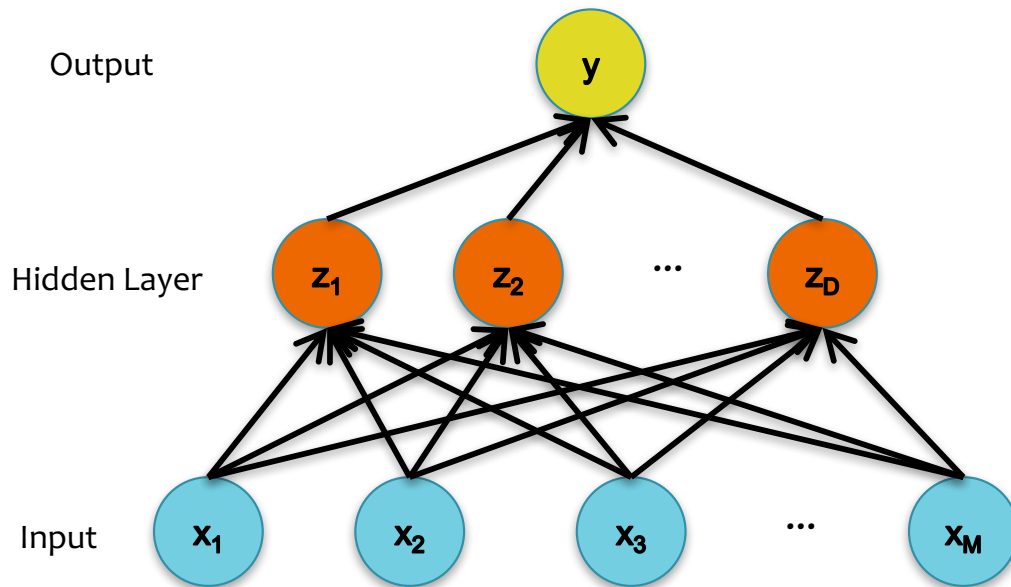
Activation Functions

Neural Network with sigmoid
activation functions



Activation Functions

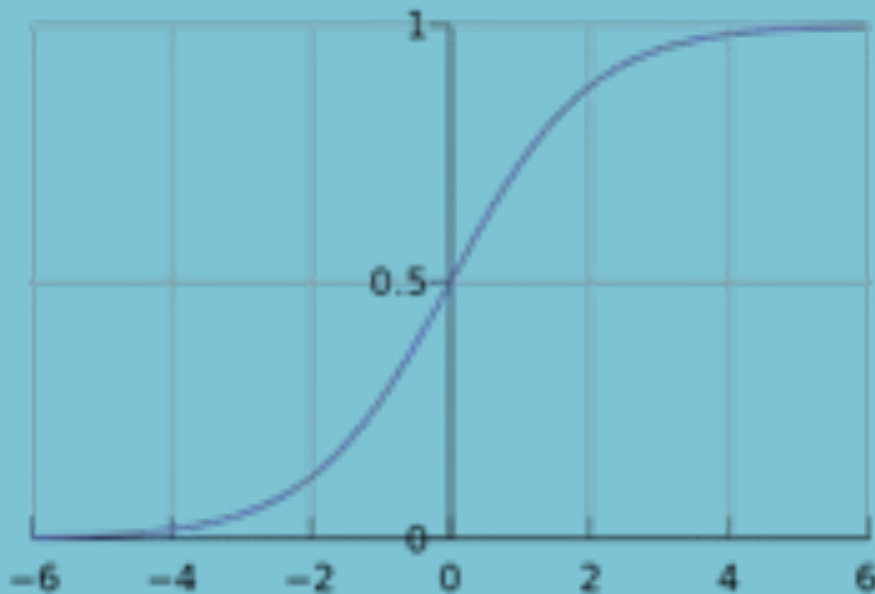
Neural Network with arbitrary nonlinear activation functions



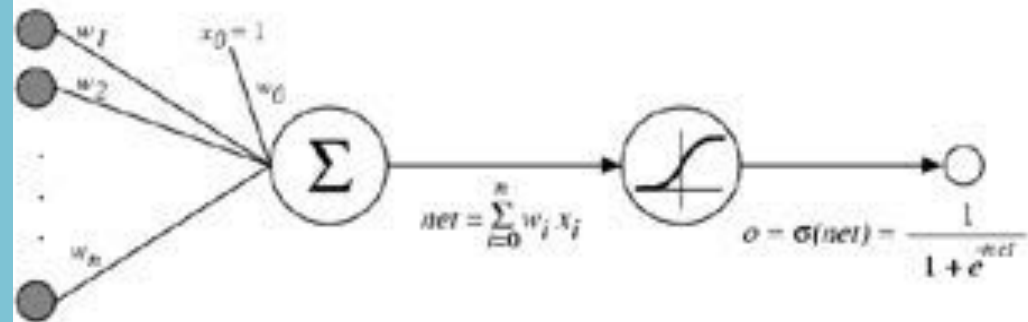
Activation Functions

Sigmoid / Logistic Function

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

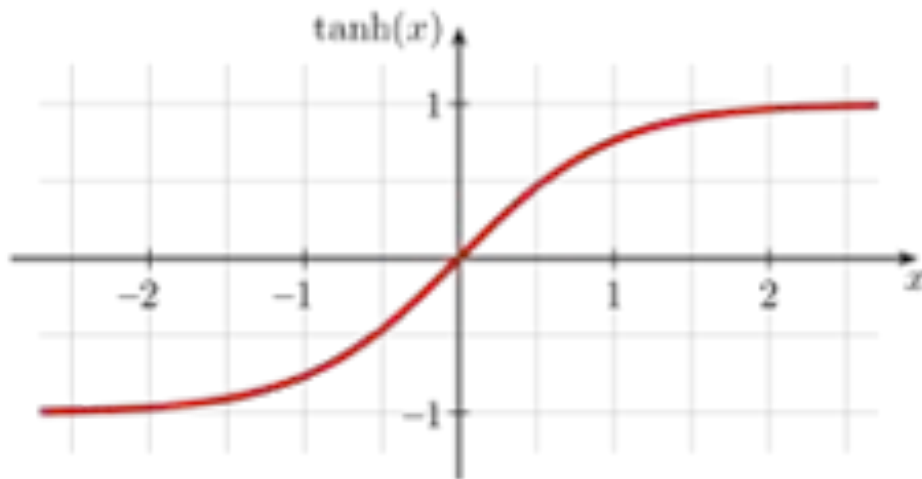


So far, we've assumed that the activation function (nonlinearity) is always the sigmoid function...



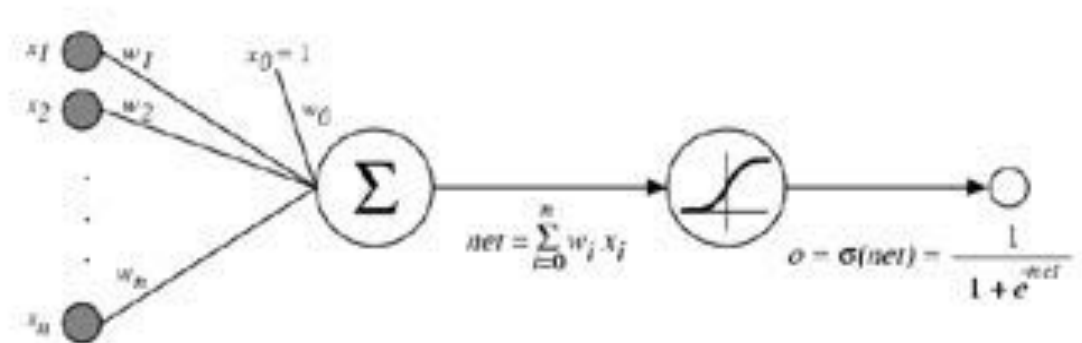
Activation Functions

- A new change: modifying the nonlinearity
 - The logistic is not widely used in modern ANNs



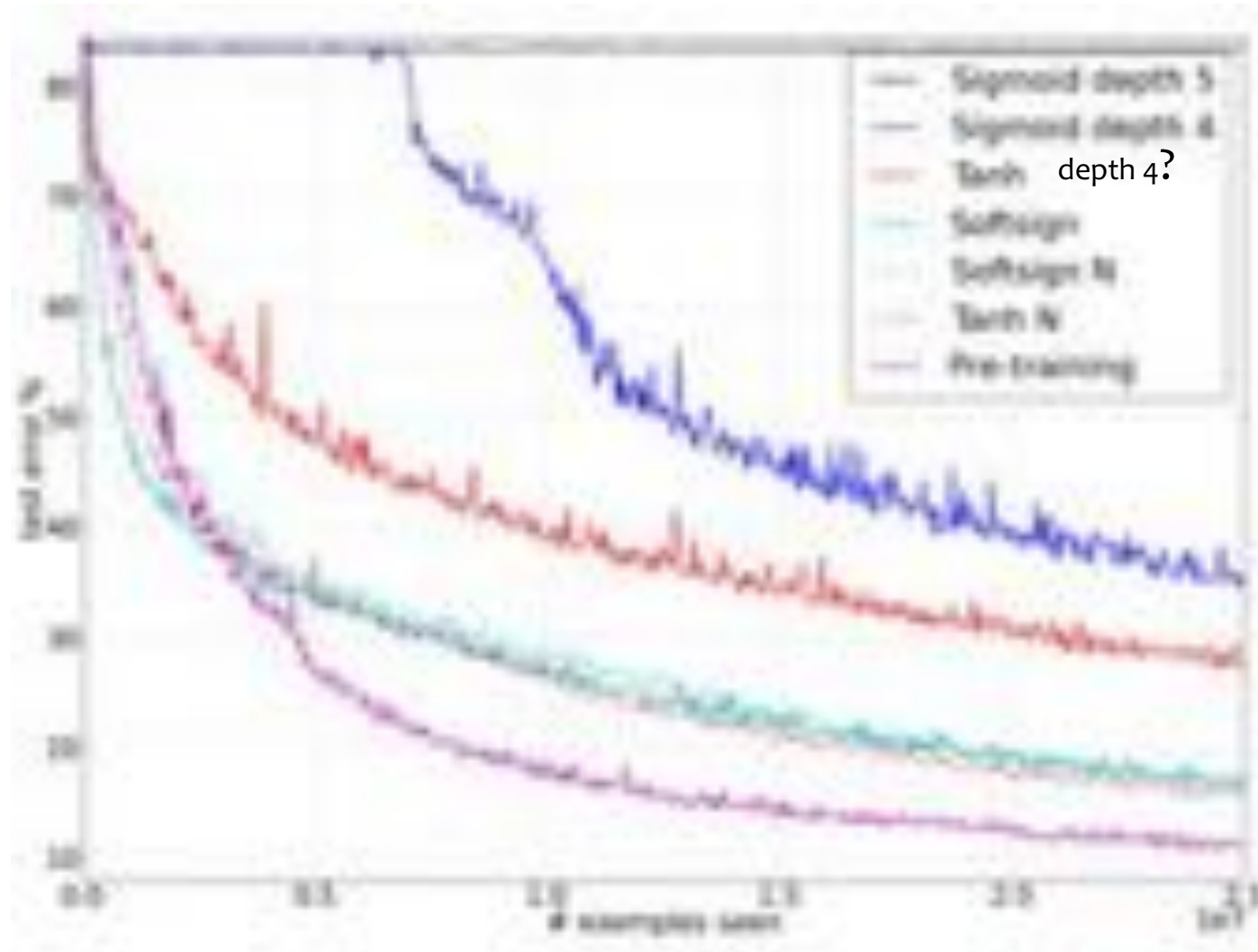
Alternate 1:
tanh

Like logistic function but
shifted to range $[-1, +1]$



Understanding the difficulty of training deep feedforward neural networks

AI Stats 2010

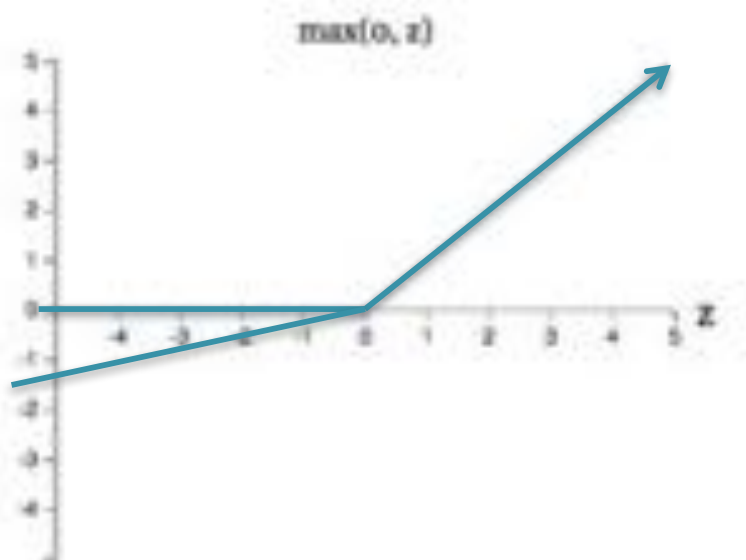


} sigmoid
vs.
tanh

Figure from Glorot & Bentio (2010)

Activation Functions

- A new change: modifying the nonlinearity
 - reLU often used in vision tasks

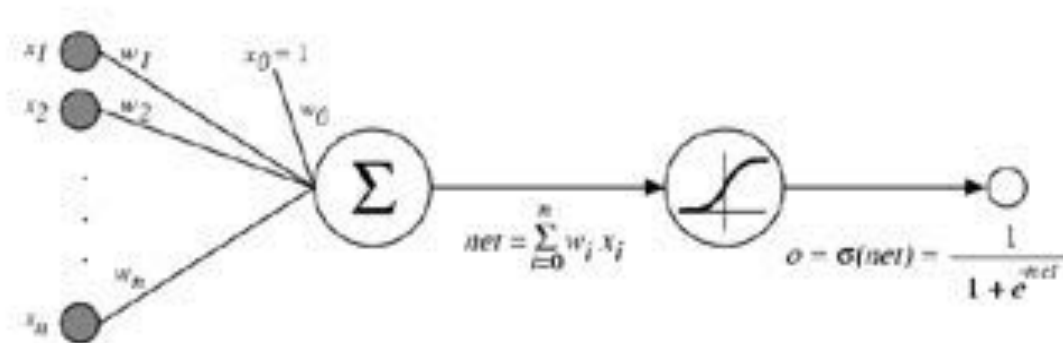


Alternate 2: rectified linear unit

Linear with a cutoff at zero

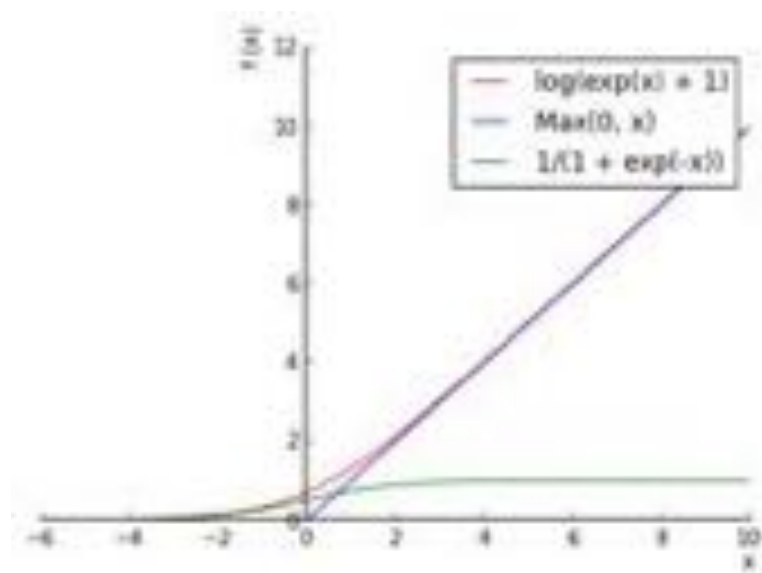
(Implementation: clip the gradient when you pass zero)

$$\max(0, w \cdot x + b).$$



Activation Functions

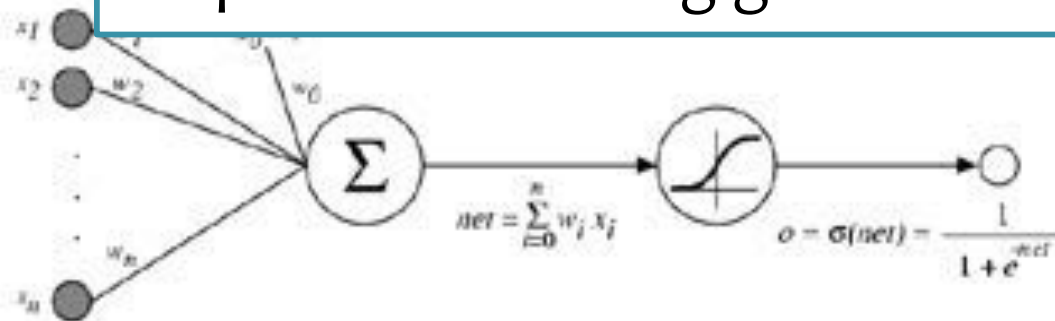
- A new change: modifying the nonlinearity
 - reLU often used in vision tasks



Alternate 2: rectified linear unit

Soft version: $\log(\exp(x)+1)$

Doesn't saturate (at one end)
Sparsifies outputs
Helps with vanishing gradient



Objective Functions for NNs

1. Quadratic Loss:

- the same objective as Linear Regression
- i.e. mean squared error

2. Cross-Entropy:

- the same objective as Logistic Regression
- i.e. negative log likelihood
- This requires probabilities, so we add an additional “softmax” layer at the end of our network

Forward

Quadratic $J = \frac{1}{2}(y - y^*)^2$

Cross Entropy $J = y^* \log(y) + (1 - y^*) \log(1 - y)$

Backward

$$\frac{dJ}{dy} = y - y^*$$

$$\frac{dJ}{dy} = y^* \frac{1}{y} + (1 - y^*) \frac{1}{y - 1}$$

Objective Functions for NNs

Cross-entropy vs. Quadratic loss

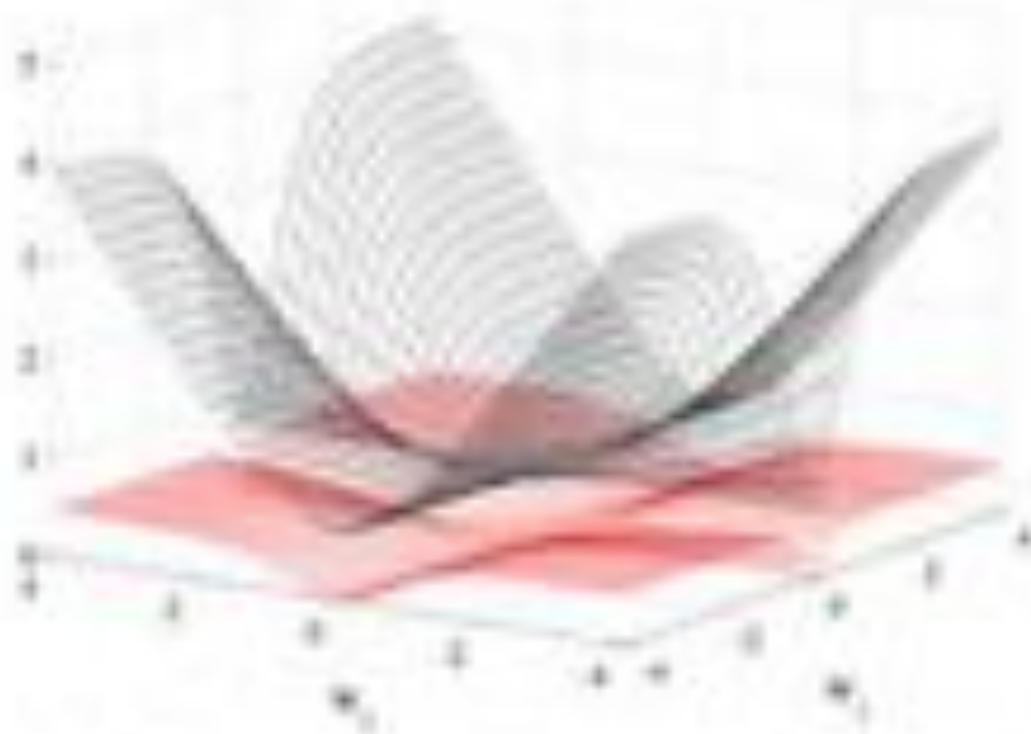
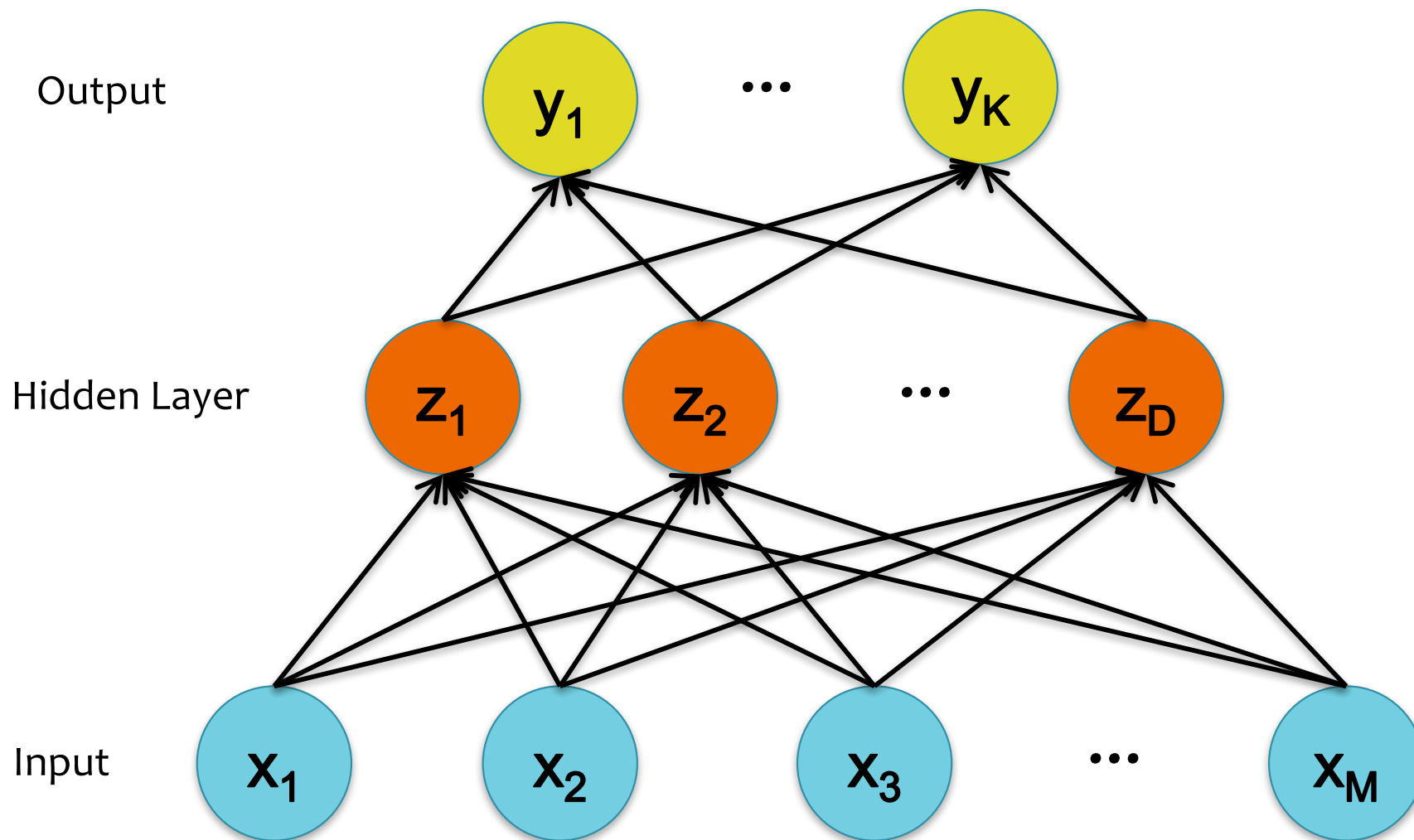


Figure 3: Cross entropy (black, surface on top) and quadratic (red, bottom surface) cost as a function of two weights (one at each layer) of a network with two layers, w_1 respectively on the first layer and w_2 on the second, output layer.

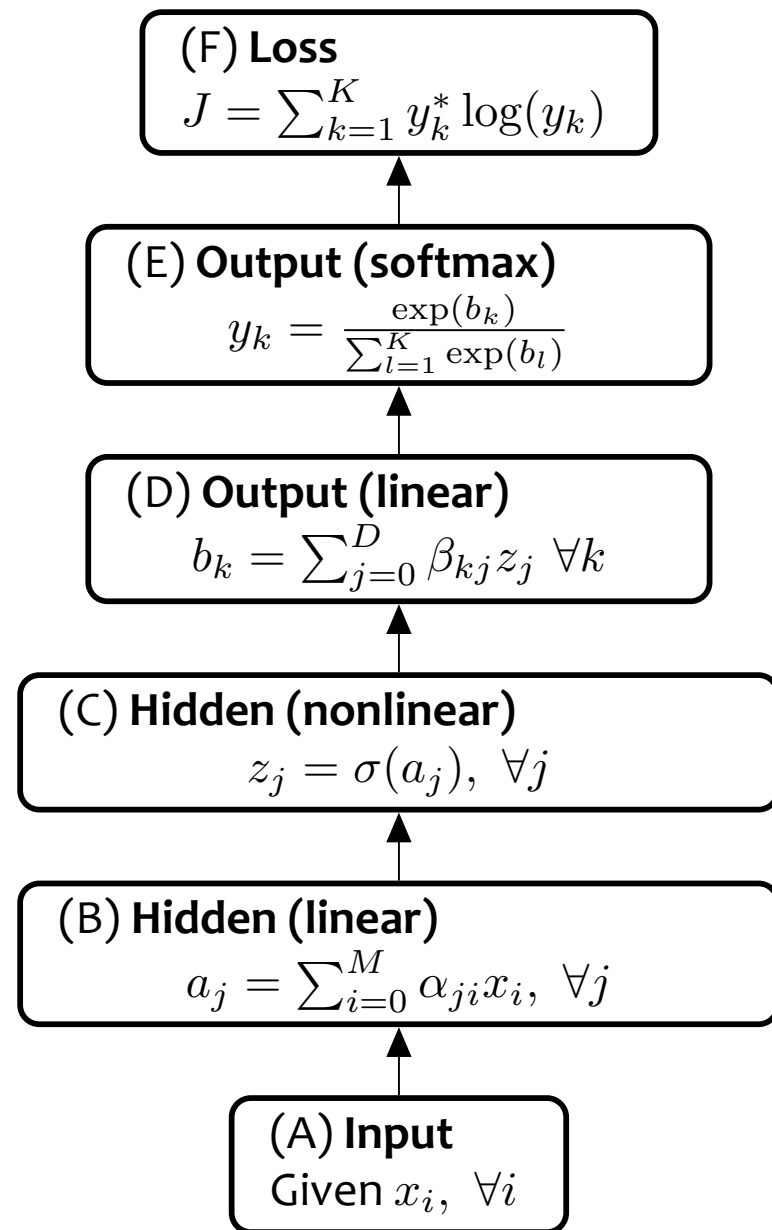
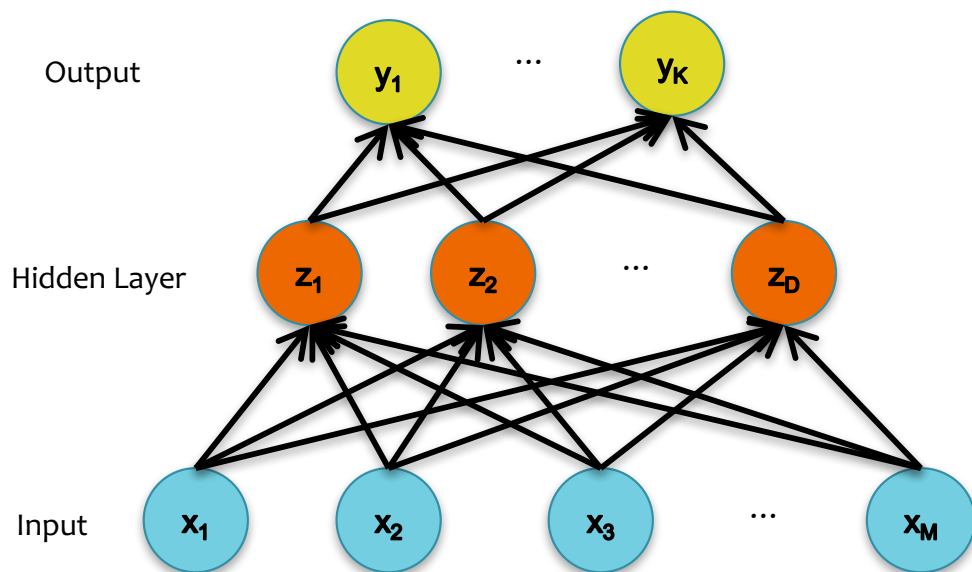
Multi-Class Output



Multi-Class Output

Softmax:

$$y_k = \frac{\exp(b_k)}{\sum_{l=1}^K \exp(b_l)}$$



Neural Networks Objectives

You should be able to...

- Explain the biological motivations for a neural network
- Combine simpler models (e.g. linear regression, binary logistic regression, multinomial logistic regression) as components to build up feed-forward neural network architectures
- Explain the reasons why a neural network can model nonlinear decision boundaries for classification
- Compare and contrast feature engineering with learning features
- Identify (some of) the options available when designing the architecture of a neural network
- Implement a feed-forward neural network