



# 10-601 Introduction to Machine Learning

Machine Learning Department  
School of Computer Science  
Carnegie Mellon University

# Neural Networks

Matt Gormley  
Lecture 11  
Feb. 21, 2018

# Reminders

- **Homework 4: Logistic Regression**
  - Out: Wed, Feb 16
  - Due: Sun, Feb 25 at 11:59pm
- **New reading on Probabilistic Learning**  
(reused later in the course for MLE/MAP)
- **Schedule changes:**
  - lecture on Friday (2/23)
  - no lecture on Monday (2/26)

# 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

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

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

# A Recipe for Machine Learning

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)

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

B

# A Recipe for

## Goals for Today's Lecture

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

2. Choose each of these.

– 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)

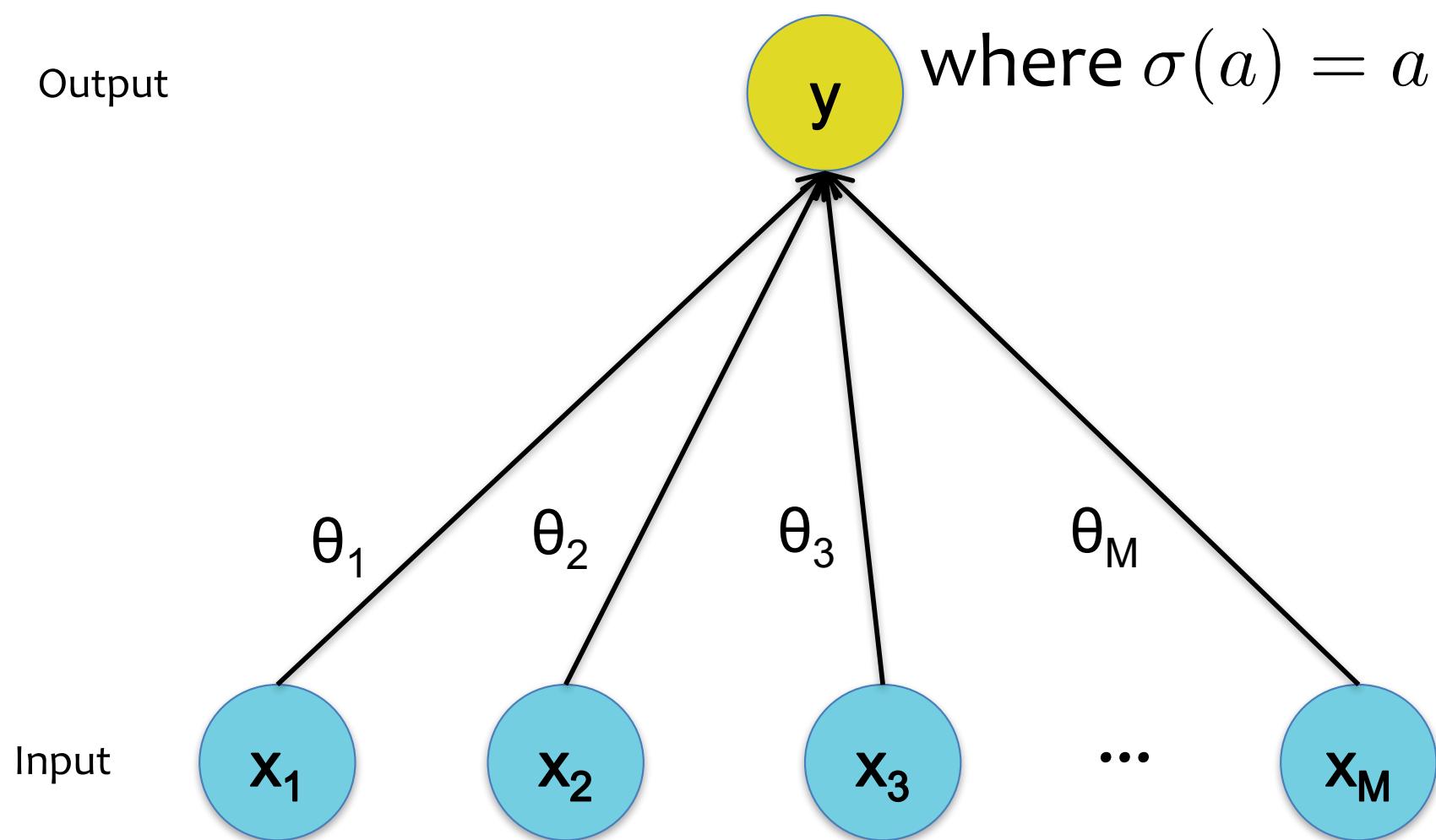
$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(x_i), y_i)$$

# Linear Regression

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

where  $\sigma(a) = a$

Output



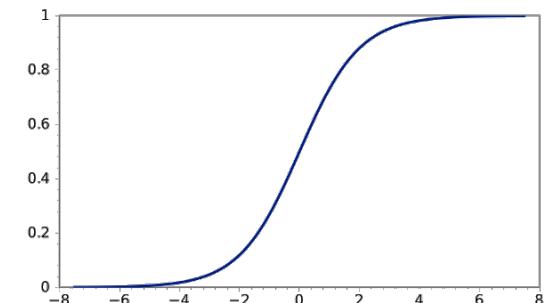
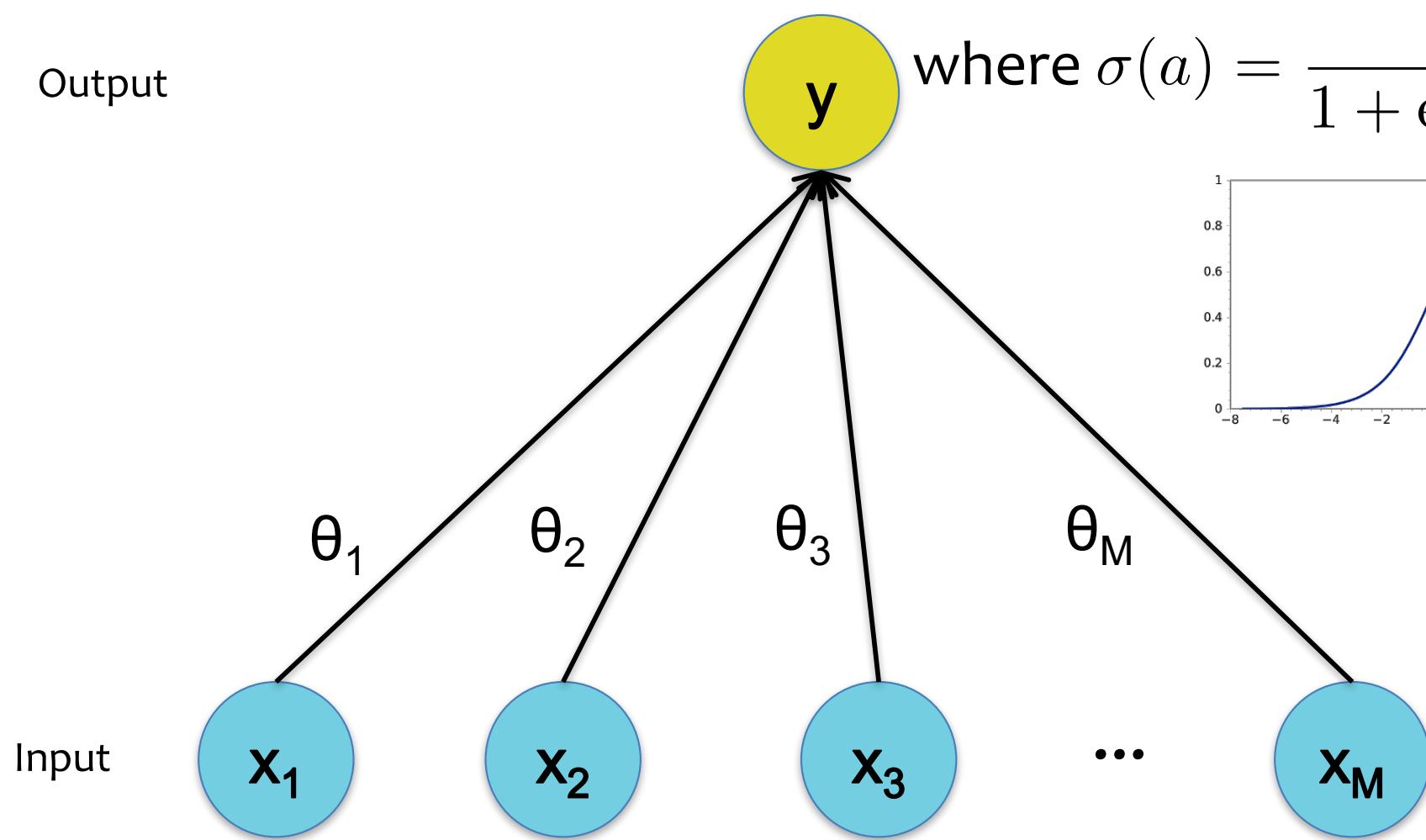
Input

# Logistic Regression

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

Output

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



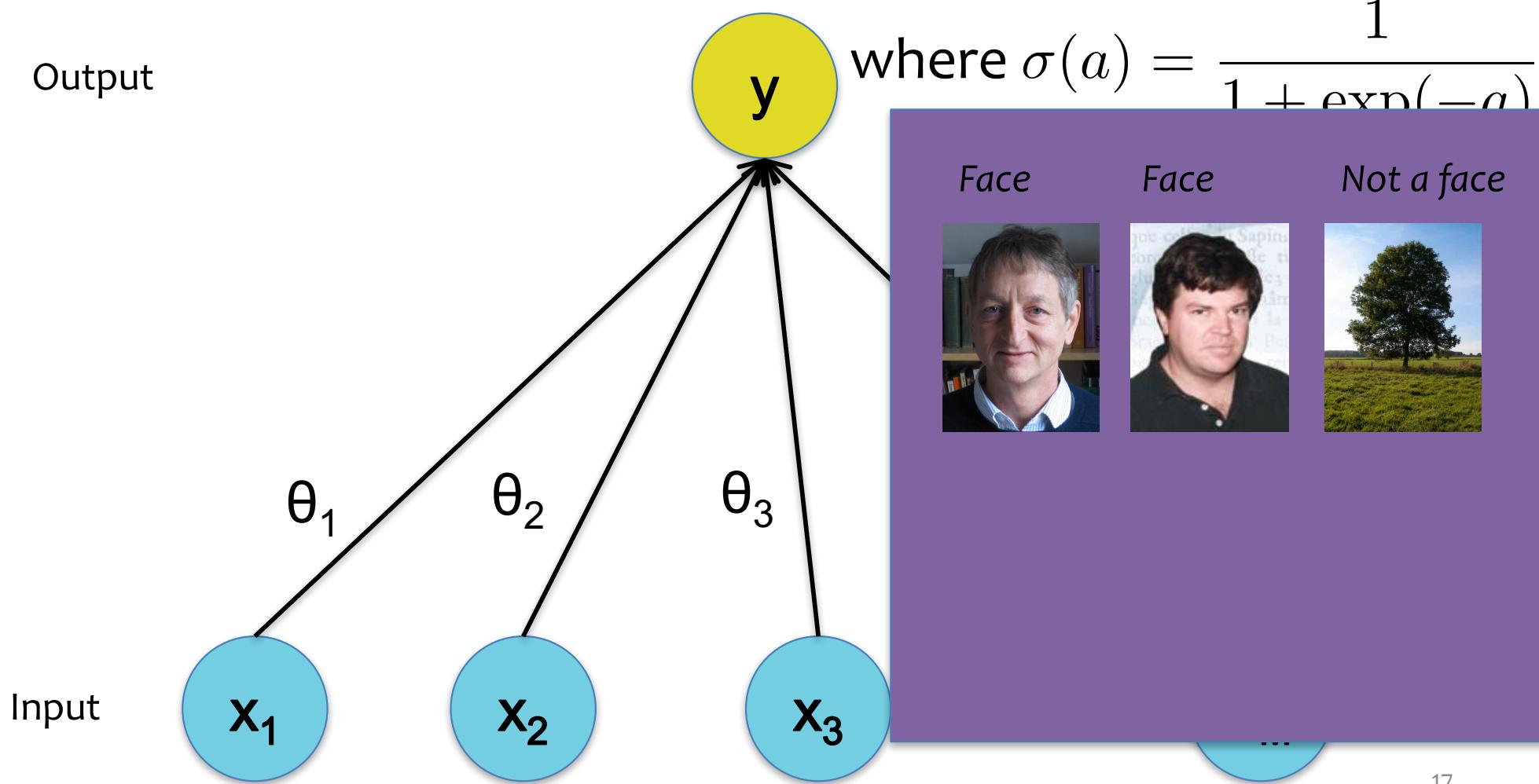
# Decision Functions

# Logistic Regression

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

Output

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



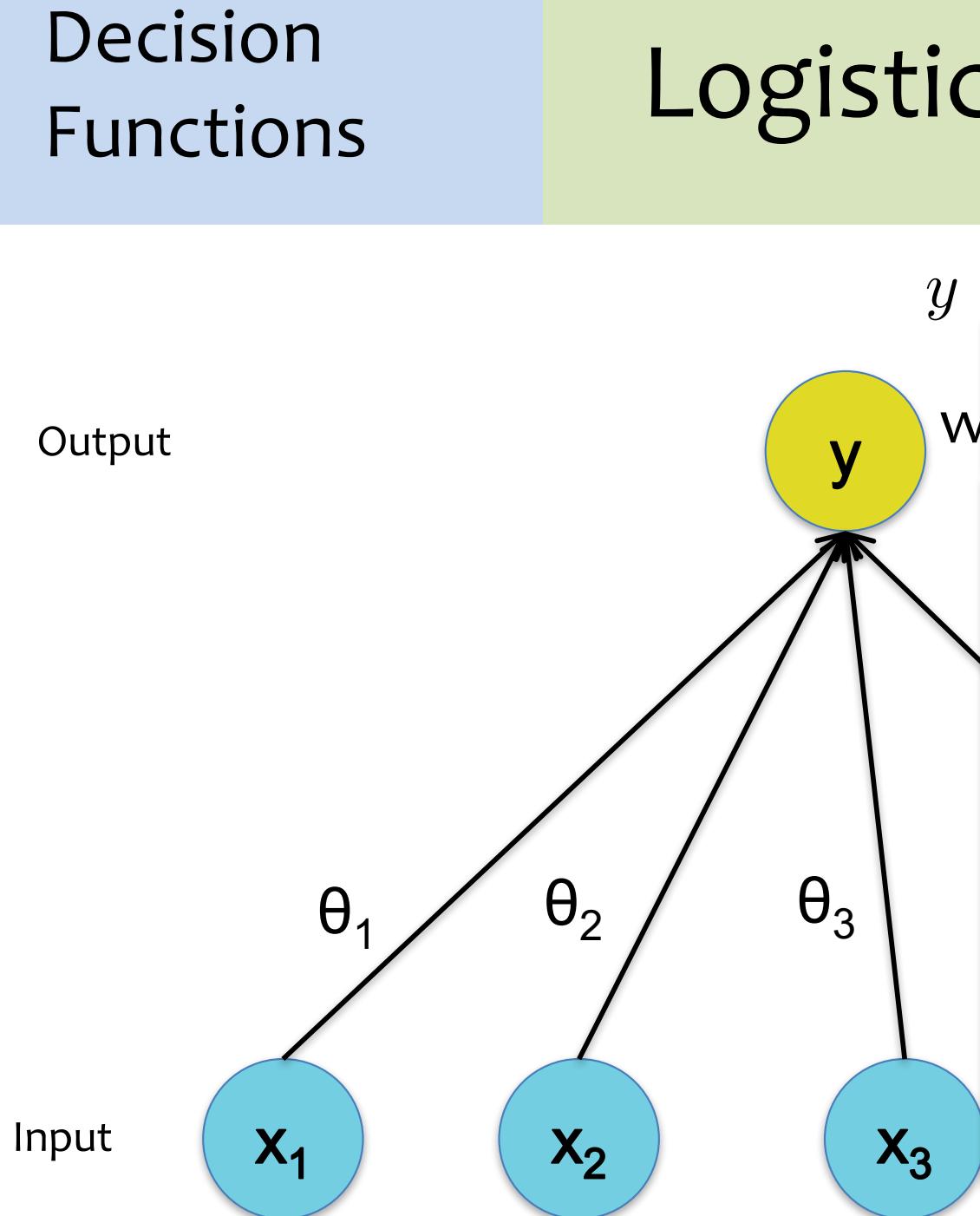
# Logistic Regression

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

Output

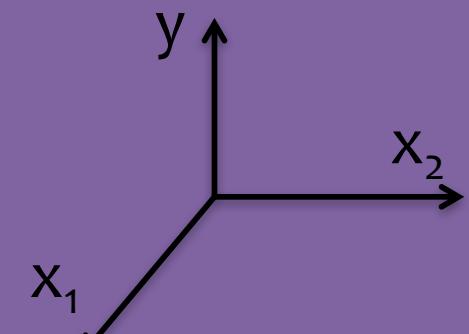
y

w



In-Class Example

1	1	0
■■	■■	■■

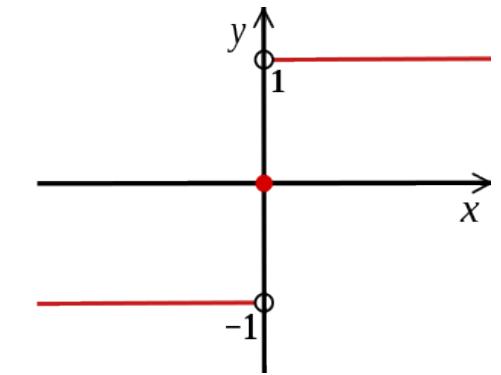
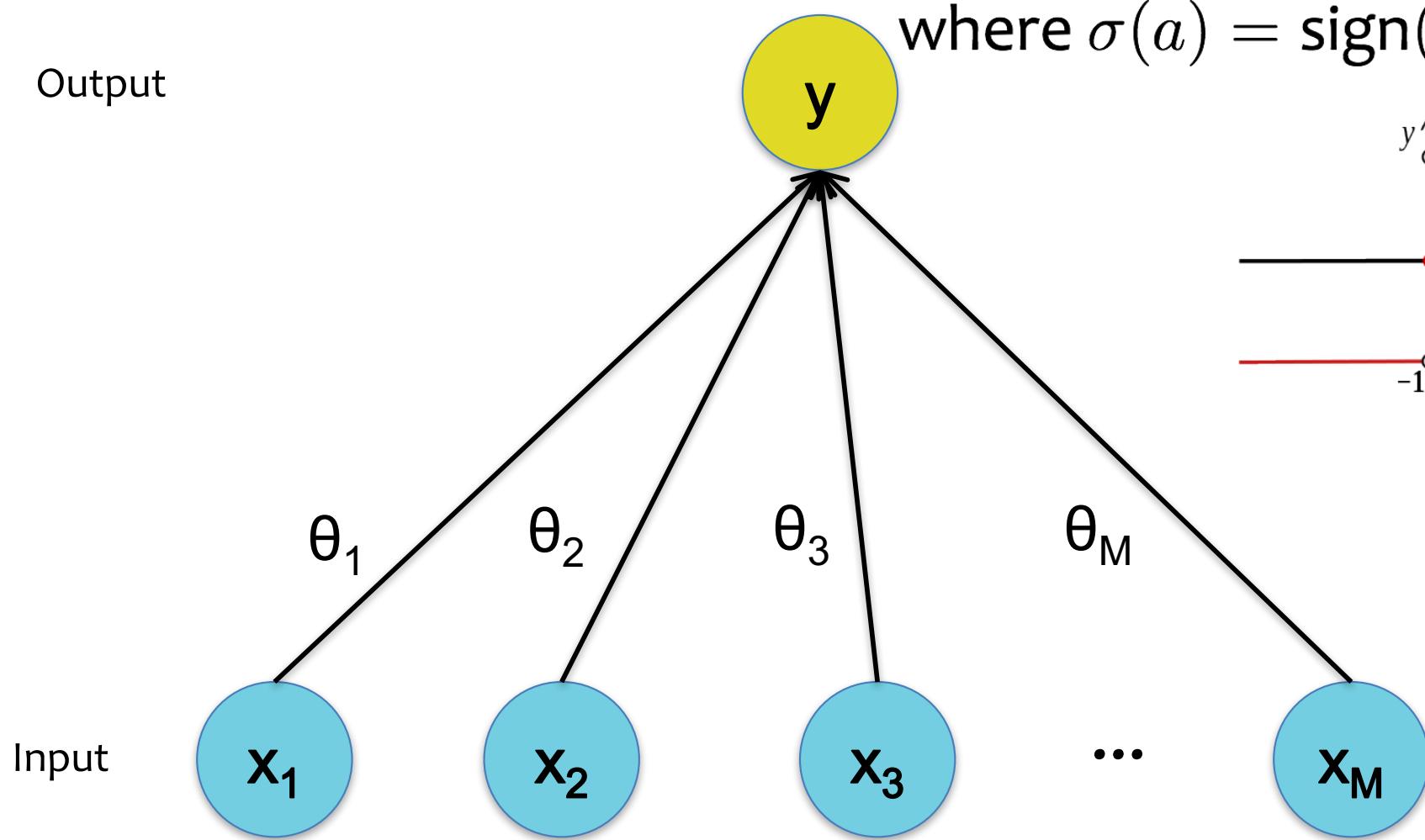


# Perceptron

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

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

Output



# From Biological to Artificial

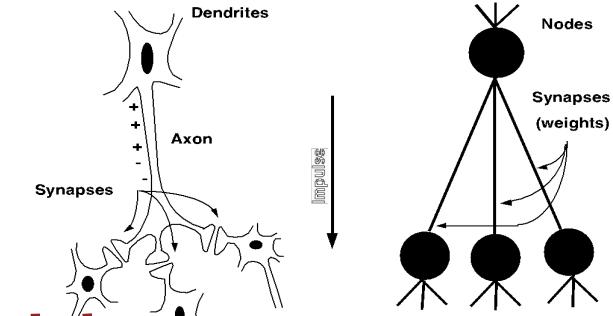
The motivation for Artificial Neural Networks comes from biology...

## 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

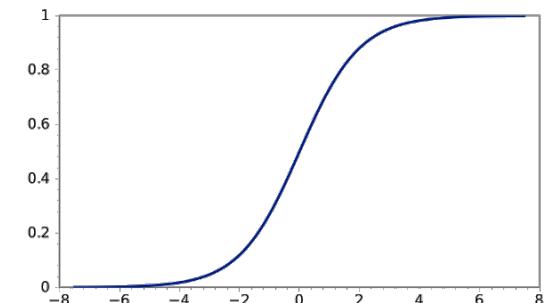
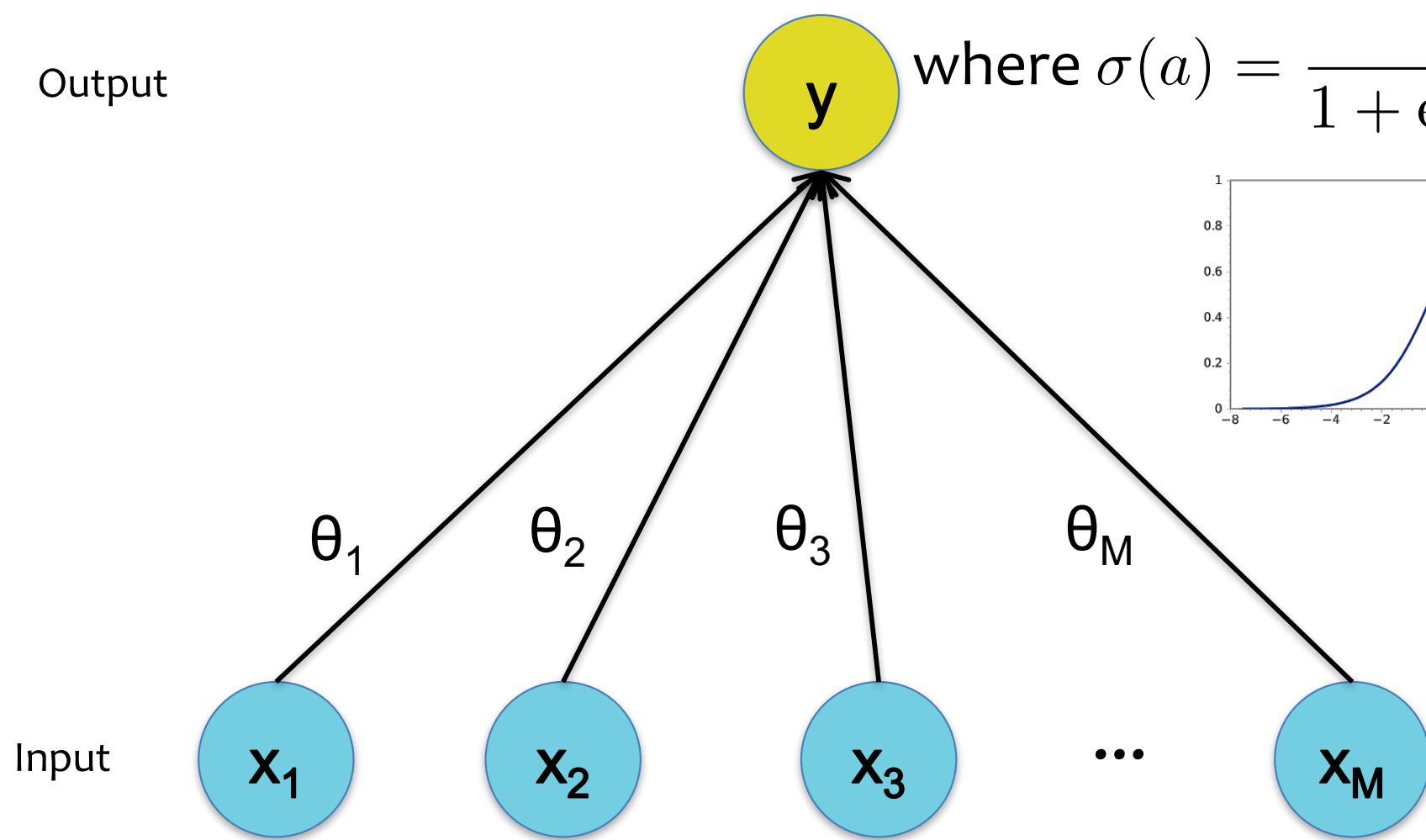
- Many neuron-like threshold switching units
- Many weighted interconnections among units
- Highly parallel, distributed processes

# Logistic Regression

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

Output

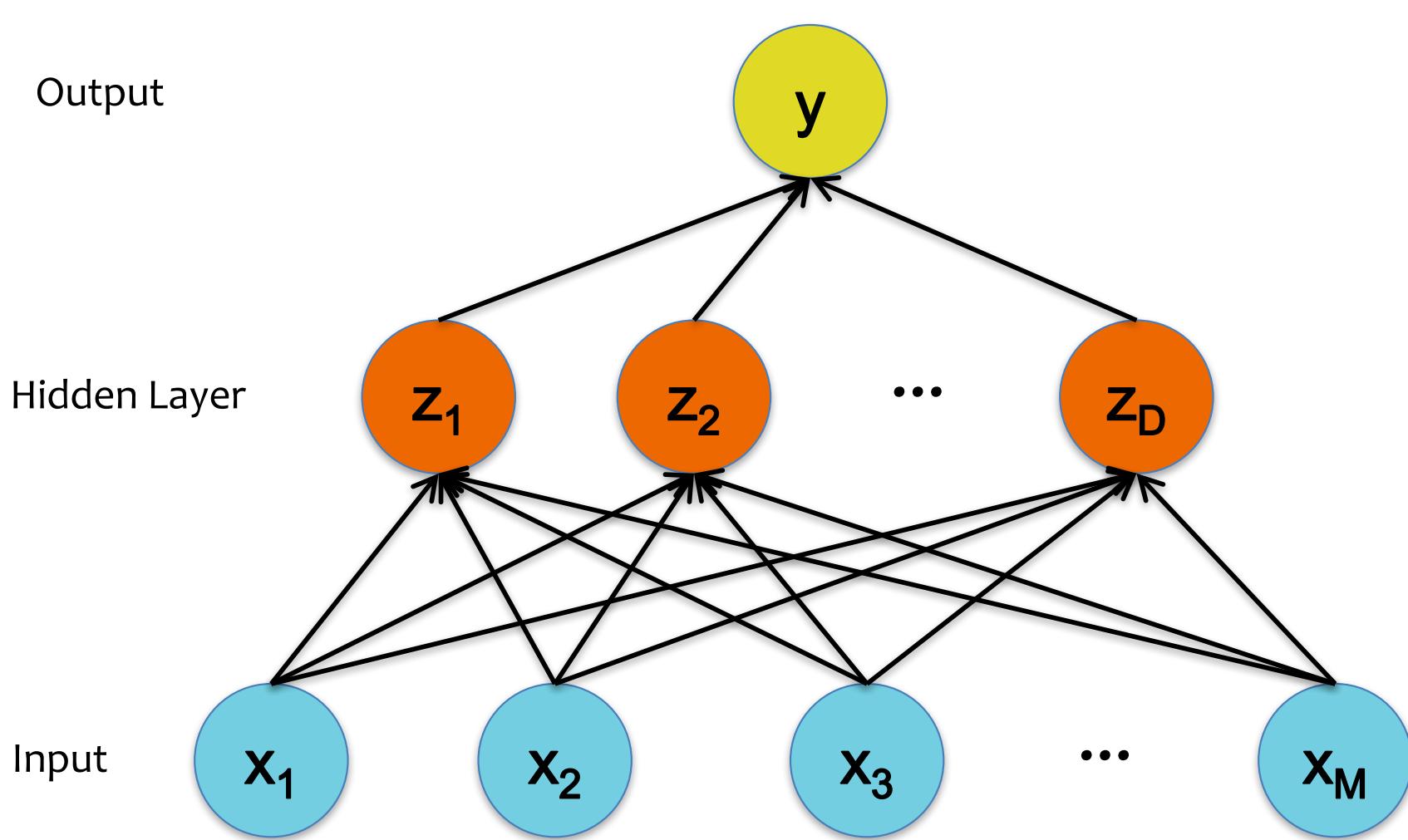
$$\text{where } \sigma(a) = \frac{1}{1 + \exp(-a)}$$



# Neural Networks

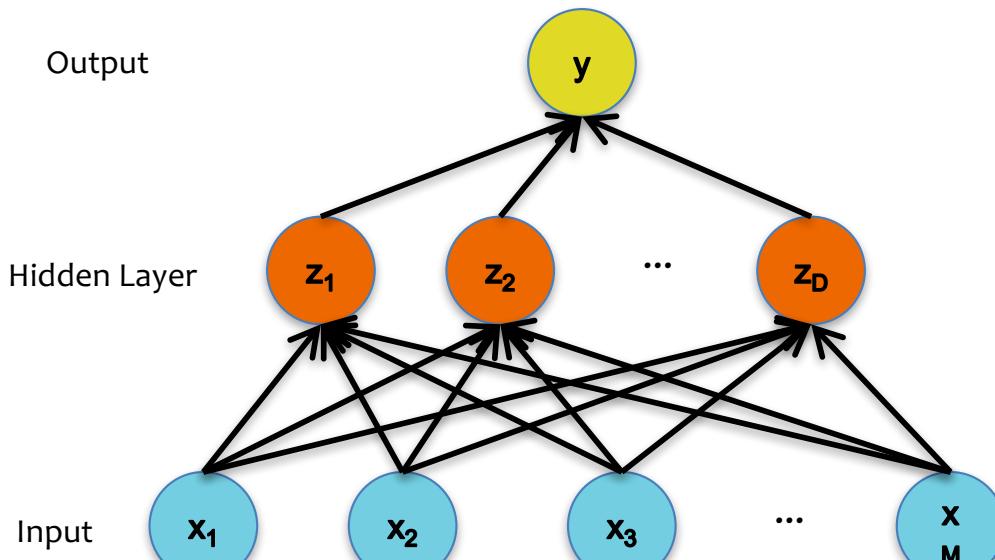
## *Chalkboard*

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



# Decision Functions

# Neural Network



(E) Output (sigmoid)

$$y = \frac{1}{1+\exp(-b)}$$

(D) Output (linear)

$$b = \sum_{j=0}^D \beta_j z_j$$

(C) Hidden (sigmoid)

$$z_j = \frac{1}{1+\exp(-a_j)}, \forall j$$

(B) Hidden (linear)

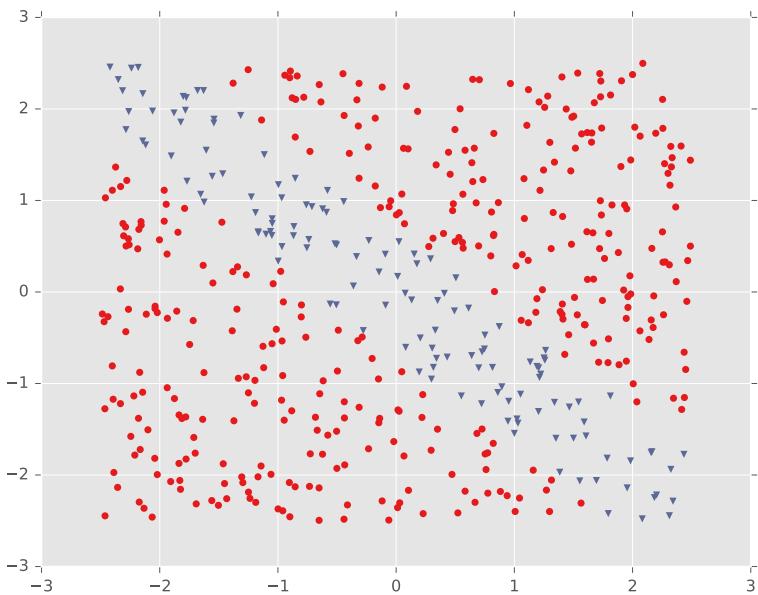
$$a_j = \sum_{i=0}^M \alpha_{ji} x_i, \forall j$$

(A) Input

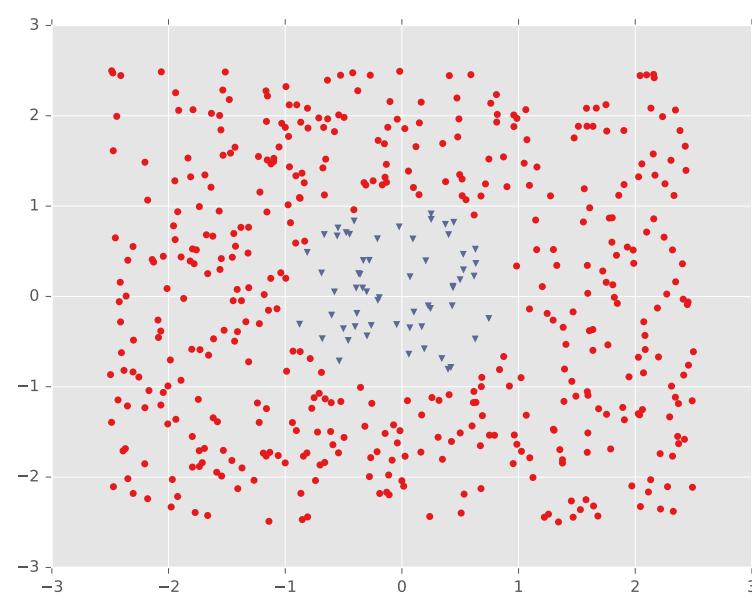
Given  $x_i, \forall i$

# **DECISION BOUNDARY EXAMPLES**

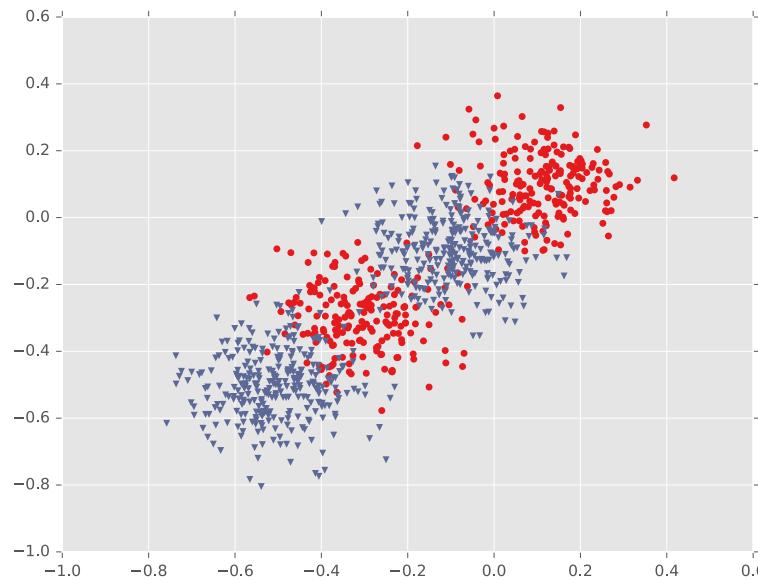
## Example #1: Diagonal Band



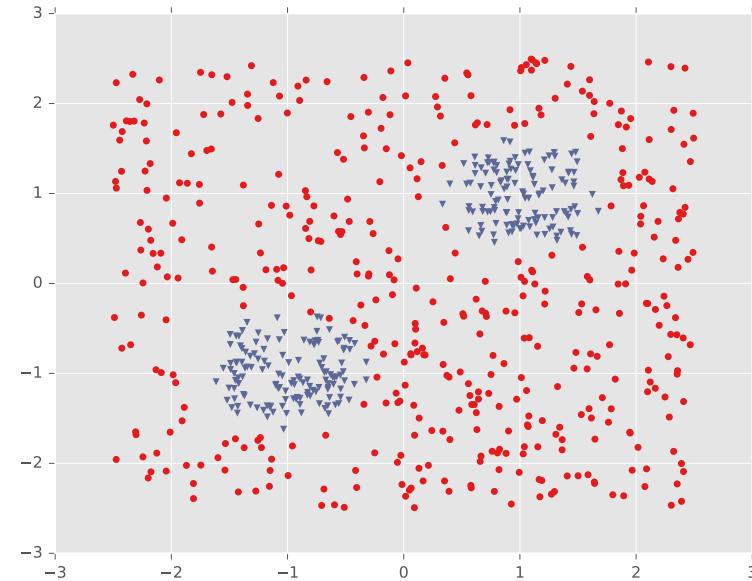
## Example #2: One Pocket



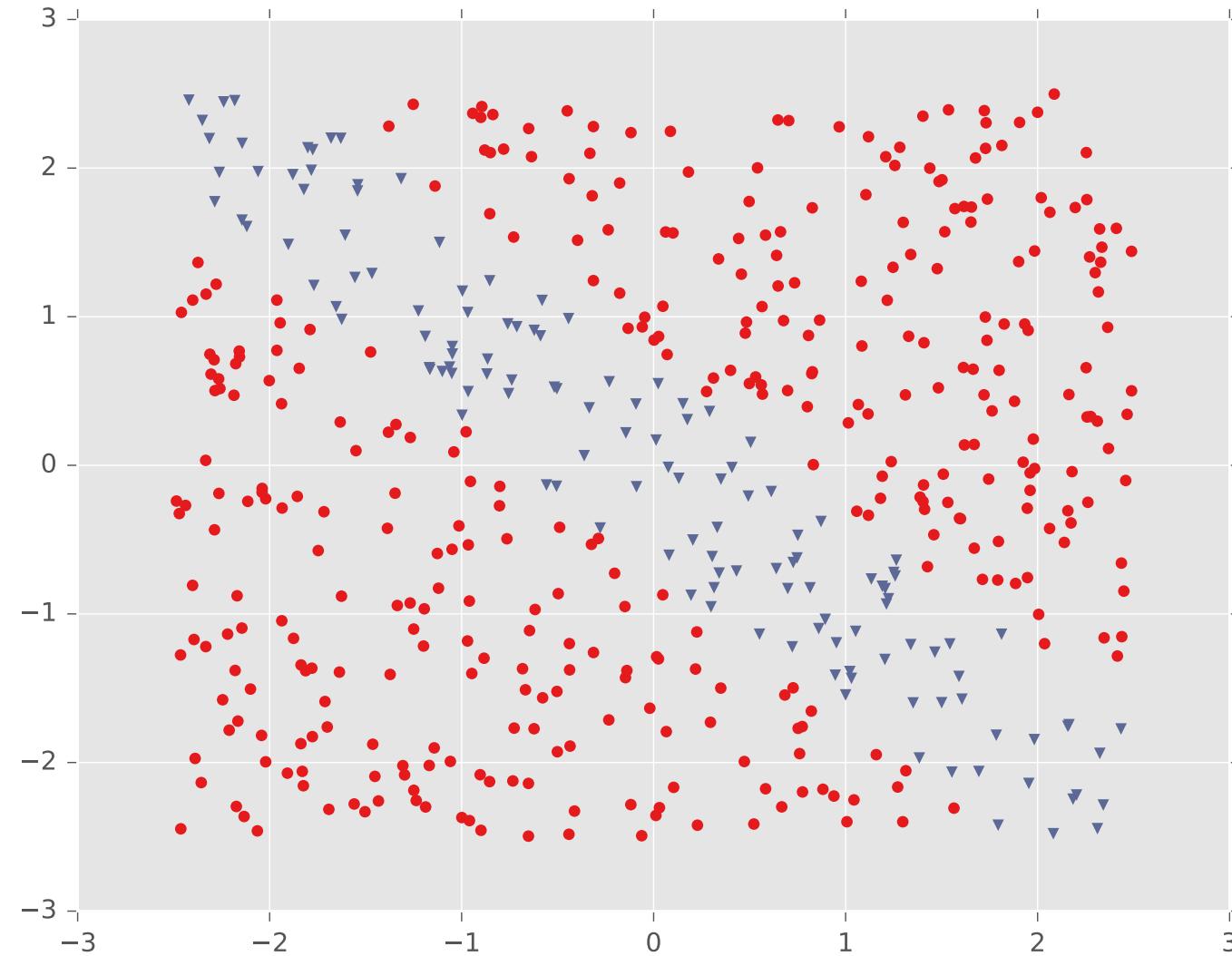
## Example #3: Four Gaussians



## Example #4: Two Pockets



# Example #1: Diagonal Band

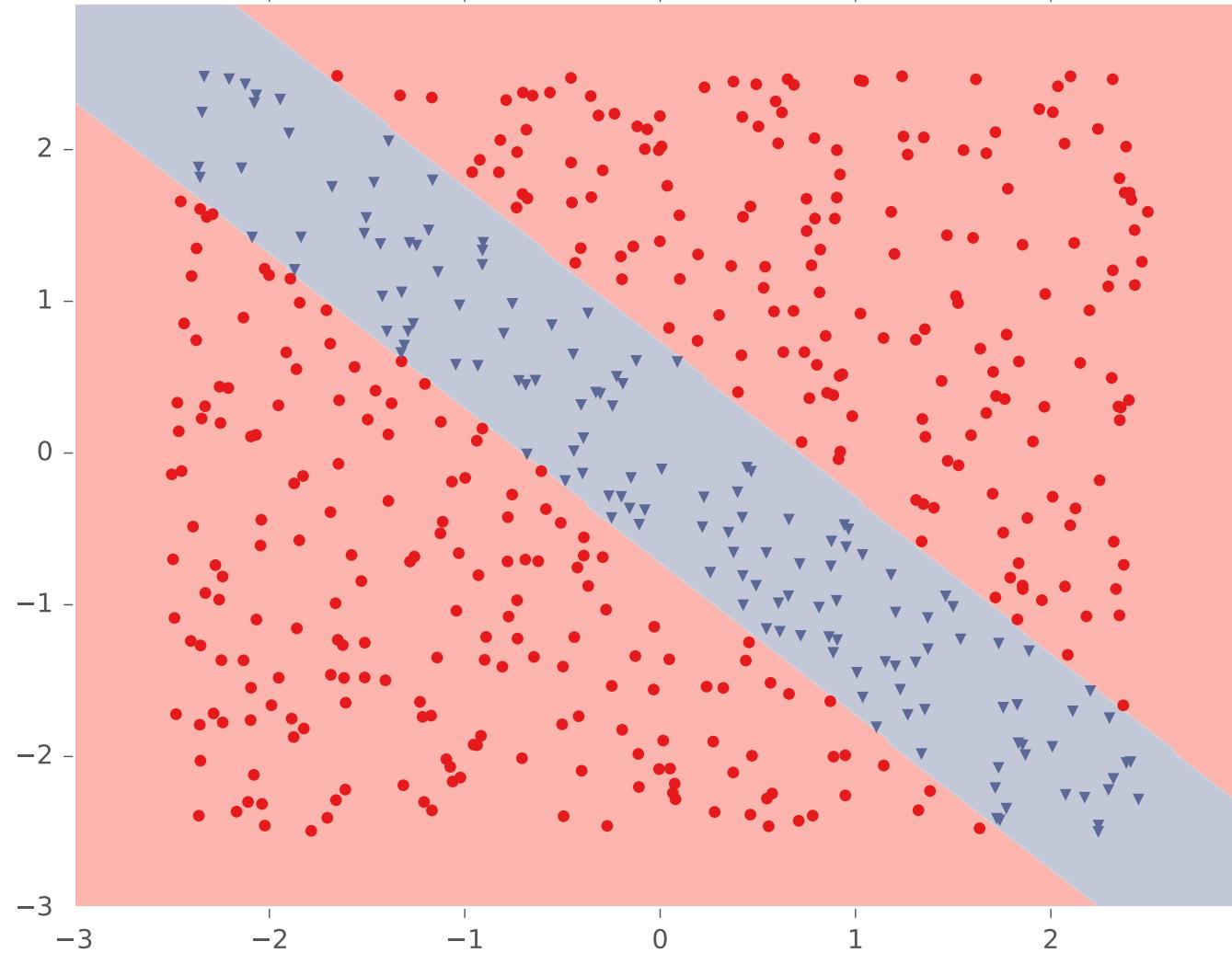


# Example #1: Diagonal Band



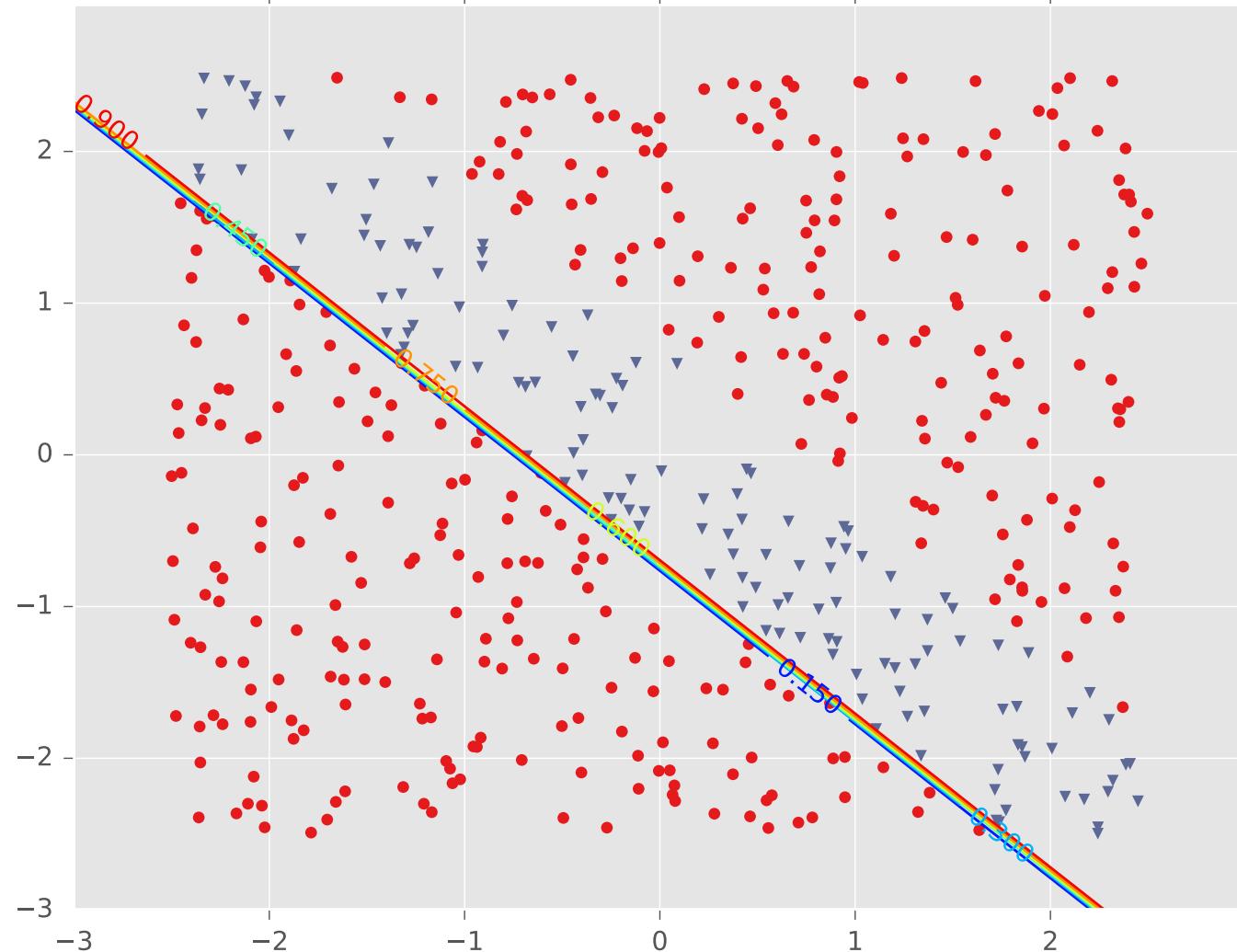
# Example #1: Diagonal Band

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



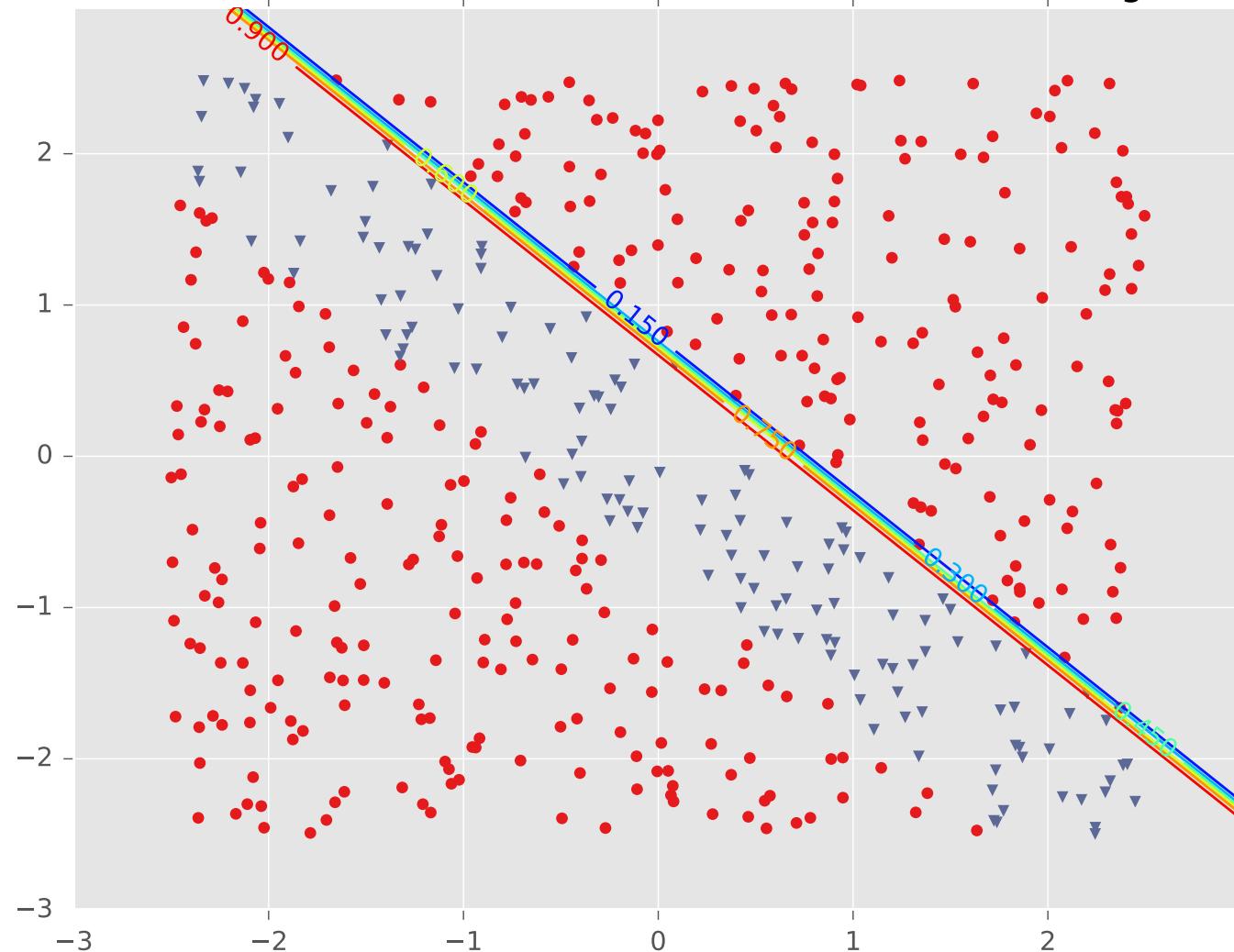
# Example #1: Diagonal Band

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



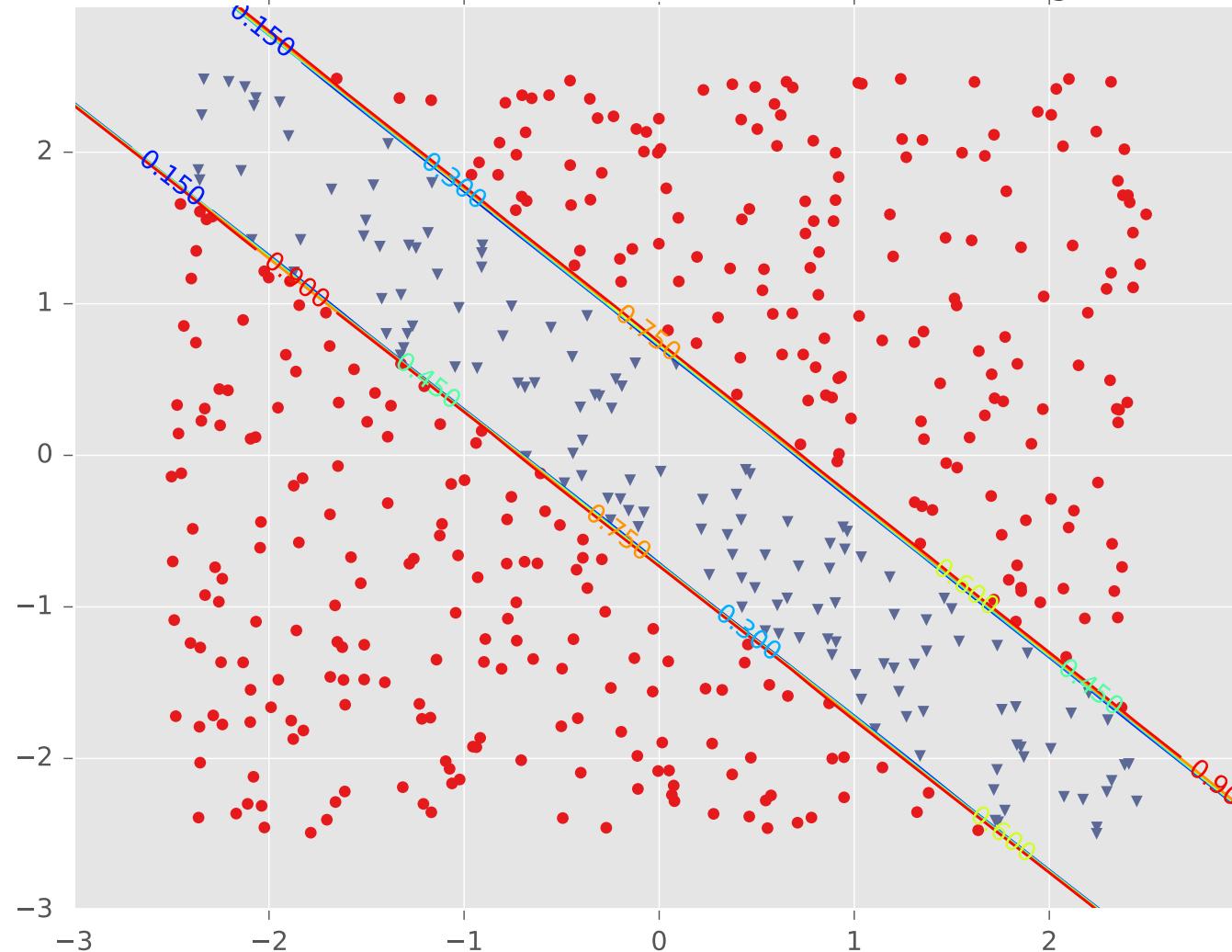
# Example #1: Diagonal Band

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

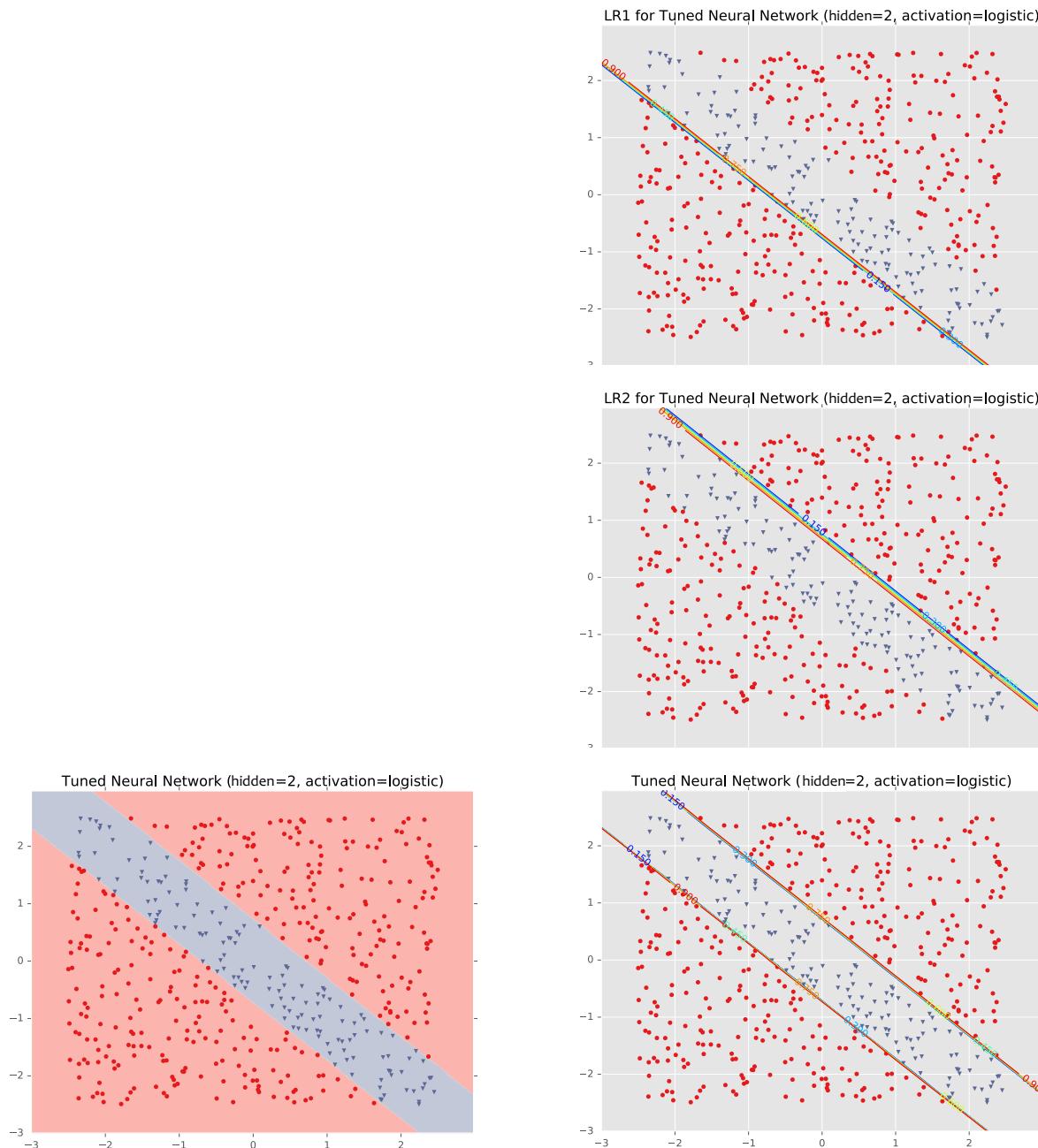


# Example #1: Diagonal Band

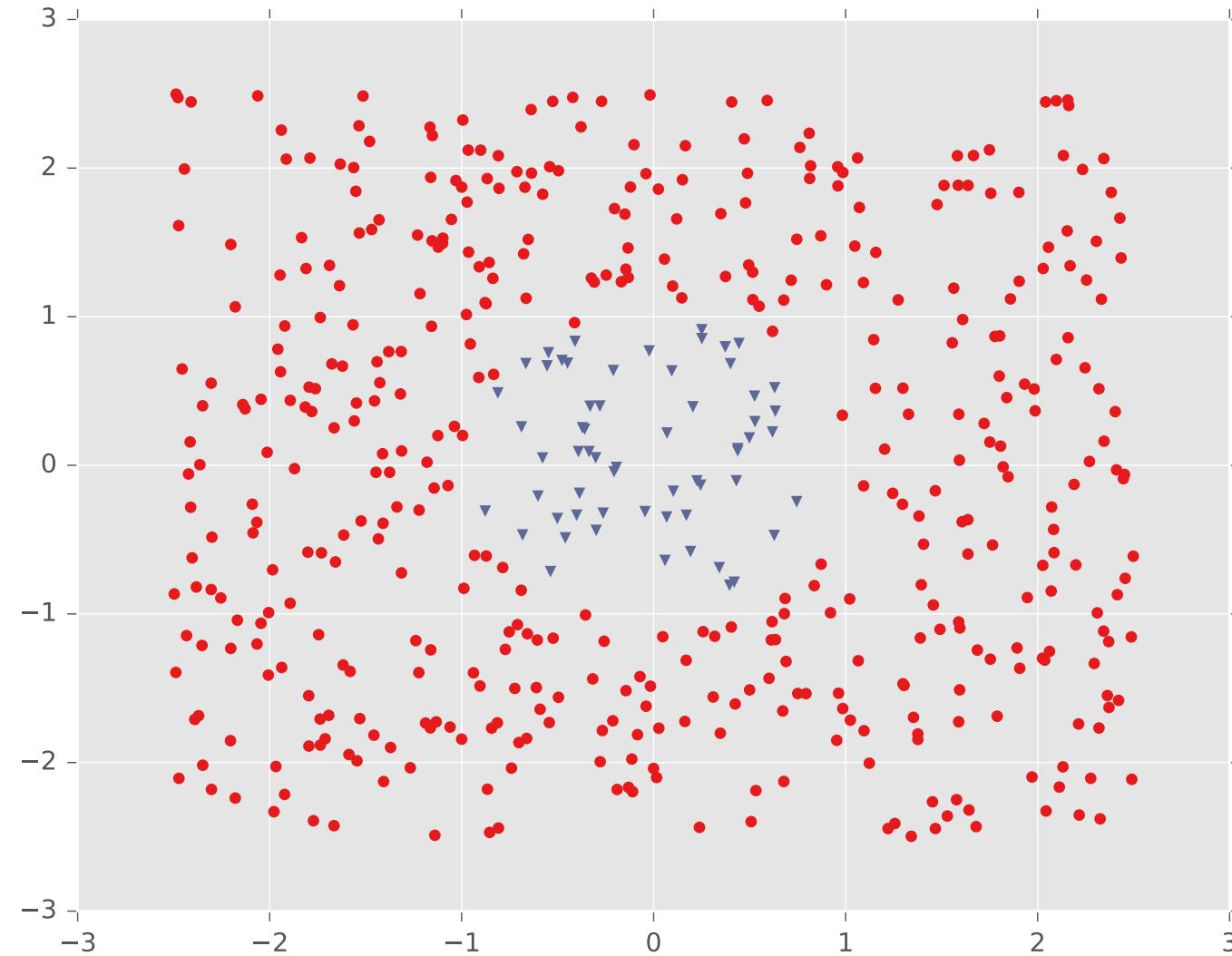
Tuned Neural Network (hidden=2, activation=logistic)



# Example #1: Diagonal Band



# Example #2: One Pocket

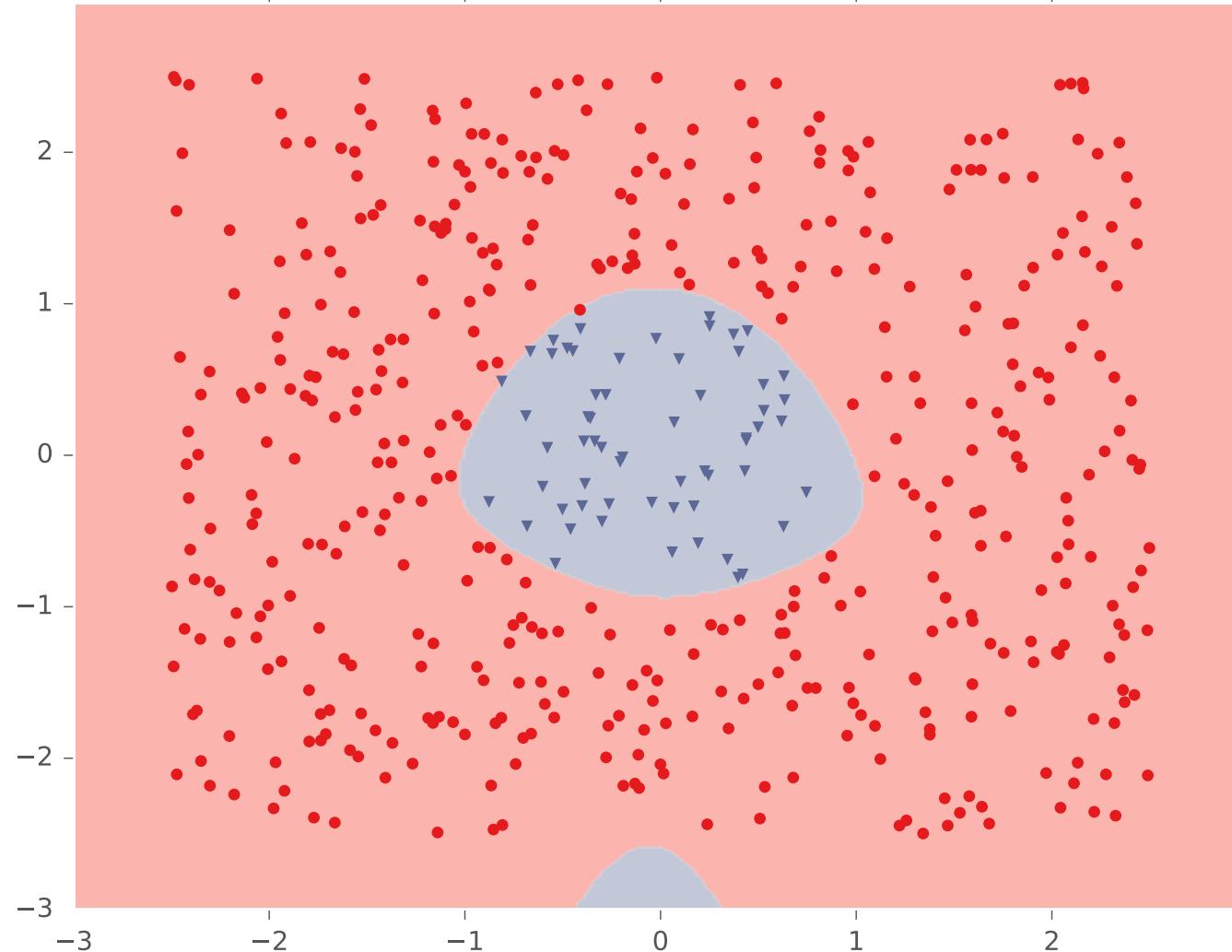


# Example #2: One Pocket



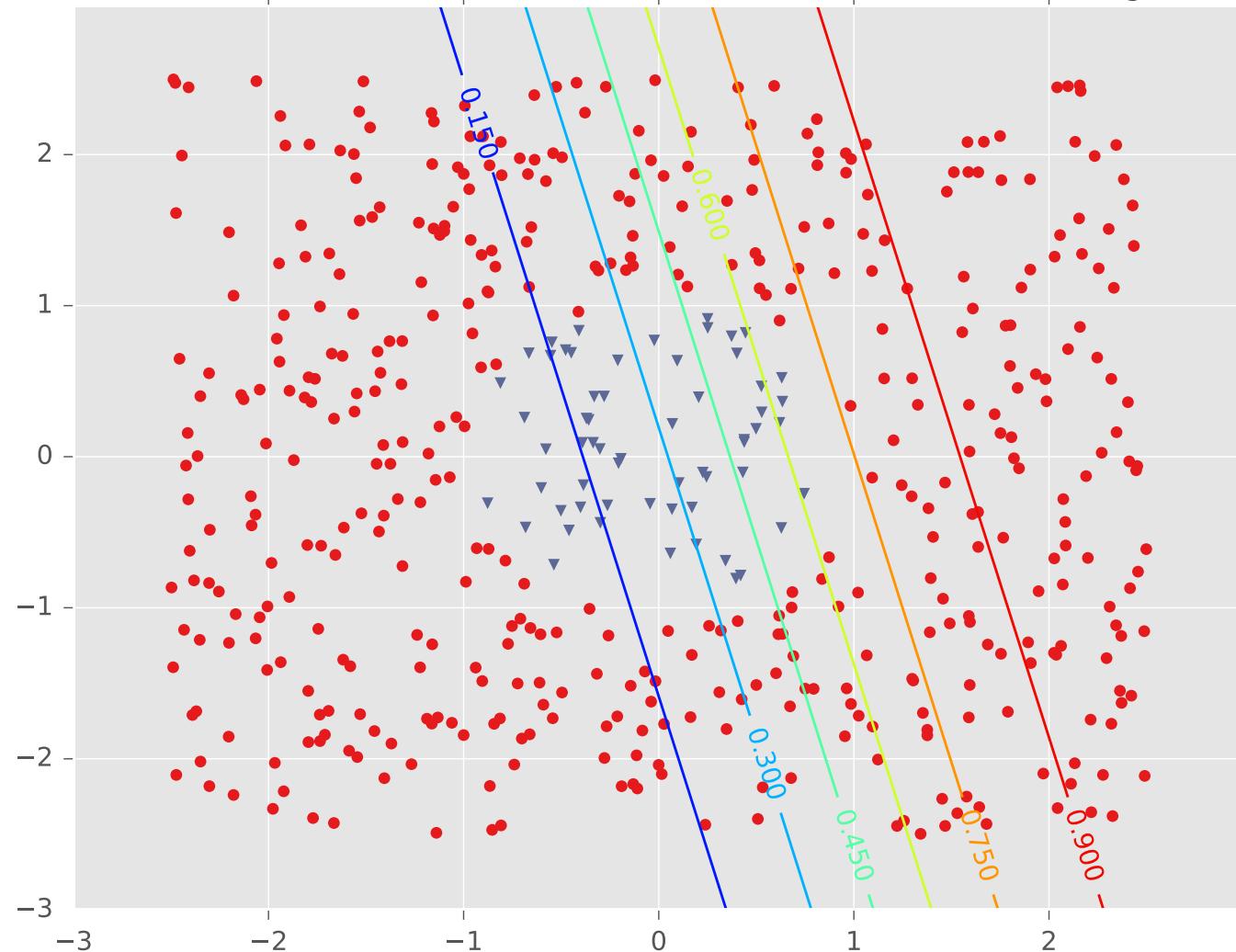
# Example #2: One Pocket

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



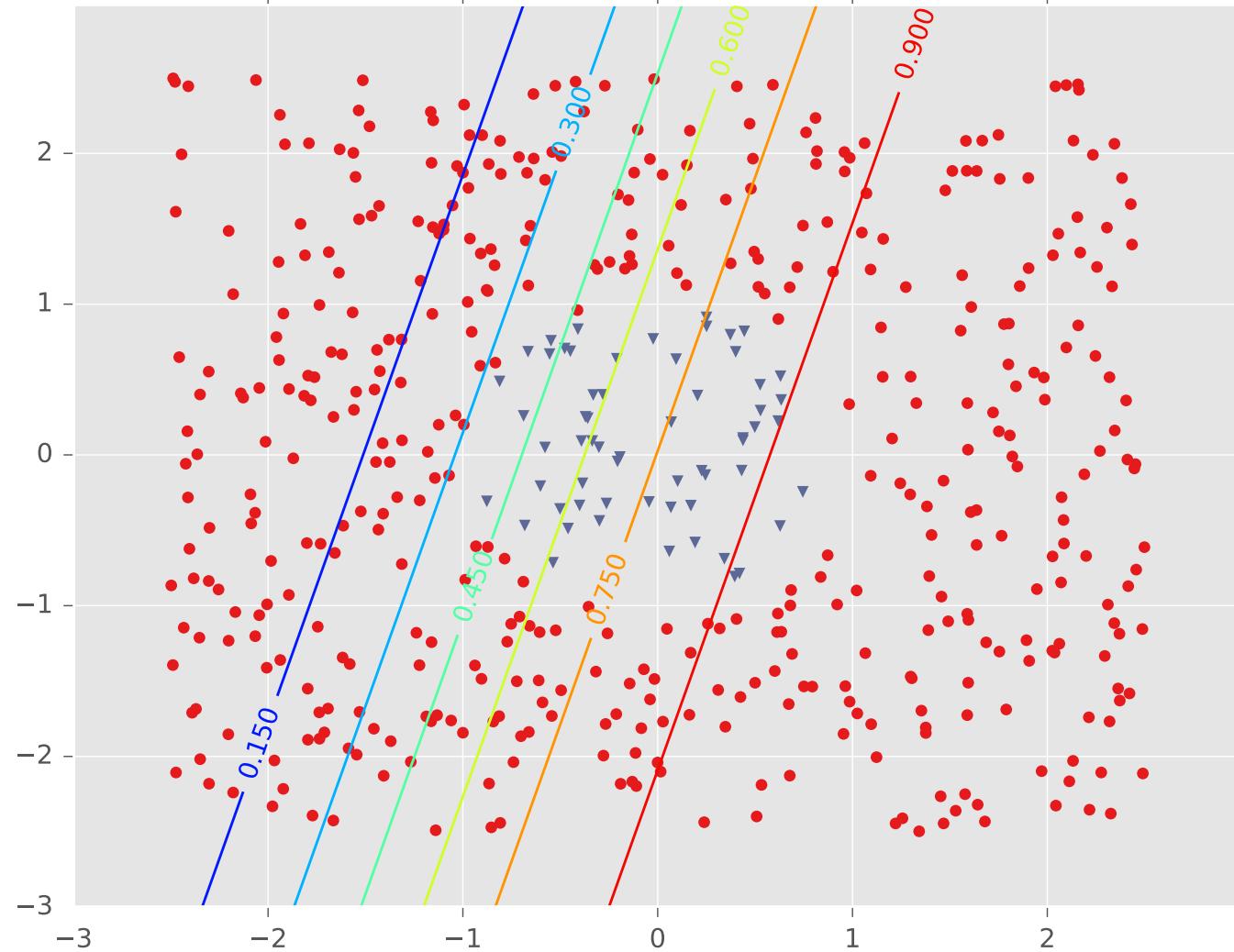
# Example #2: One Pocket

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



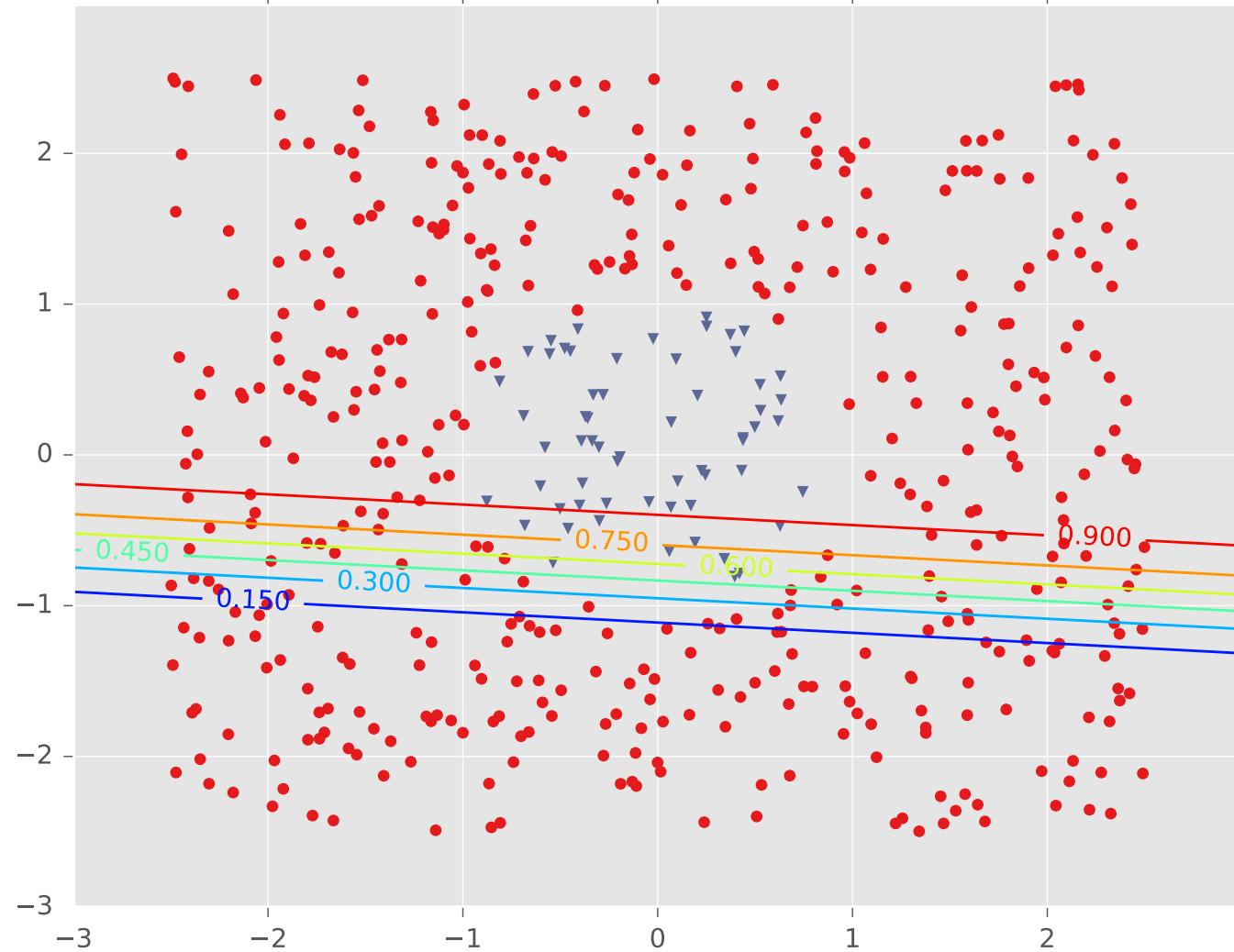
# Example #2: One Pocket

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



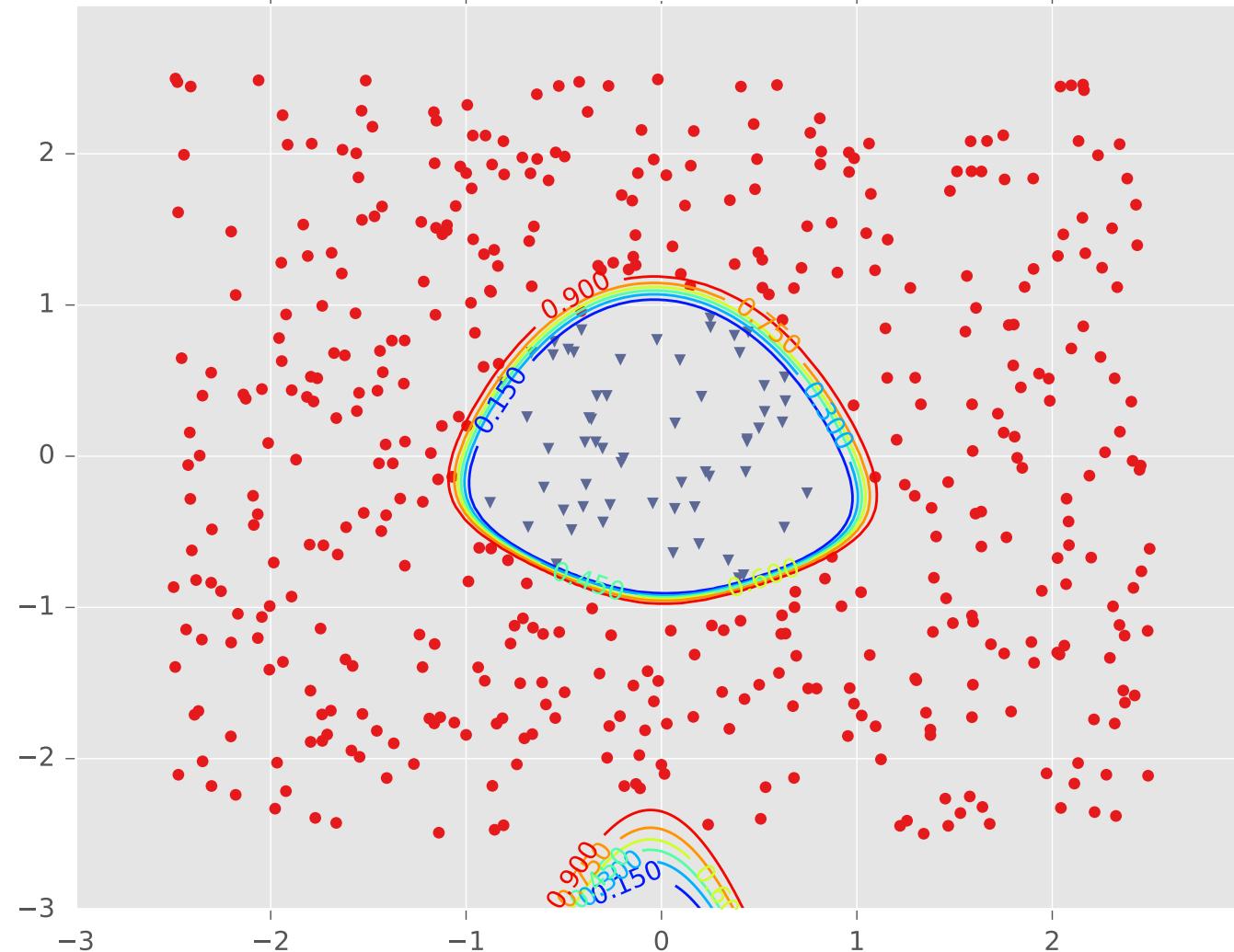
# Example #2: One Pocket

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

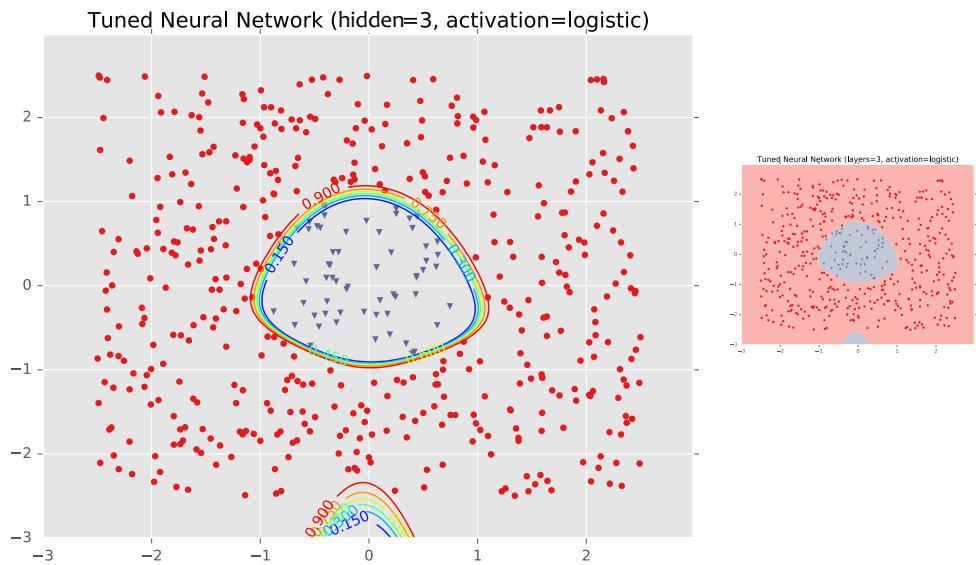
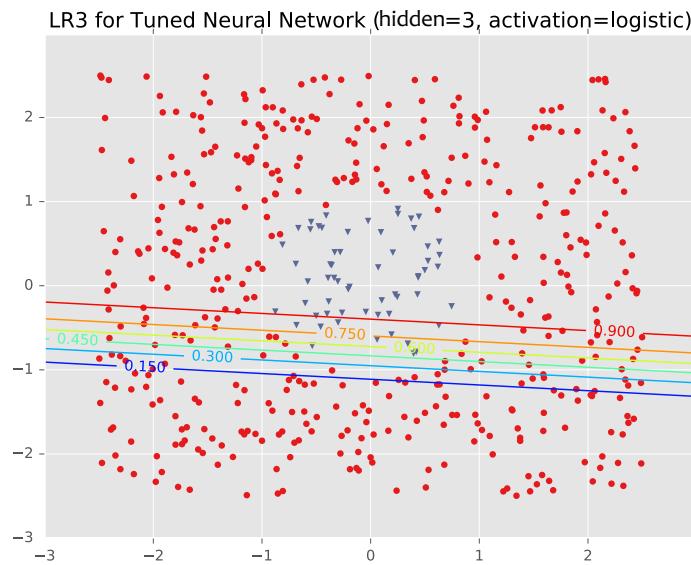
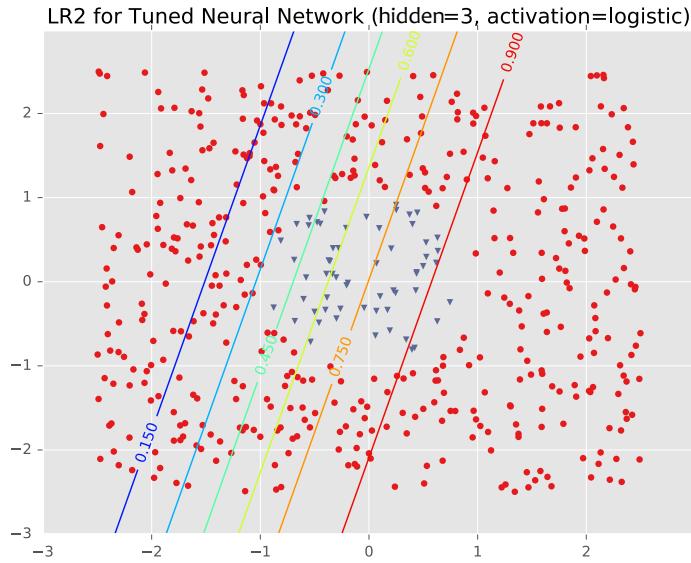
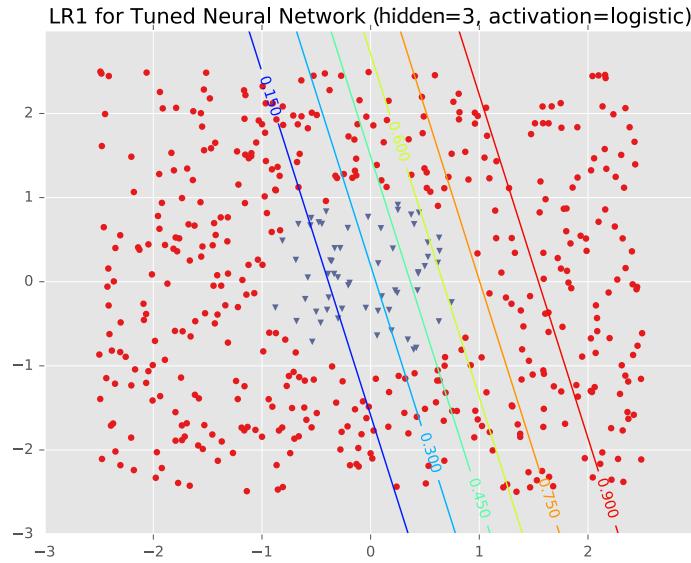


# Example #2: One Pocket

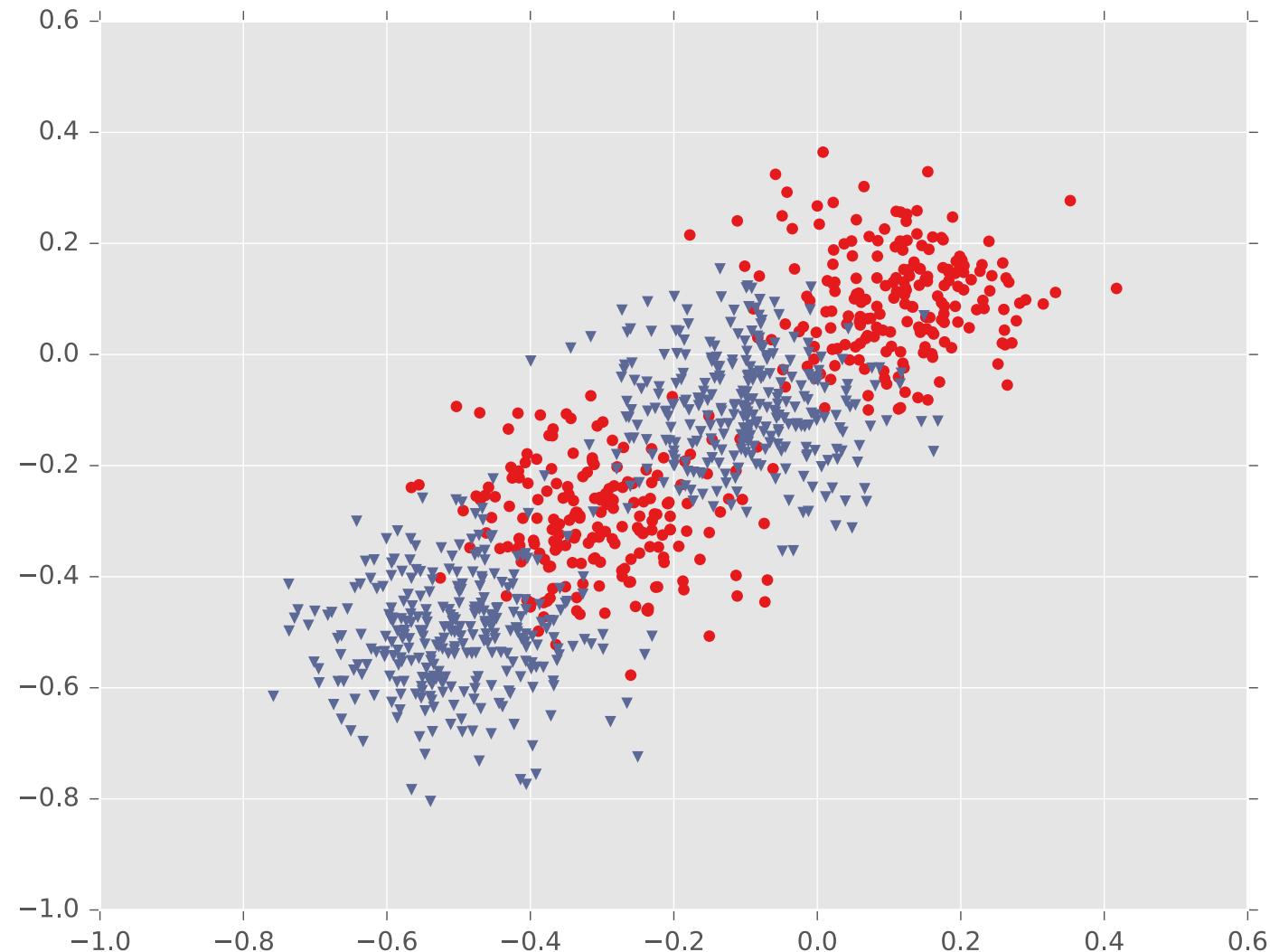
Tuned Neural Network (hidden=3, activation=logistic)



# Example #2: One Pocket



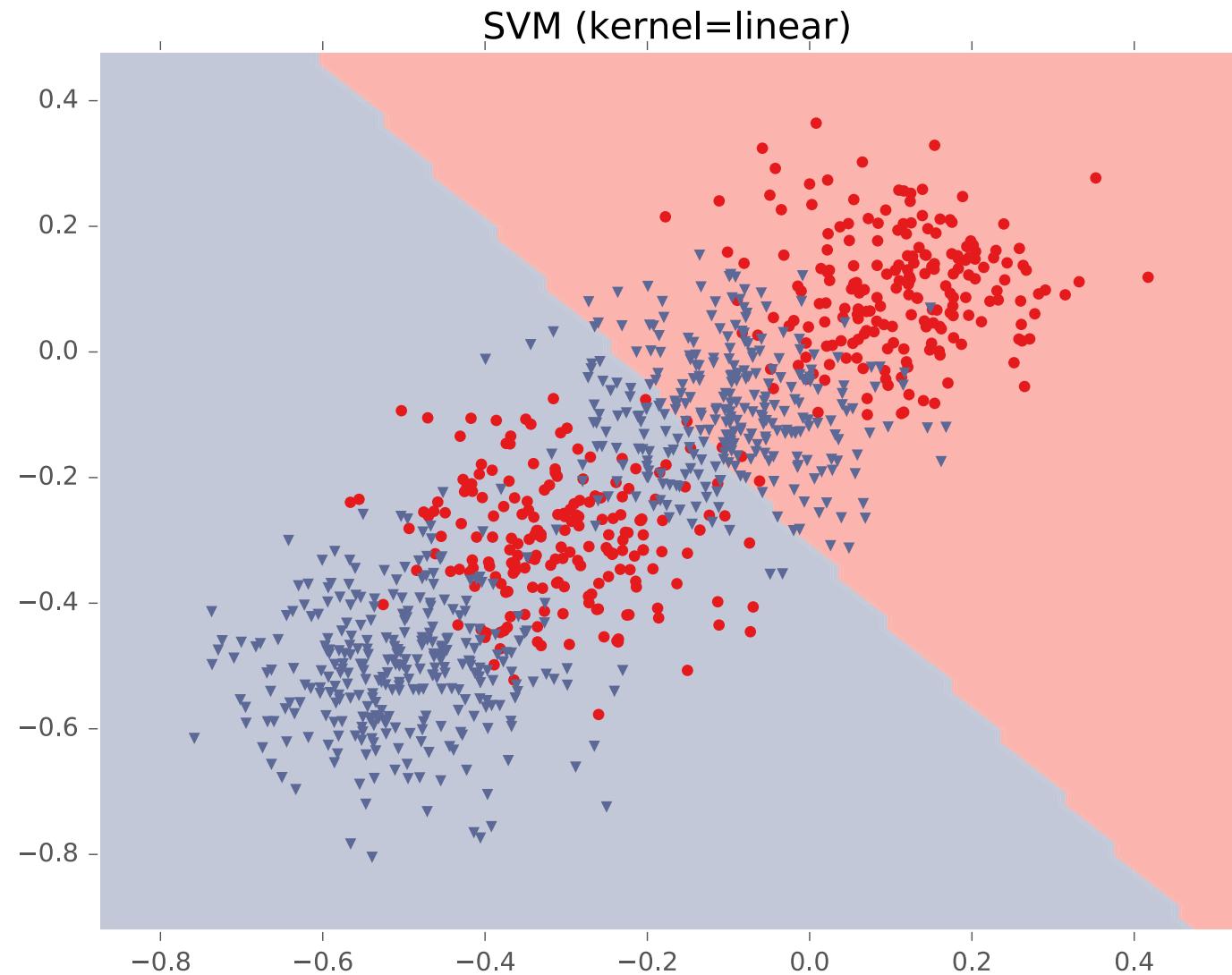
# Example #3: Four Gaussians



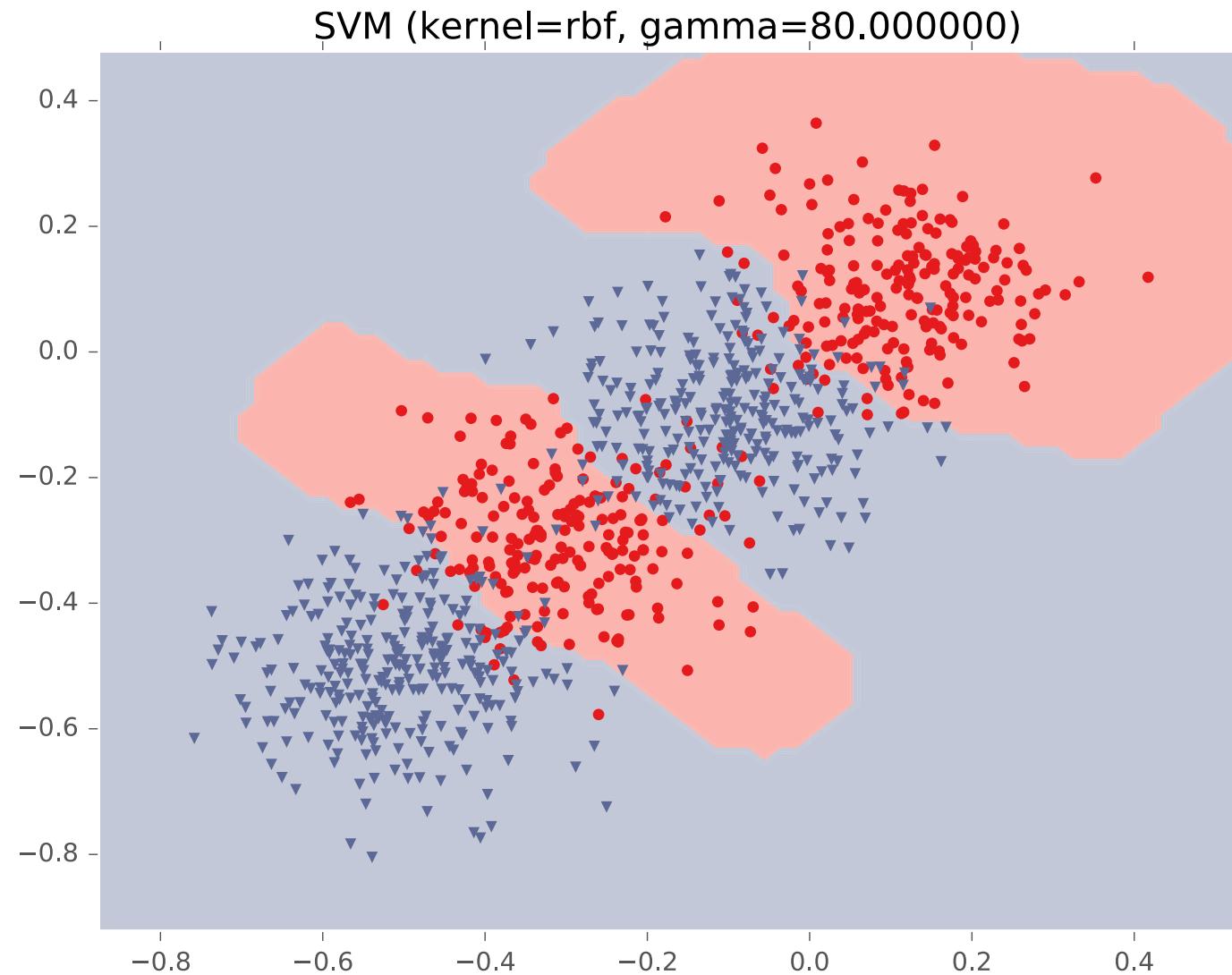
# Example #3: Four Gaussians



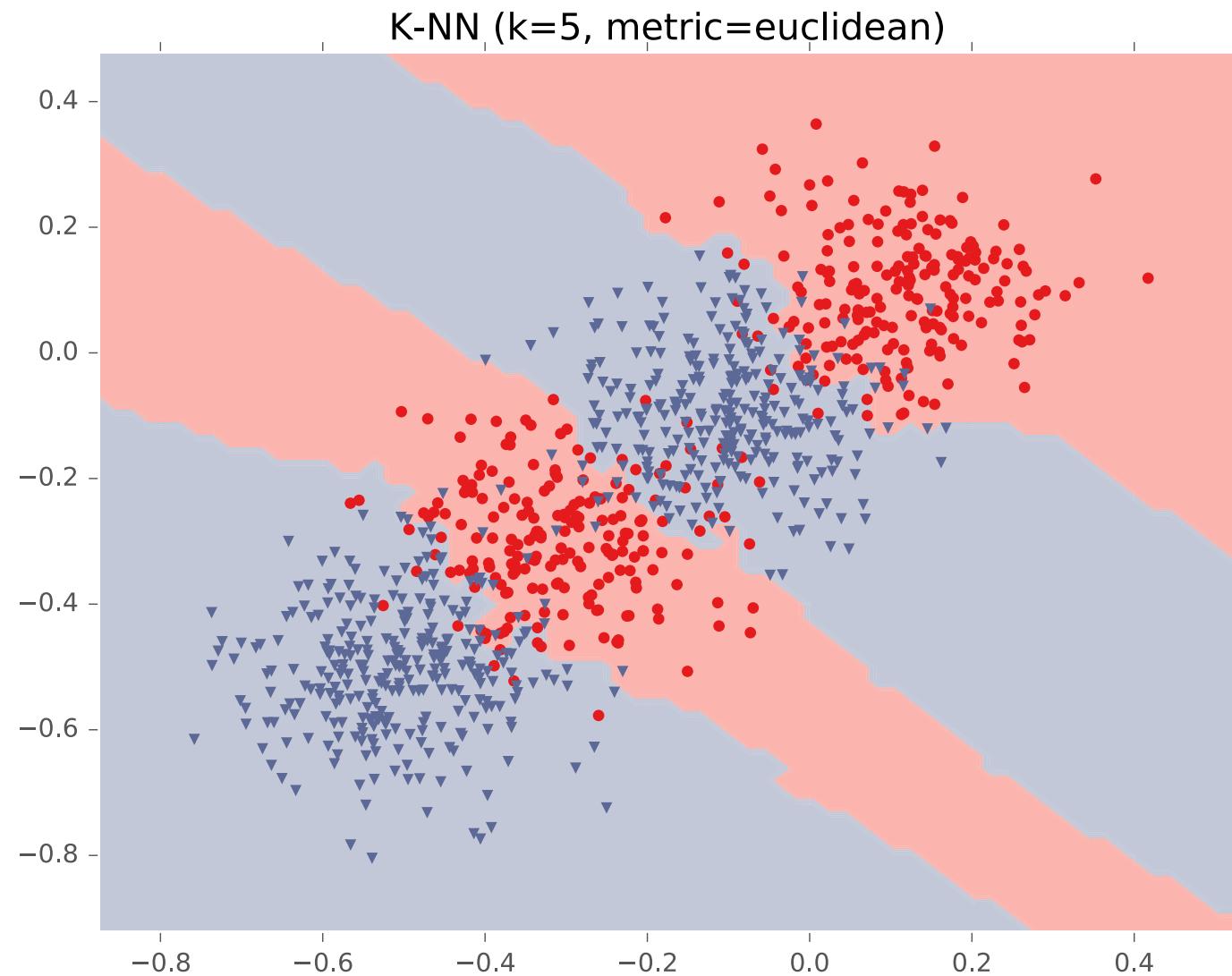
# Example #3: Four Gaussians



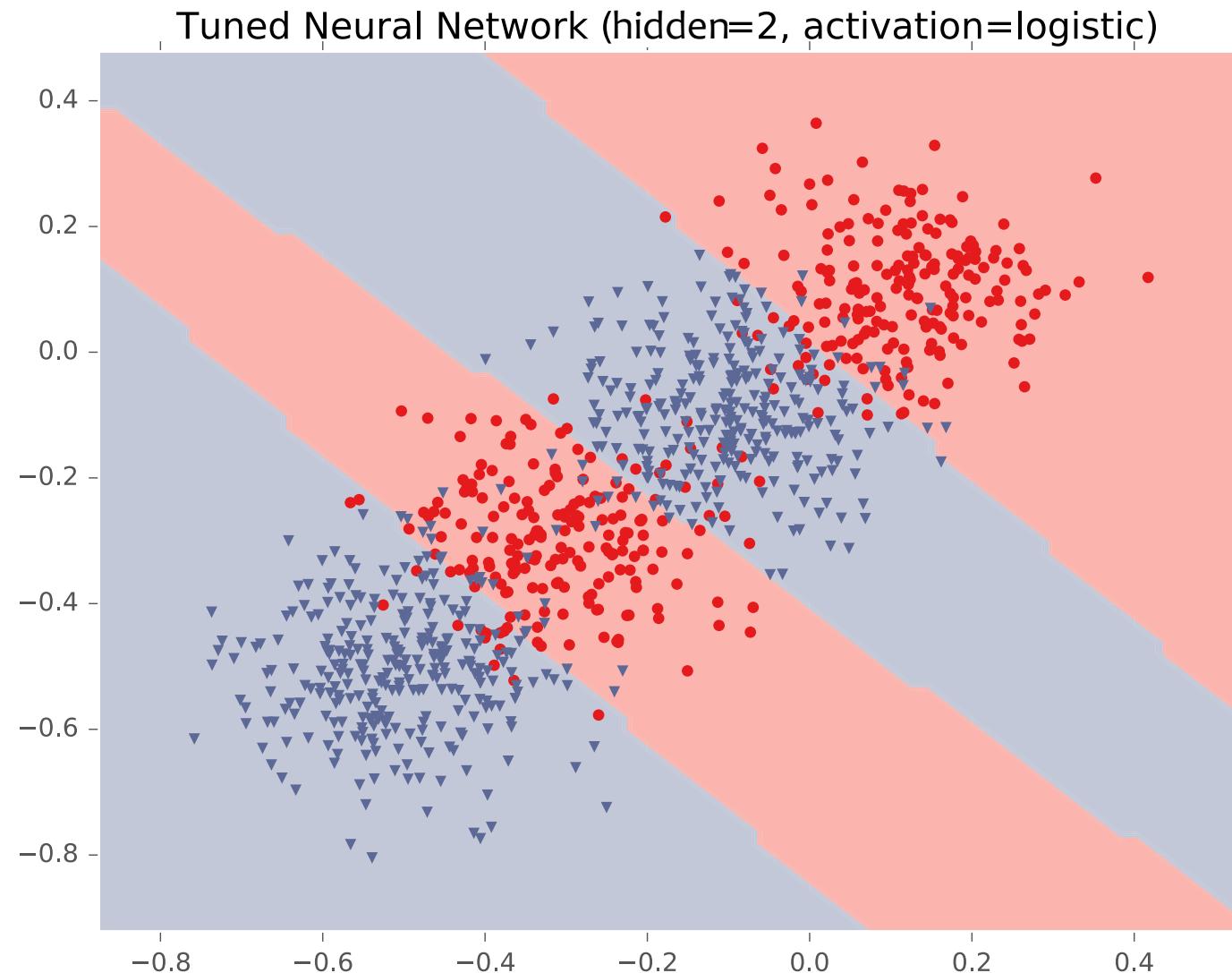
# Example #3: Four Gaussians



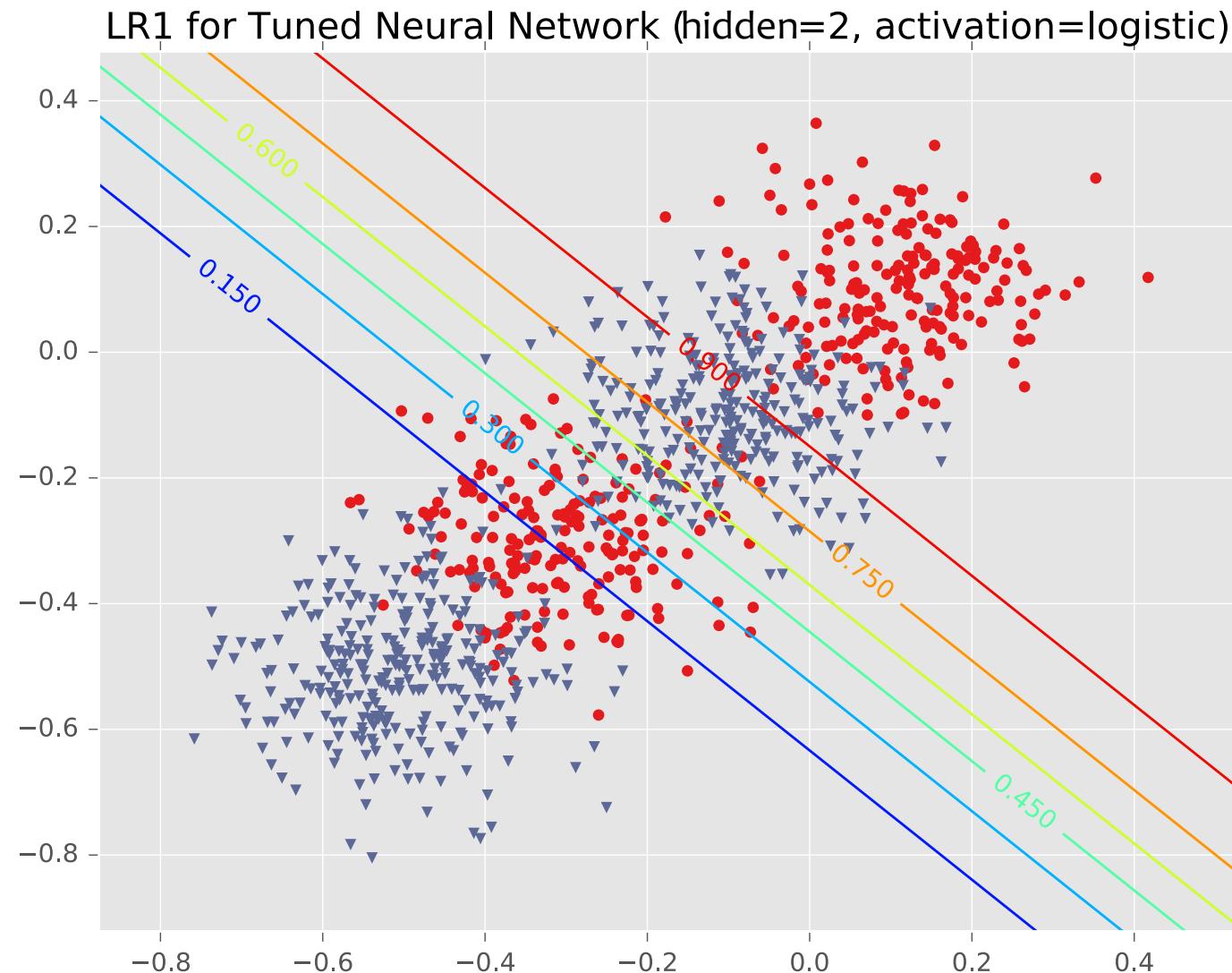
# Example #3: Four Gaussians



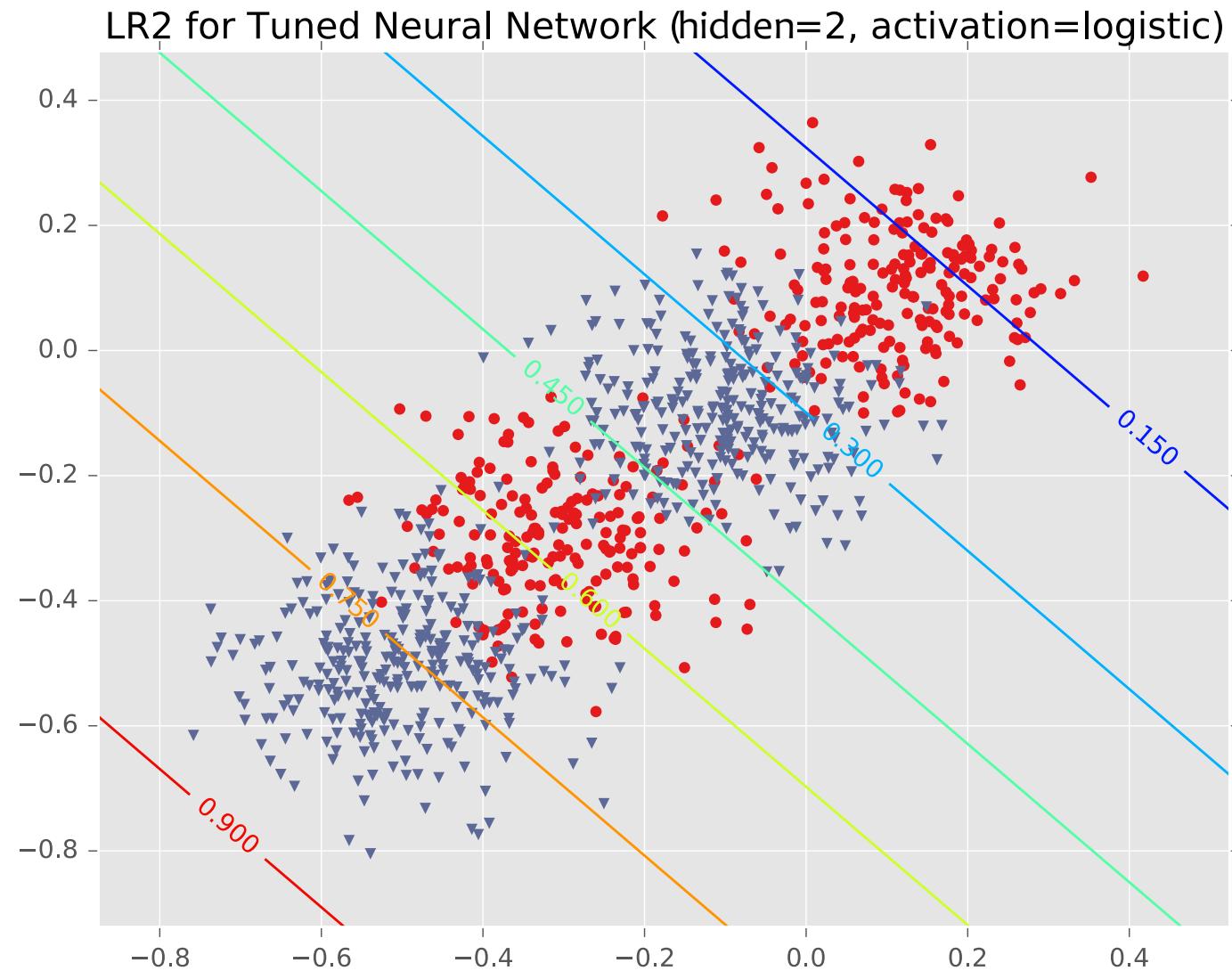
# Example #3: Four Gaussians



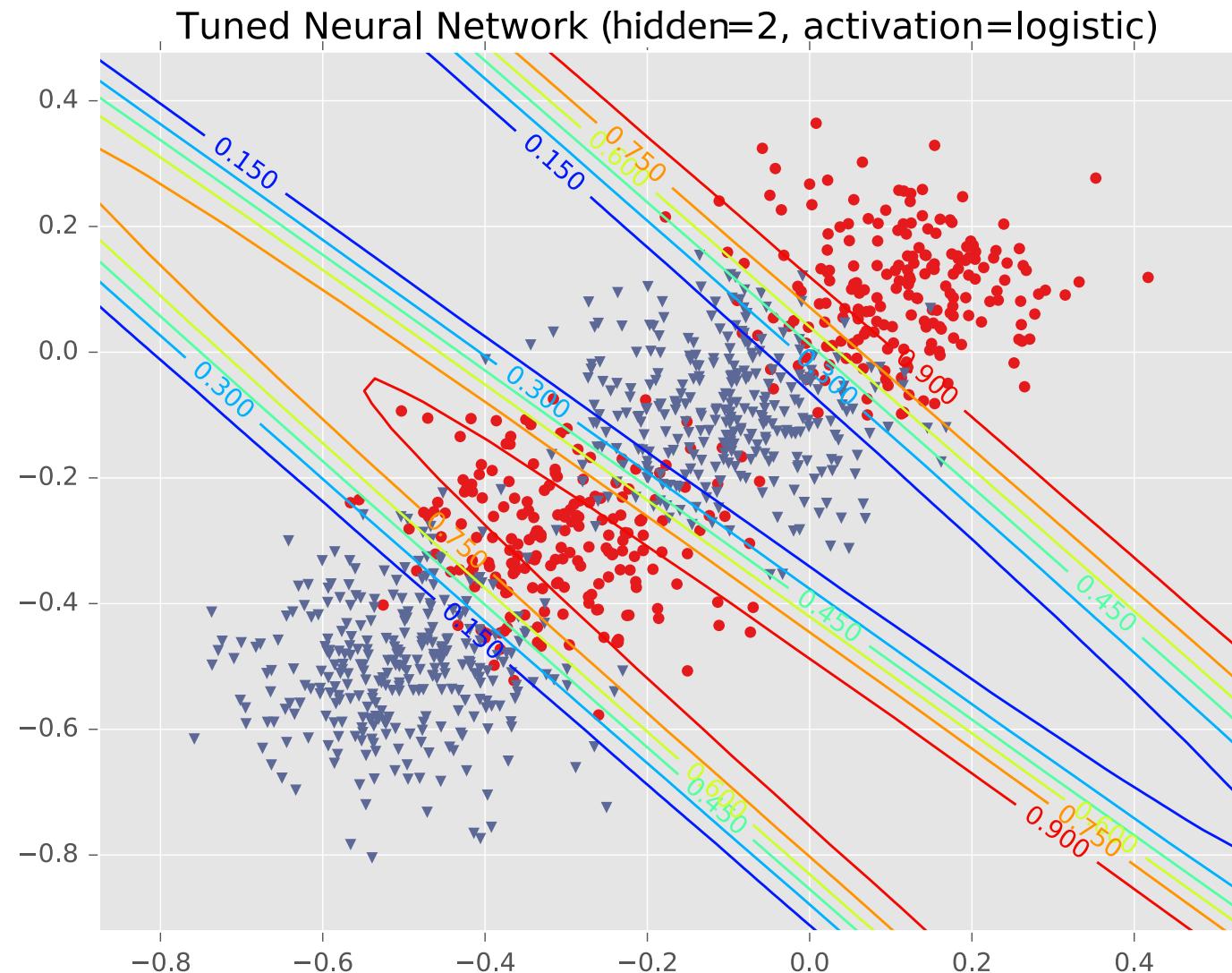
# Example #3: Four Gaussians



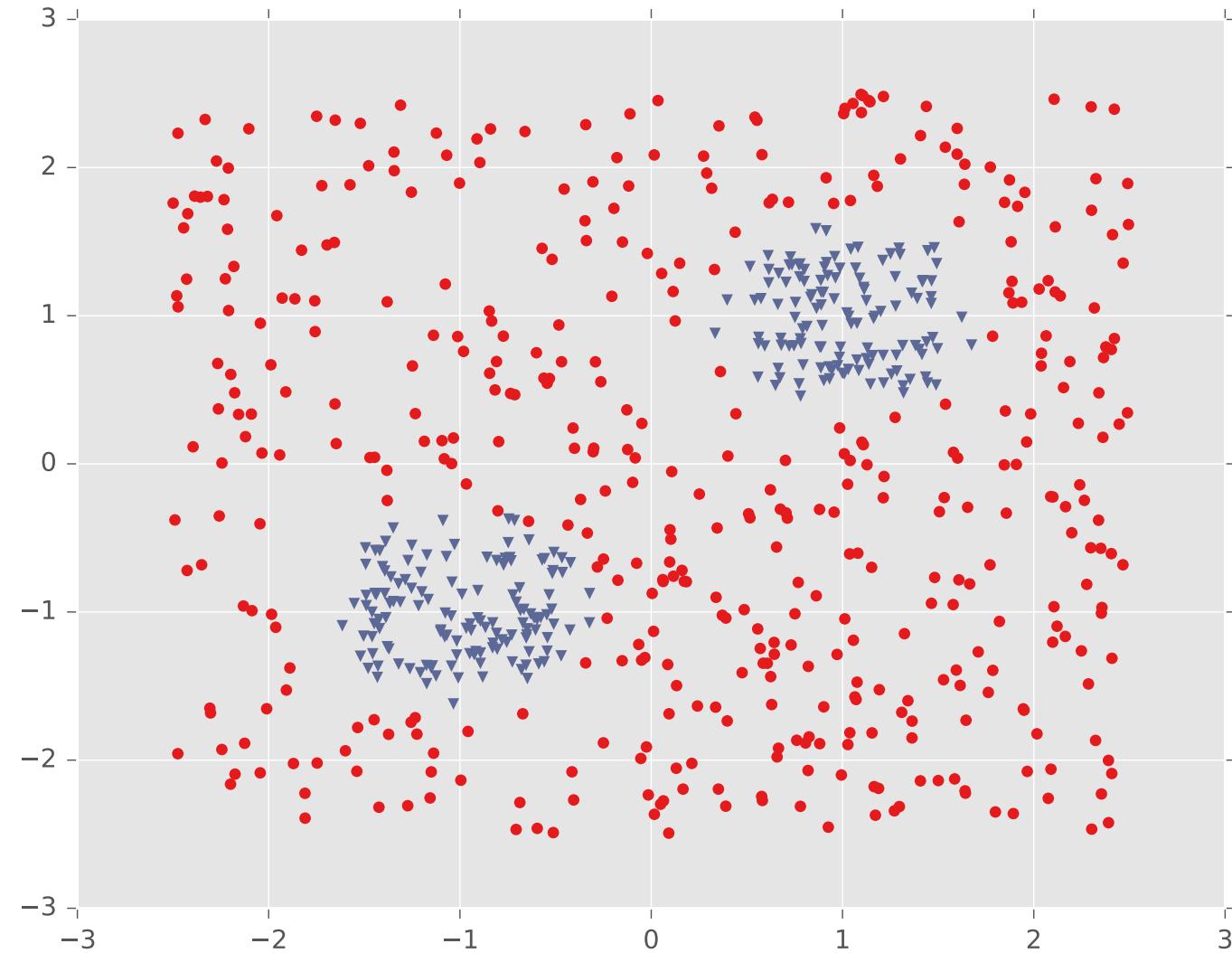
# Example #3: Four Gaussians



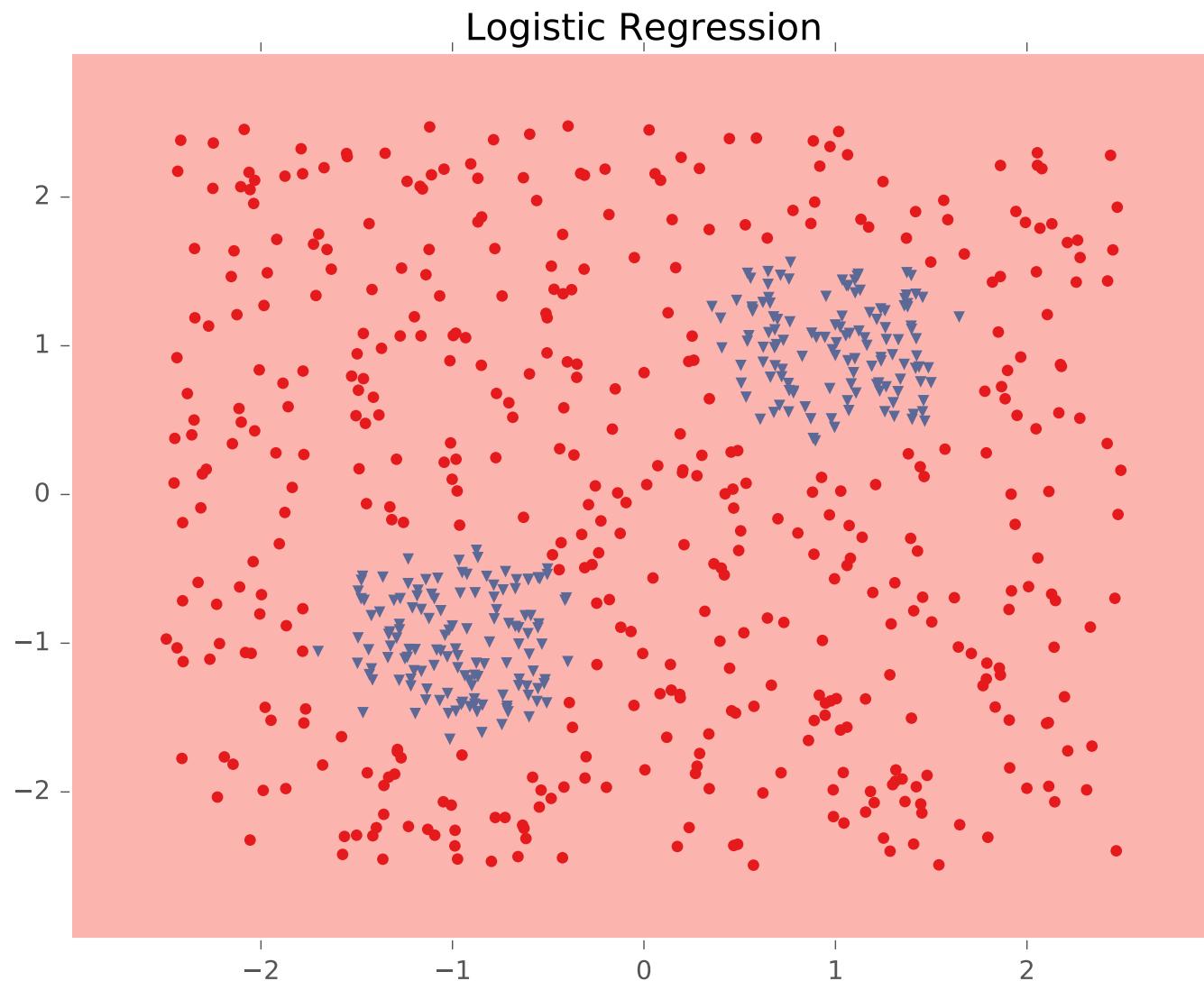
# Example #3: Four Gaussians



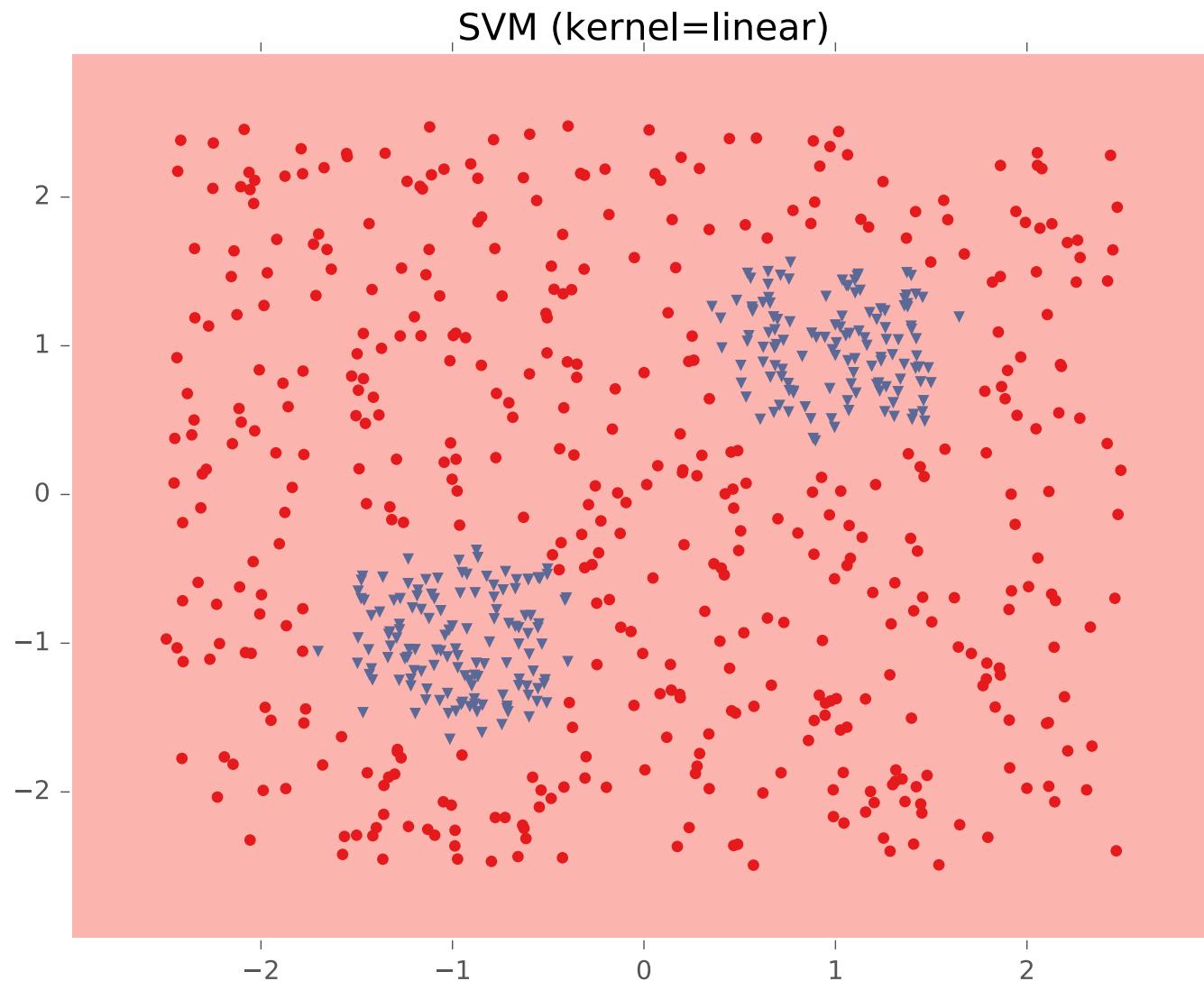
# Example #4: Two Pockets



# Example #4: Two Pockets

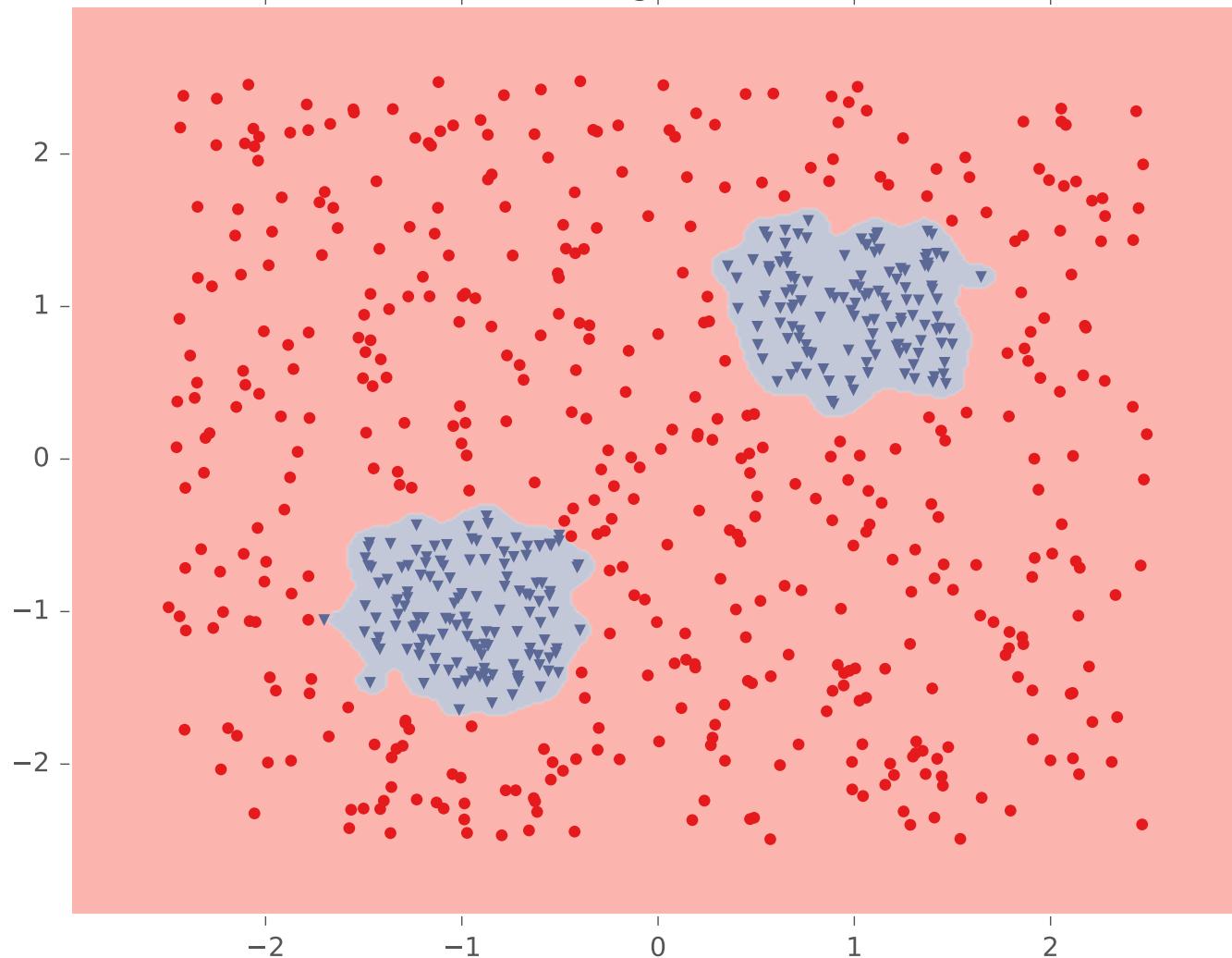


# Example #4: Two Pockets



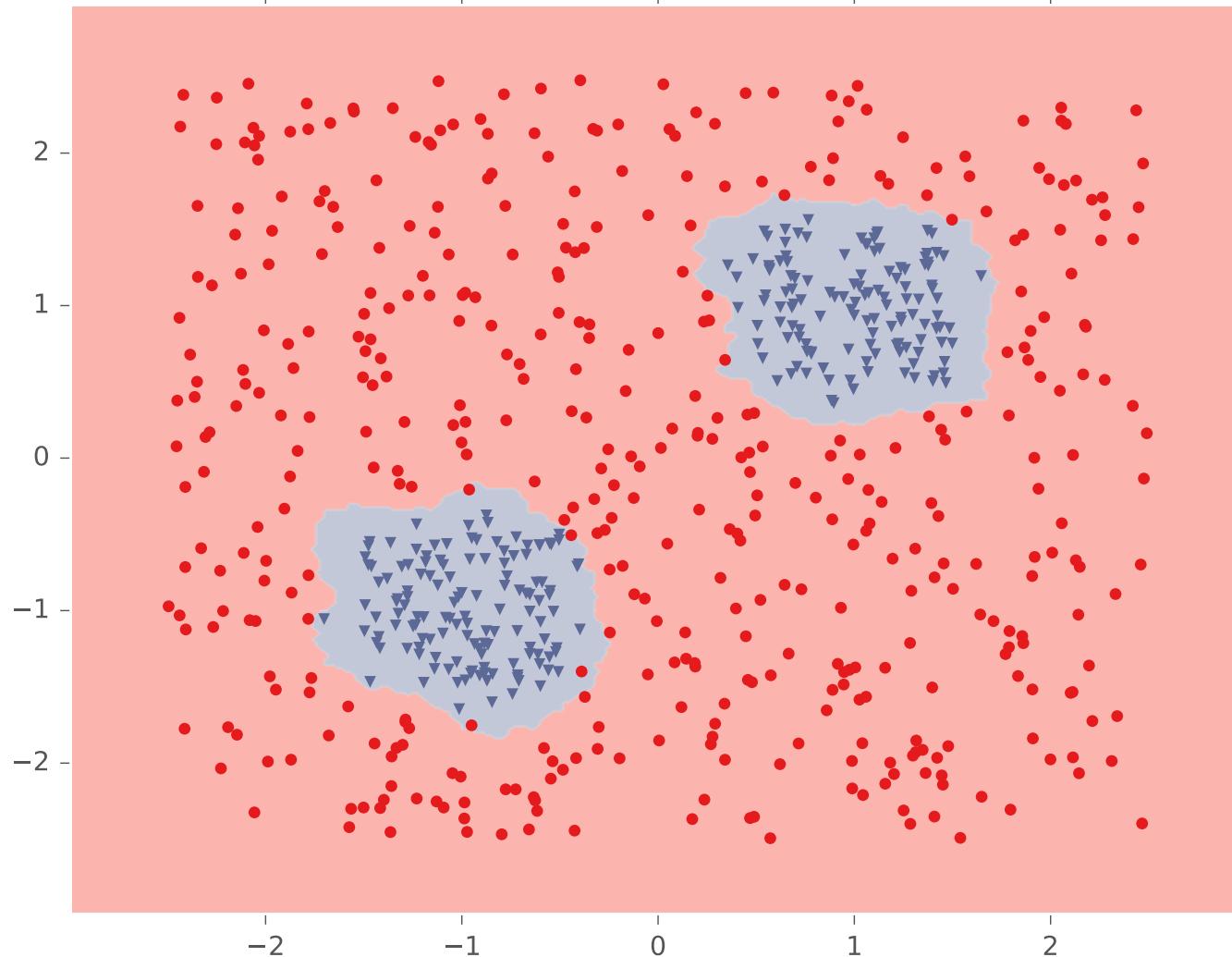
# Example #4: Two Pockets

SVM (kernel=rbf, gamma=80.000000)



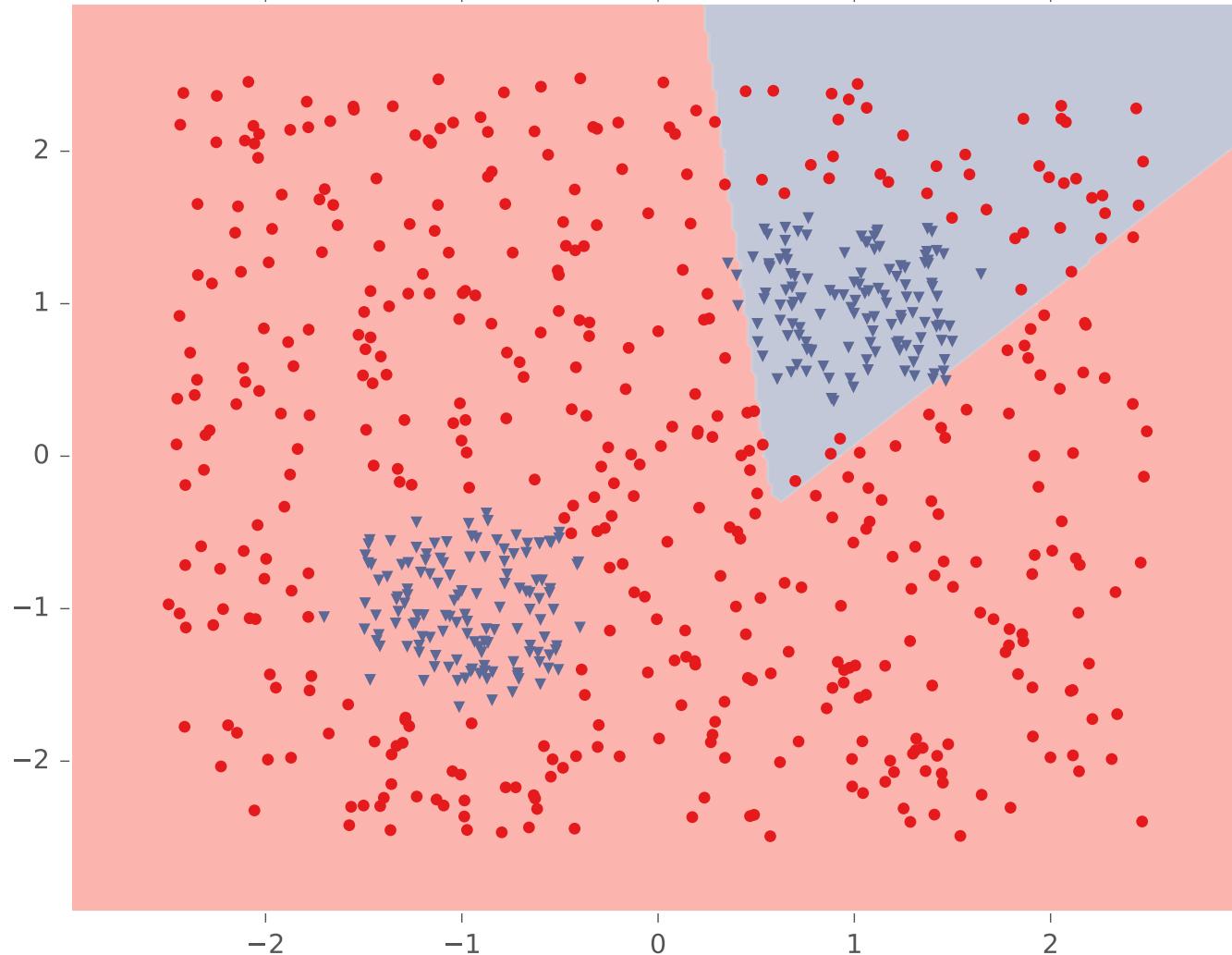
# Example #4: Two Pockets

K-NN (k=5, metric=euclidean)



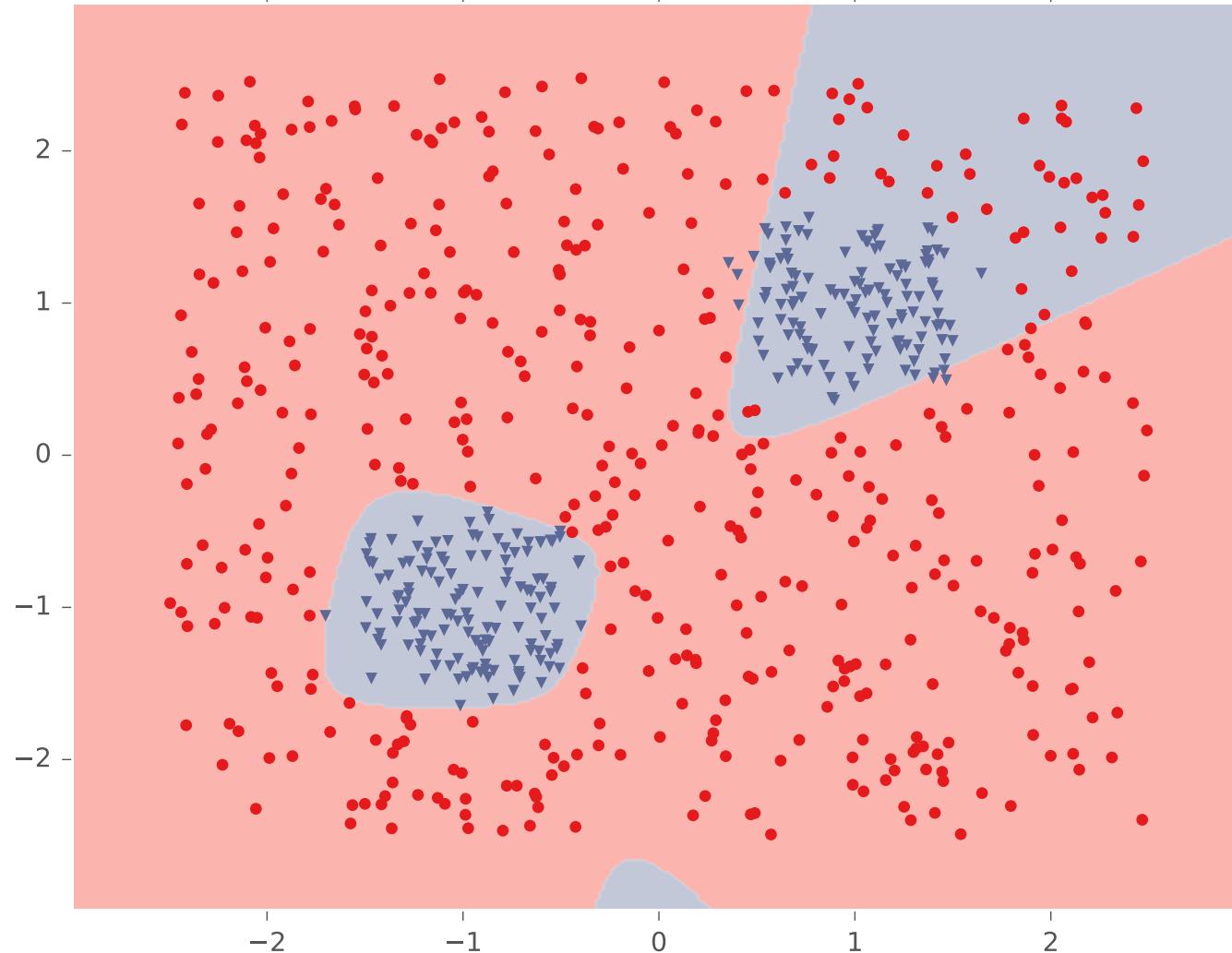
# Example #4: Two Pockets

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



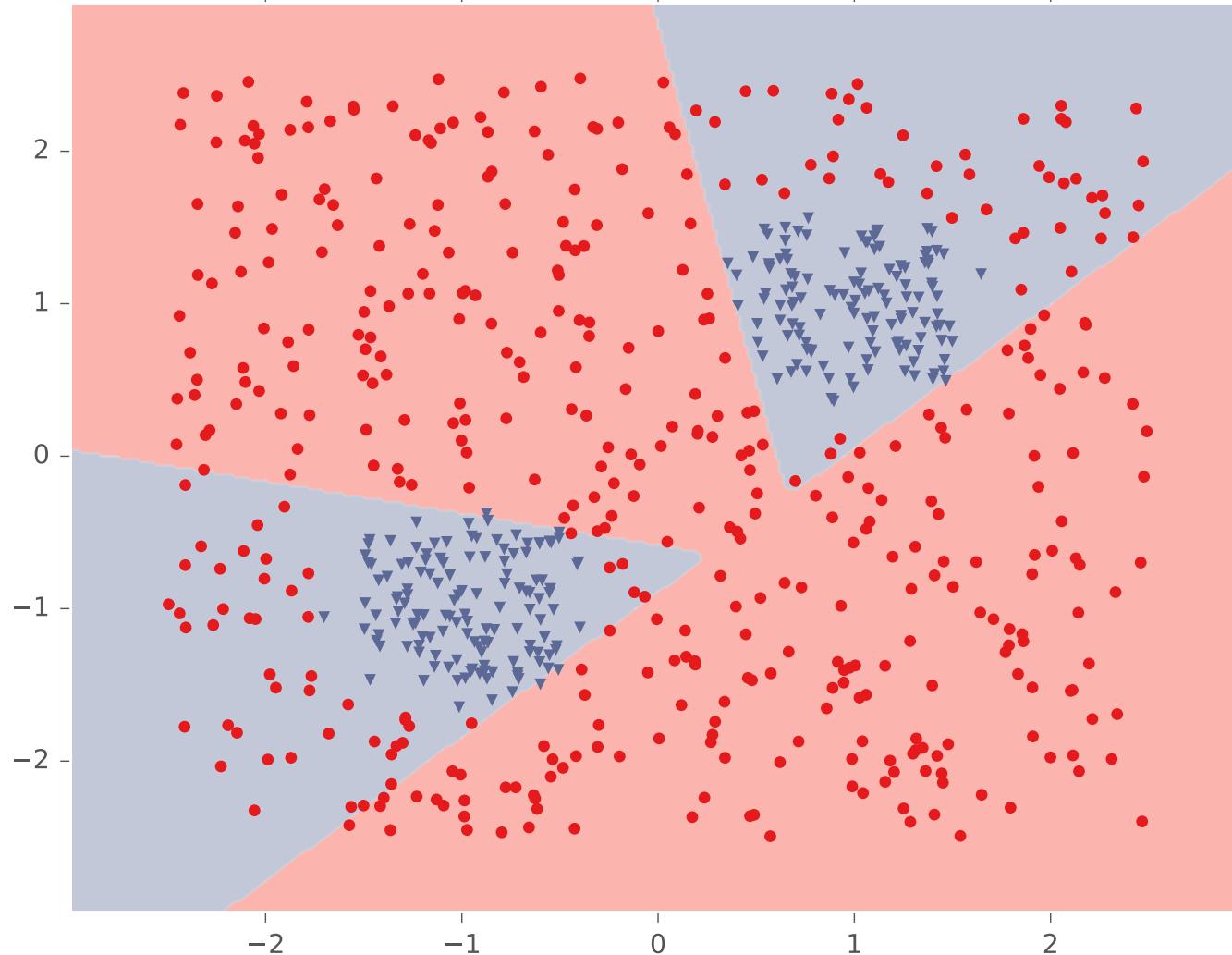
# Example #4: Two Pockets

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



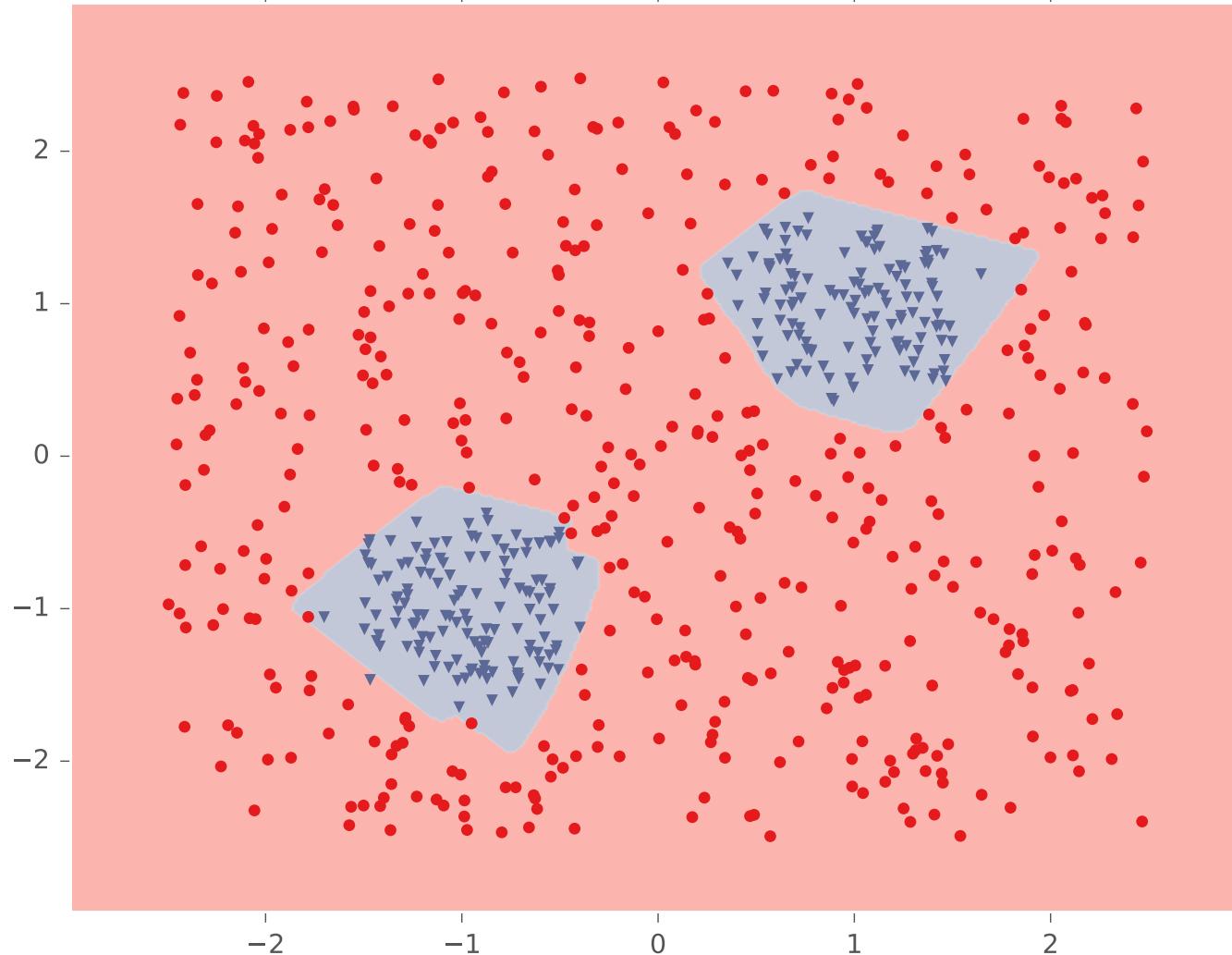
# Example #4: Two Pockets

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



# Example #4: Two Pockets

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



# 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