**10-601 Introduction to Machine Learning**

Machine Learning Department
School of Computer Science
Carnegie Mellon University

# Neural Networks

# +

# Backpropagation

Matt Gormley
Lecture 13
Oct. 7, 2019

# Reminders

- **Homework 4: Logistic Regression**
  - **Out: Wed, Sep. 25**
  - **Due: Fri, Oct. 11 at 11:59pm**
- **Homework 5: Neural Networks**
  - **Out: Fri, Oct. 11**
  - **Due: Fri, Oct. 25 at 11:59pm**
- **Today's In-Class Poll**
  - **http://p13.mlcourse.org**

# Q&A

**Q:** What is mini-batch SGD?

**A:** A variant of SGD…

# Mini-Batch SGD

- **Gradient Descent**:
  Compute true gradient exactly from all N examples

- **Mini-Batch SGD**:
  Approximate true gradient by the average gradient of K randomly chosen examples

- **Stochastic Gradient Descent (SGD)**:
  Approximate true gradient by the gradient of one randomly chosen example

# Mini-Batch SGD

**while** not converged: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \lambda \mathbf{g}$

**Three variants of first-order optimization:**

Gradient Descent: $\mathbf{g} = \nabla J(\boldsymbol{\theta}) = \dfrac{1}{N} \sum\limits_{i=1}^{N} \nabla J^{(i)}(\boldsymbol{\theta})$

SGD: $\mathbf{g} = \nabla J^{(i)}(\boldsymbol{\theta})$      where $i$ sampled uniformly

Mini-batch SGD: $\mathbf{g} = \dfrac{1}{S} \sum\limits_{s=1}^{S} \nabla J^{(i_s)}(\boldsymbol{\theta})$      where $i_s$ sampled uniformly $\forall s$

# NEURAL NETWORKS

# Neural Networks

*Chalkboard*

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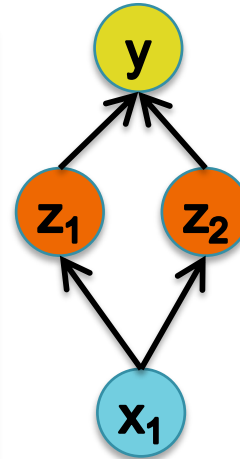- – Example: Feed Forward Neural Network

# Neural Network Parameters

**Question:**

Suppose you are training a one-hidden layer neural network with sigmoid activations for binary classification.



**True or False**: There is a unique set of parameters that maximize the likelihood of the dataset above.
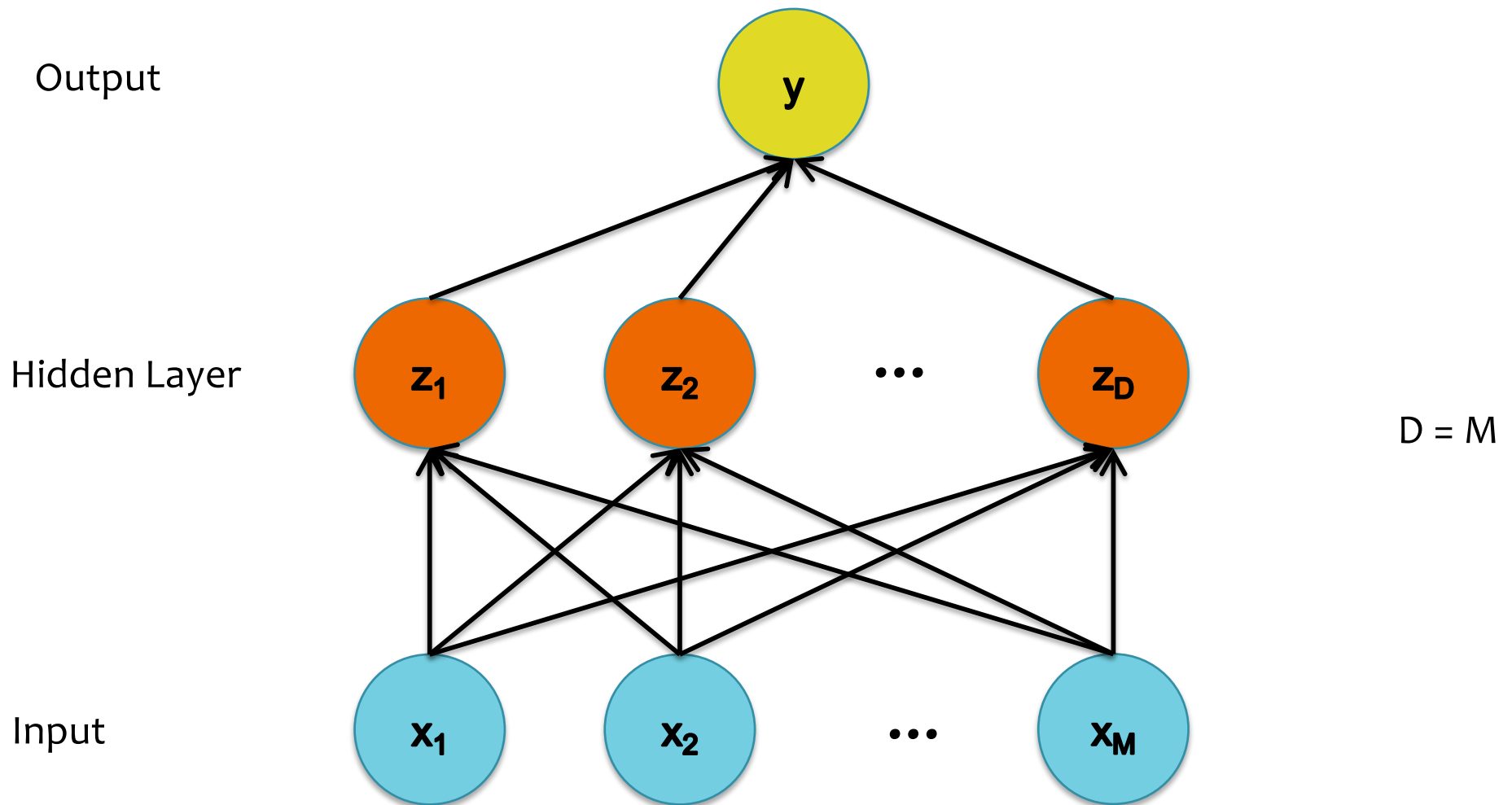
**Answer:**

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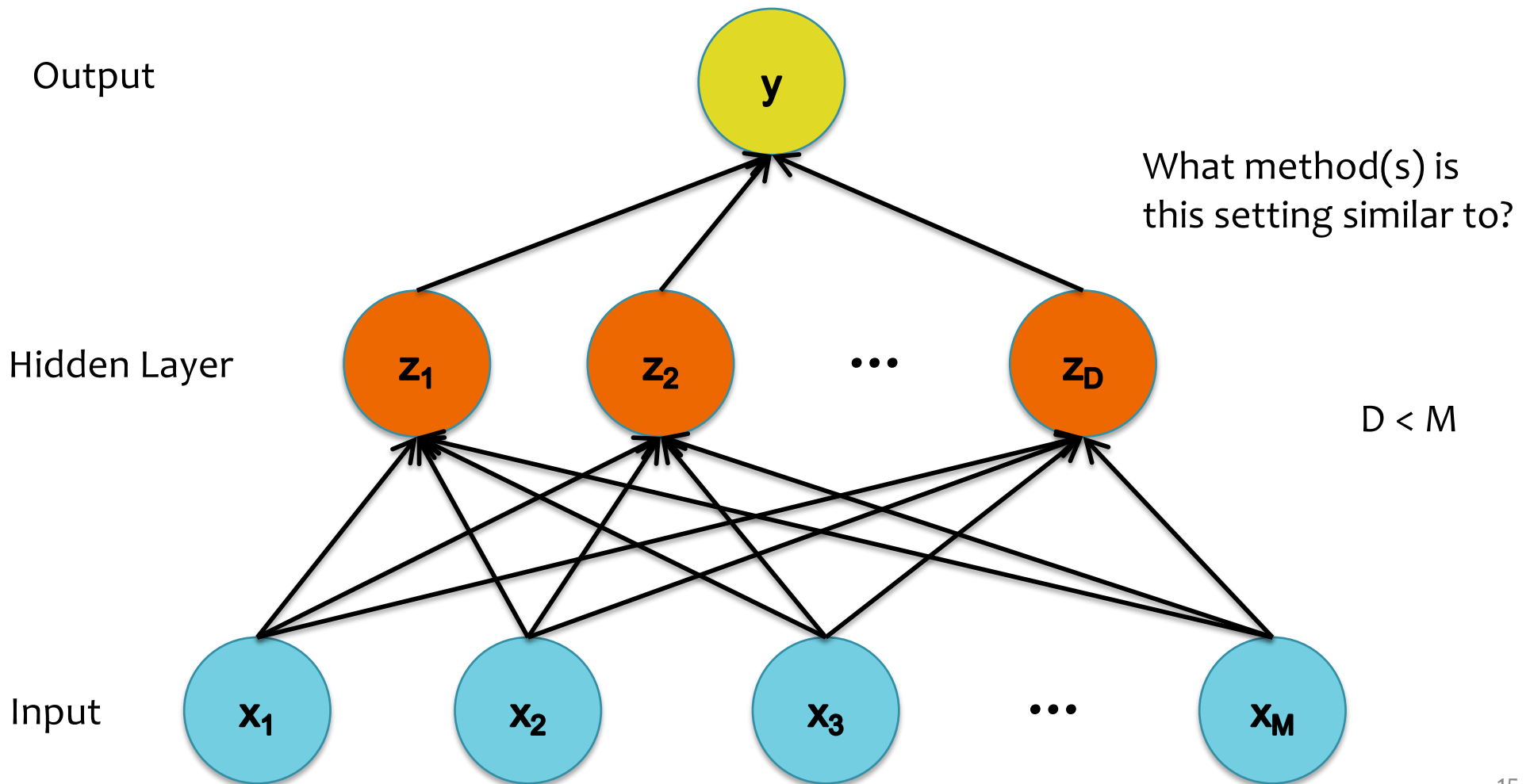
*Q: How many hidden units, D, should we use?*

Output

Hidden Layer

$D = M$

Input

# Building a Neural Net

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

Output

Hidden Layer

Input

D = M

# Building a Neural Net

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

Output

Hidden Layer

Input

$x_1$  $x_2$  $x_3$  ...  $x_M$

$z_1$  $z_2$  ...  $z_D$

y

What method(s) is this setting similar to?

D < M

# Building a Neural Net

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



Output

Hidden Layer

$D > M$

What method(s) is
this setting similar to?

Input

16

# Deeper Networks

*Q: How many layers should we use?*

# Deeper Networks

*Q: How many layers should we use?*

# Deeper Networks

*Q: How many layers should we use?*

# 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.

# Different Levels of Abstraction

- 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

Example from Honglak Lee (NIPS 2010)

# Different Levels of Abstraction

**Face Recognition:**

– Deep Network can build up increasingly higher levels of abstraction

– Lines, parts, regions

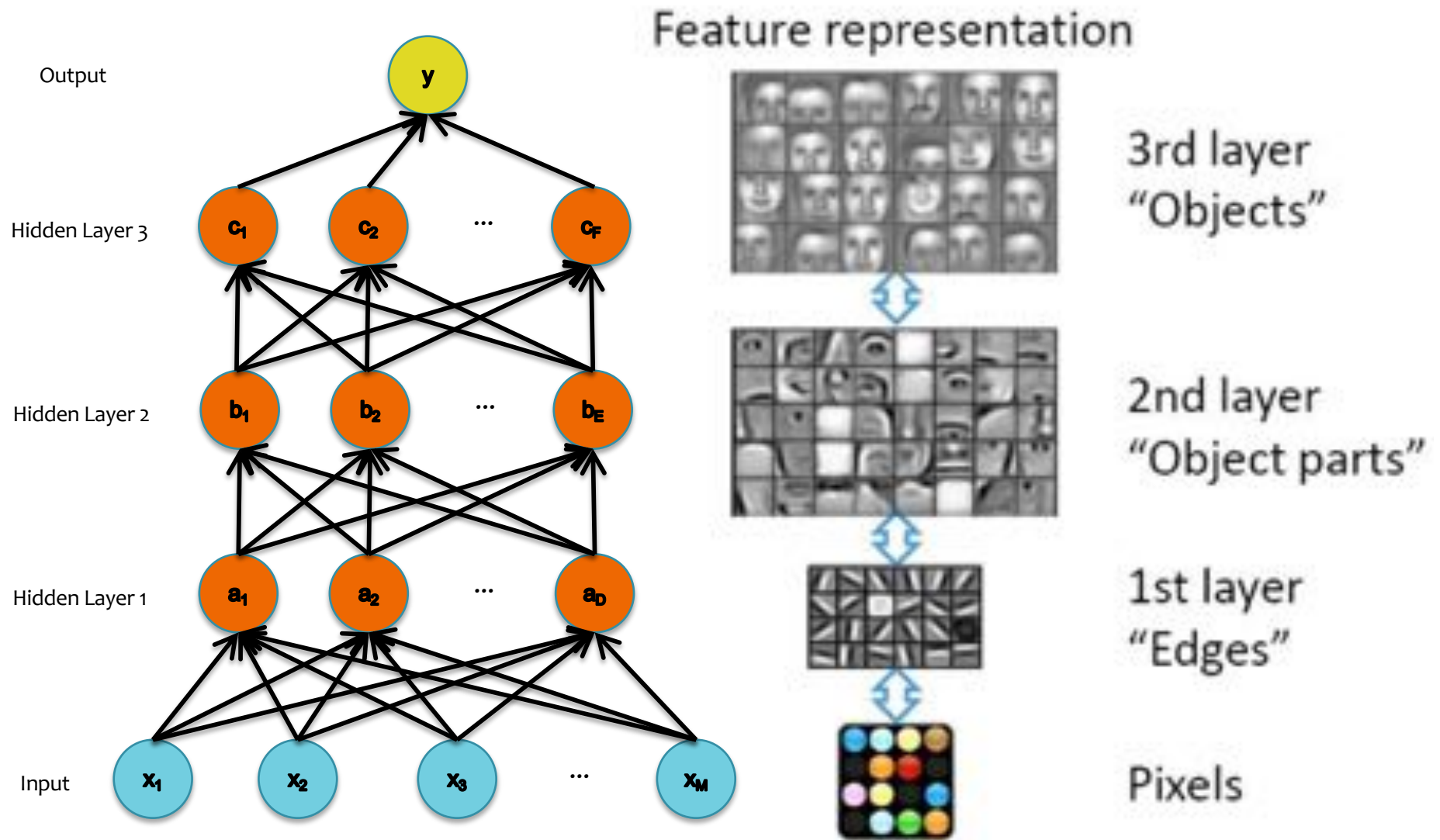Feature representation

3rd layer "Objects"

2nd layer "Object parts"

1st layer "Edges"

Pixels

Example from Honglak Lee (NIPS 2010)

# Different Levels of Abstraction



Output — y

Hidden Layer 3 — $c_1$ $c_2$ ... $c_F$

Hidden Layer 2 — $b_1$ $b_2$ ... $b_E$

Hidden Layer 1 — $a_1$ $a_2$ ... $a_D$

Input — $x_1$ $x_2$ $x_3$ ... $x_M$

Feature representation

3rd layer "Objects"

2nd layer "Object parts"

1st layer "Edges"

Pixels

Example from Honglak Lee (NIPS 2010)

# Activation Functions

Neural Network with sigmoid activation functions



Output

Hidden Layer

Input

**(F) Loss**
$$J = \tfrac{1}{2}(y - y^*)^2$$

**(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 $x_i, \ \forall i$

# Activation Functions



Neural Network with arbitrary nonlinear activation functions

Output

Hidden Layer

Input

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

(E) **Output (nonlinear)**
$y = \sigma(b)$

(D) **Output (linear)**
$b = \sum_{j=0}^{D} \beta_j z_j$

(C) **Hidden (nonlinear)**
$z_j = \sigma(a_j), \ \forall j$

(B) **Hidden (linear)**
$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$

(A) **Input**
Given $x_i, \ \forall i$

28

# 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]

$$net = \sum_{i=0}^{n} w_i x_i$$

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

# Understanding the difficulty of training deep feedforward neural networks

AI Stats 2010



Figure from Glorot & Bentio (2010)

# Activation Functions

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



$$\max(0, z)$$

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

Alternate 2: rectified linear unit

Linear with a cutoff at zero

(Implementation: clip the gradient when you pass zero)

$$net = \sum_{i=0}^{n} w_i\, x_i$$

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

Slide from William Cohen

# 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

# Decision Functions

# Neural Network

**Neural Network for Classification**



Output

Hidden Layer

Input

(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 $x_i, \ \forall i$

34

# Decision Functions

# Neural Network

**Neural Network for Regression**



**Output** — y

**Hidden Layer** — $z_1$, $z_2$, ..., $z_D$

**Input** — $x_1$, $x_2$, $x_3$, ..., $x_M$

(D) **Output (linear)**
$$y = \sum_{j=0}^{D} \beta_j z_j$$

(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 $x_i, \ \forall i$

# 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 | Backward |
|---|---|---|
| Quadratic | $J = \dfrac{1}{2}(y - y^*)^2$ | $\dfrac{dJ}{dy} = y - y^*$ |
| Cross Entropy | $J = y^* \log(y) + (1 - y^*) \log(1 - y)$ | $\dfrac{dJ}{dy} = y^* \dfrac{1}{y} + (1 - y^*) \dfrac{1}{y - 1}$ |

# Objective Functions for NNs

**Cross-entropy vs. Quadratic loss**



Figure 5: 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.

Figure from Glorot & Bentio (2010)

# Multi-Class Output



Output

Hidden Layer

Input

# Multi-Class Output

Softmax:

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

(F) **Loss**
$$J = \sum_{k=1}^{K} y_k^* \log(y_k)$$

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

(D) **Output (linear)**
$$b_k = \sum_{j=0}^{D} \beta_{kj} z_j \ \forall k$$

(C) **Hidden (nonlinear)**
$$z_j = \sigma(a_j), \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

Output  $y_1$ ... $y_K$

Hidden Layer  $z_1$  $z_2$ ... $z_D$

Input  $x_1$  $x_2$  $x_3$ ... $x_M$

# Neural Network Errors

**Question A:** On which of the datasets below could a one-hidden layer neural network achieve zero *classification* error? **Select all that apply.**

**Question B:** On which of the datasets below could a one-hidden layer neural network for *regression* achieve *nearly* zero MSE? **Select all that apply.**

# DECISION BOUNDARY EXAMPLES

# Example #1: Diagonal Band
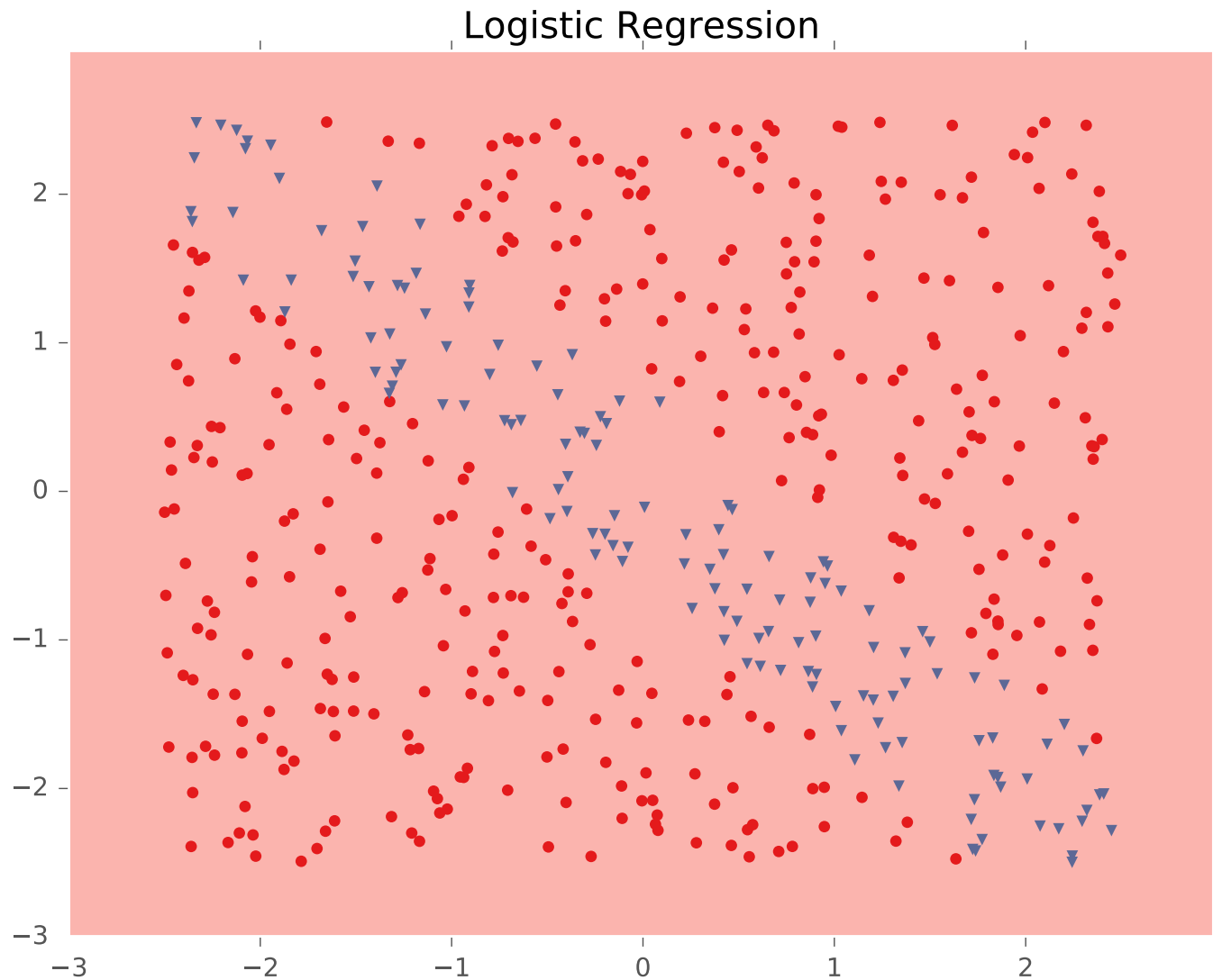


# Example #2: One Pocket



# Example #3: Four Gaussians



# Example #4: Two Pockets

# Example #1: Diagonal Band

# Example #1: Diagonal Band



Logistic Regression

# Example #1: Diagonal Band

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

# Example #1: Diagonal Band

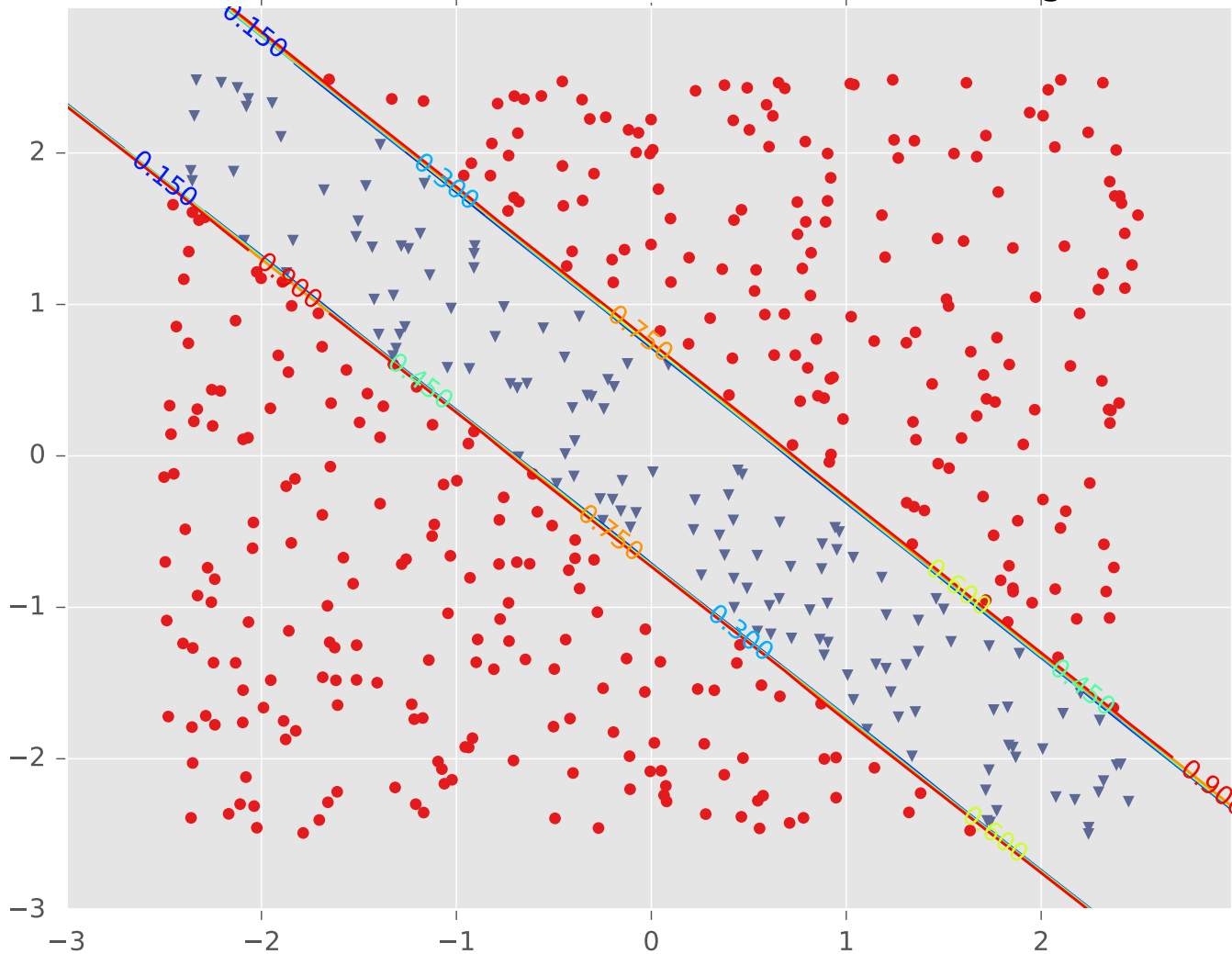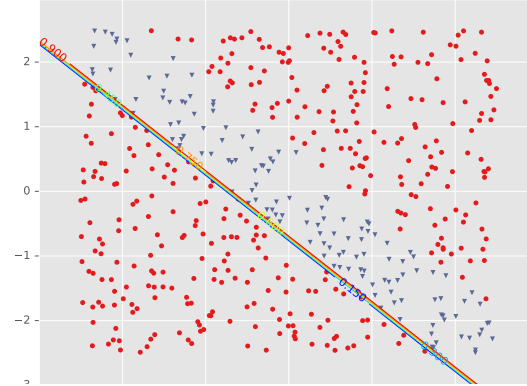LR1 for Tuned Neural Network (hidden=2, activation=logistic)

# Example #1: Diagonal Band

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

# Example #1: Diagonal Band



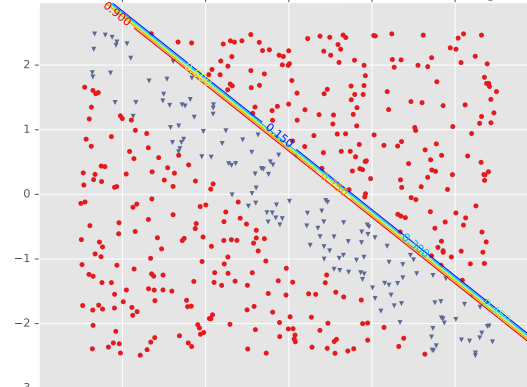Tuned Neural Network (hidden=2, activation=logistic)
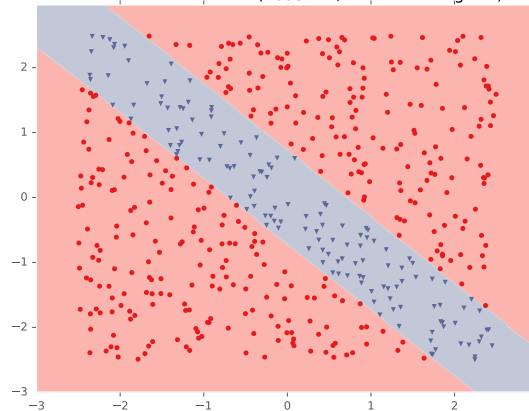
# Example #1: Diagonal Band



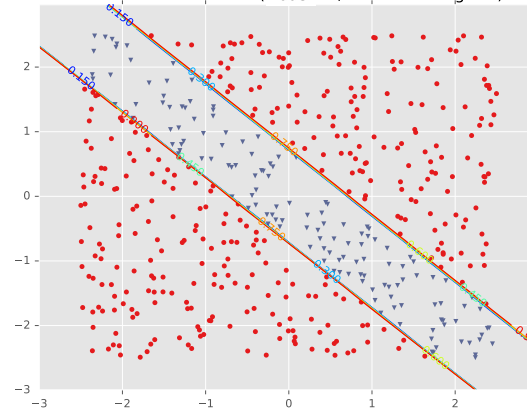LR1 for Tuned Neural Network (hidden=2, activation=logistic)

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

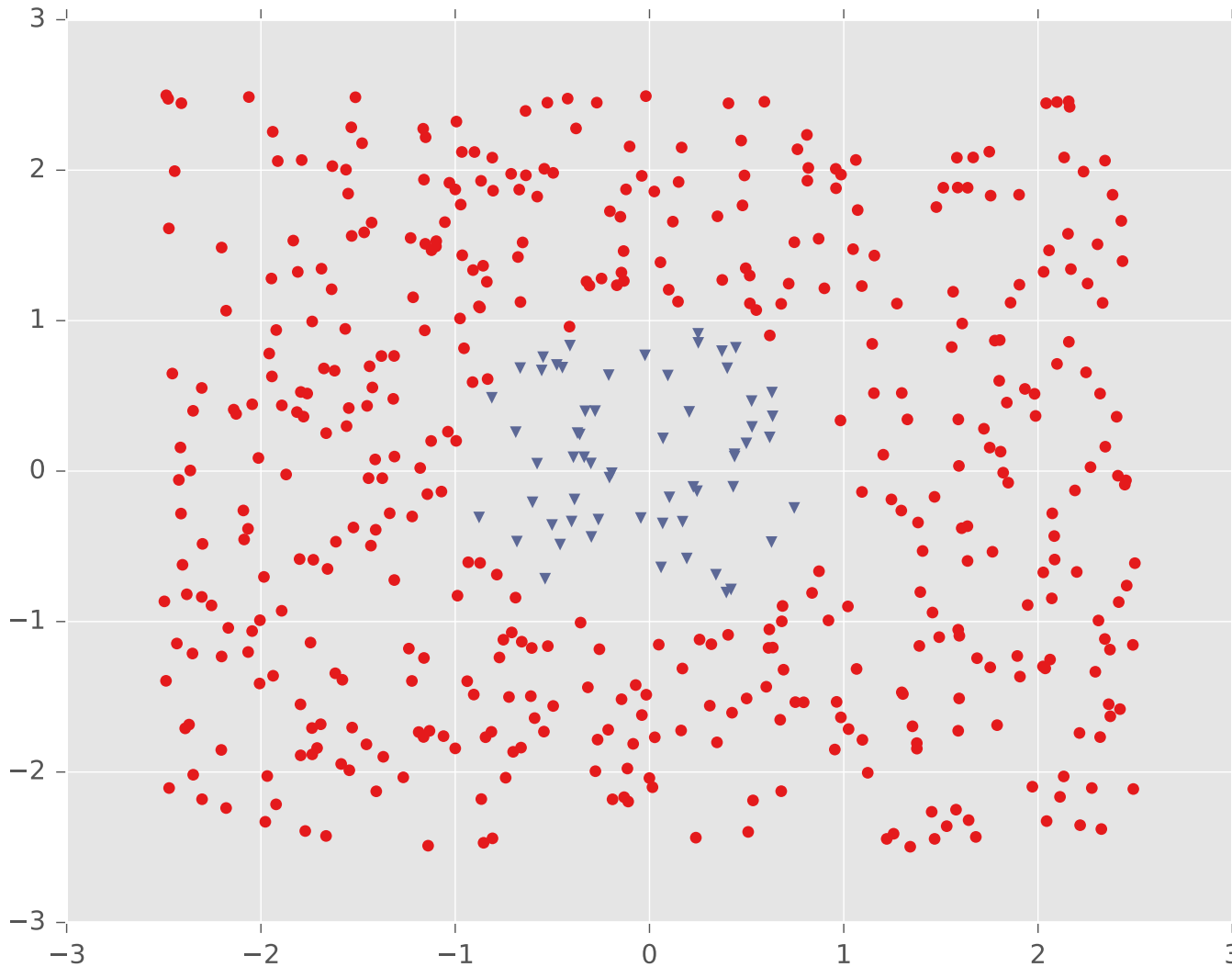Tuned Neural Network (hidden=2, activation=logistic)

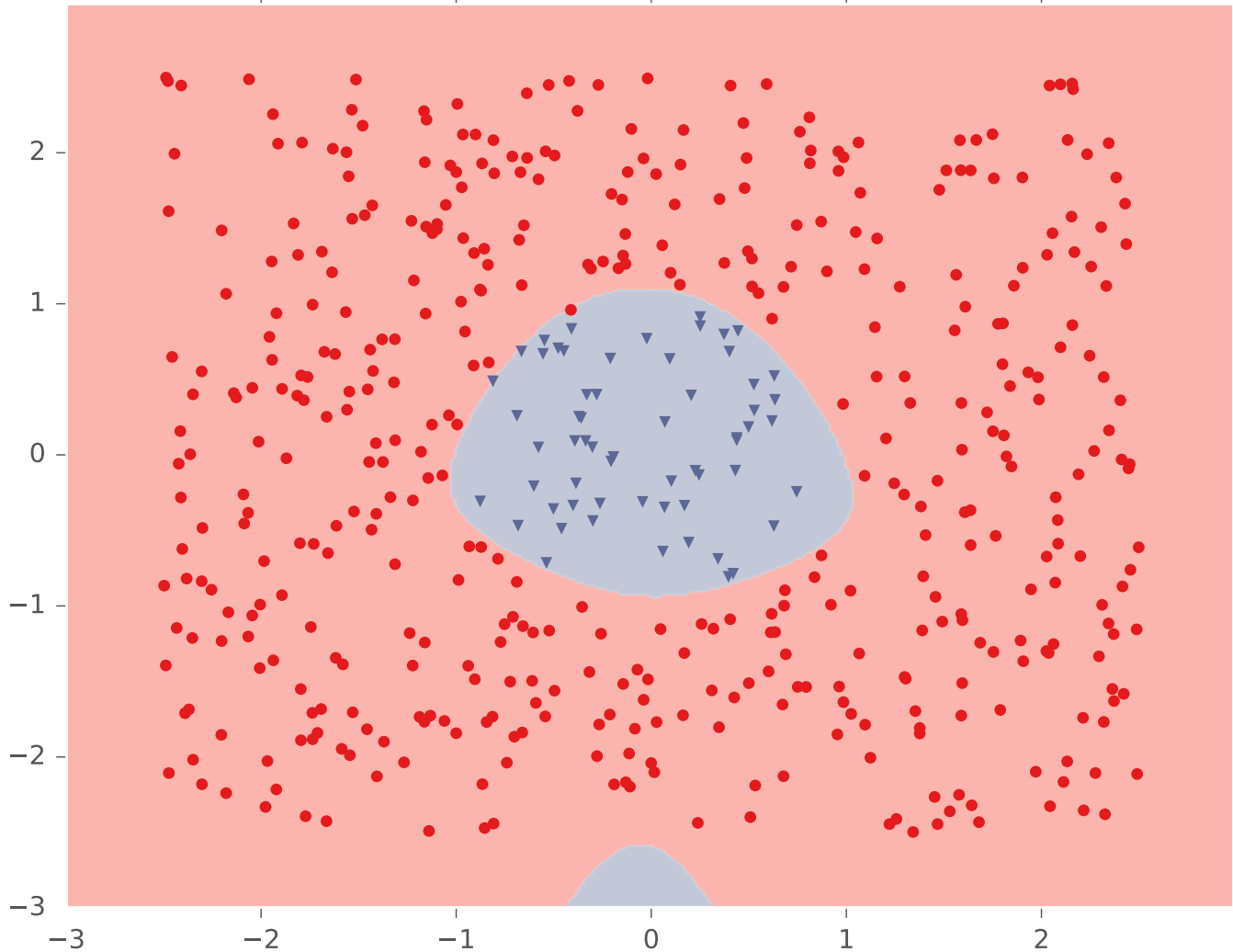Tuned Neural Network (hidden=2, activation=logistic)

49

# Example #2: One Pocket
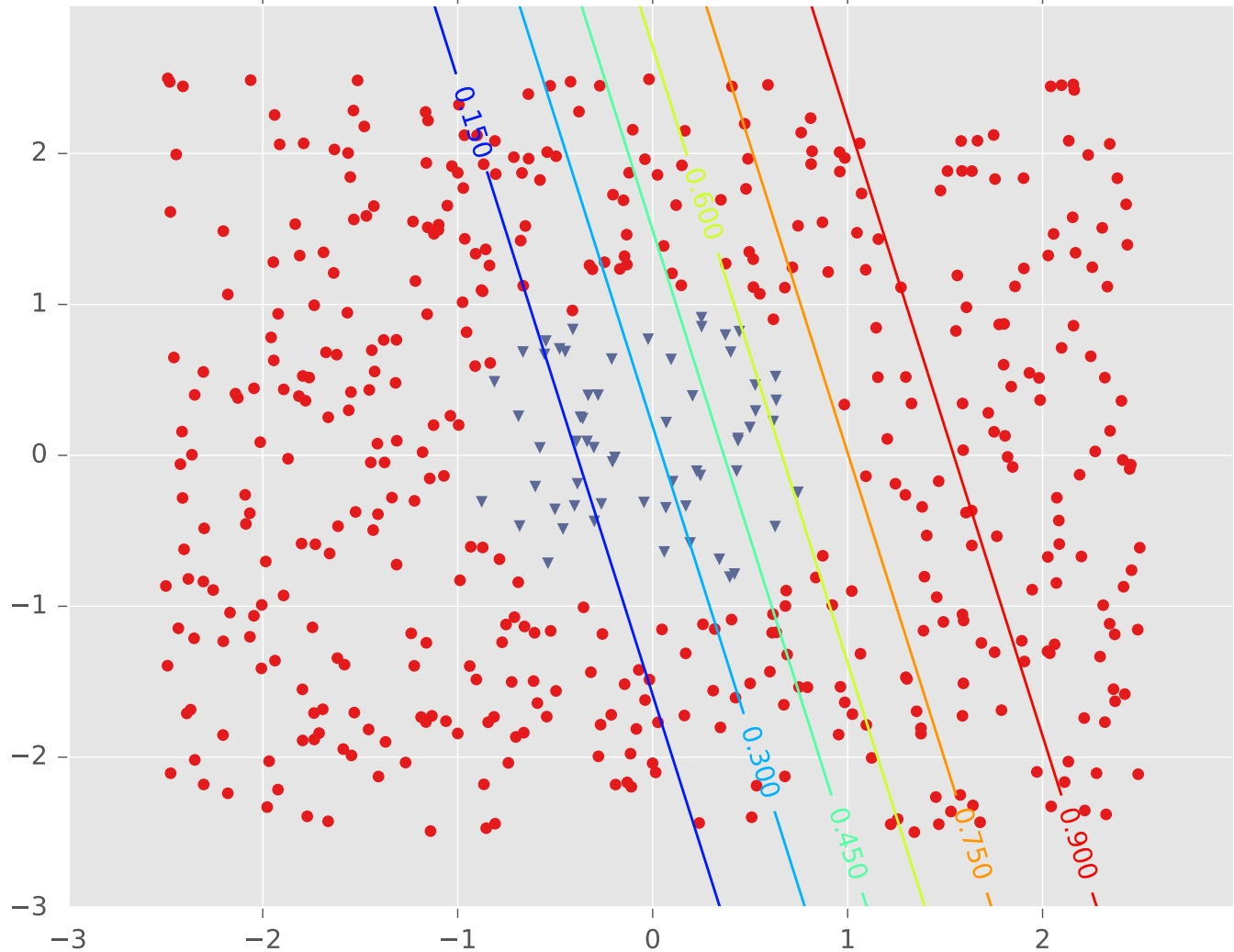
# Example #2: One Pocket



Logistic Regression

# Example #2: One Pocket

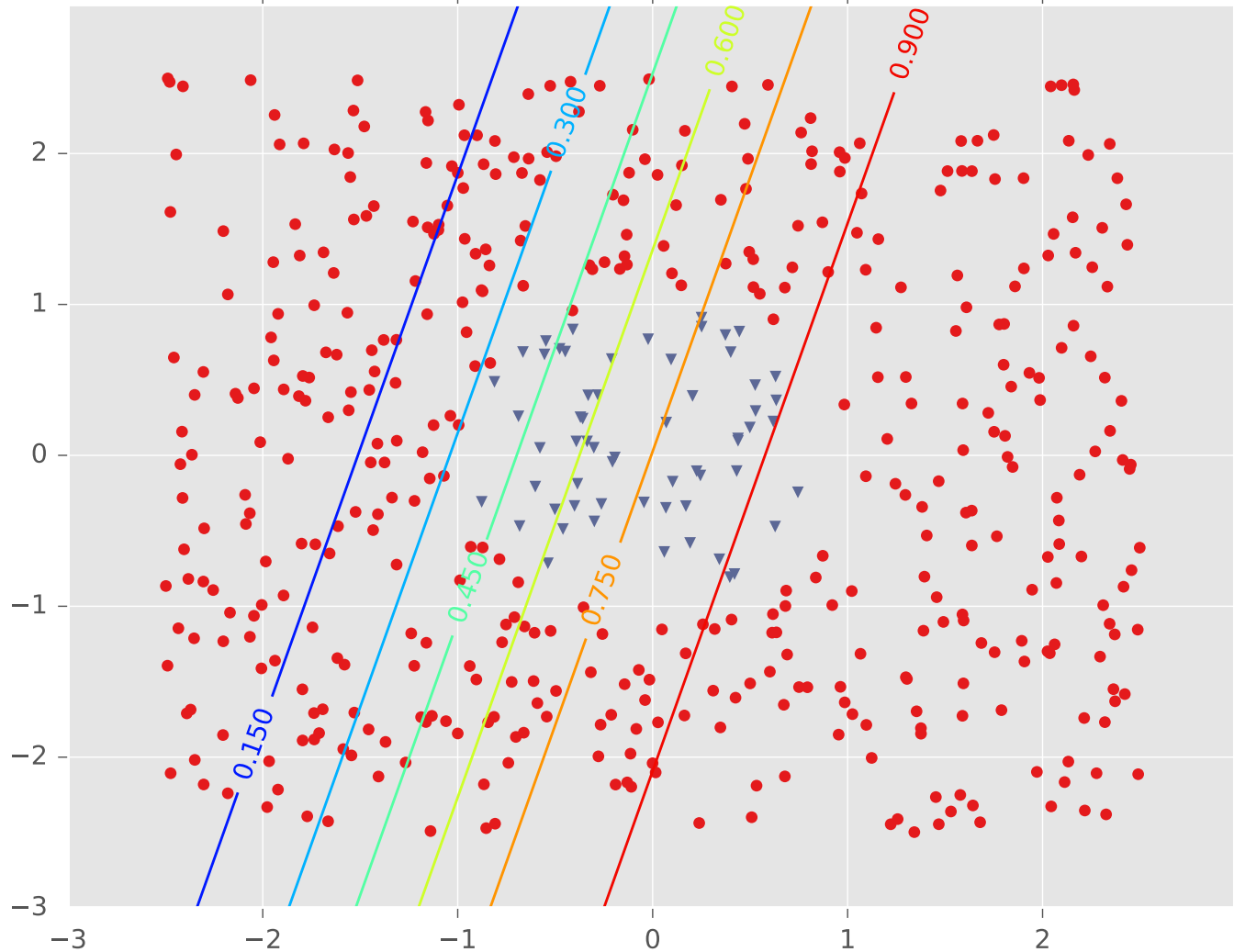Tuned Neural Network (hidden=3, activation=logistic)

# Example #2: One Pocket



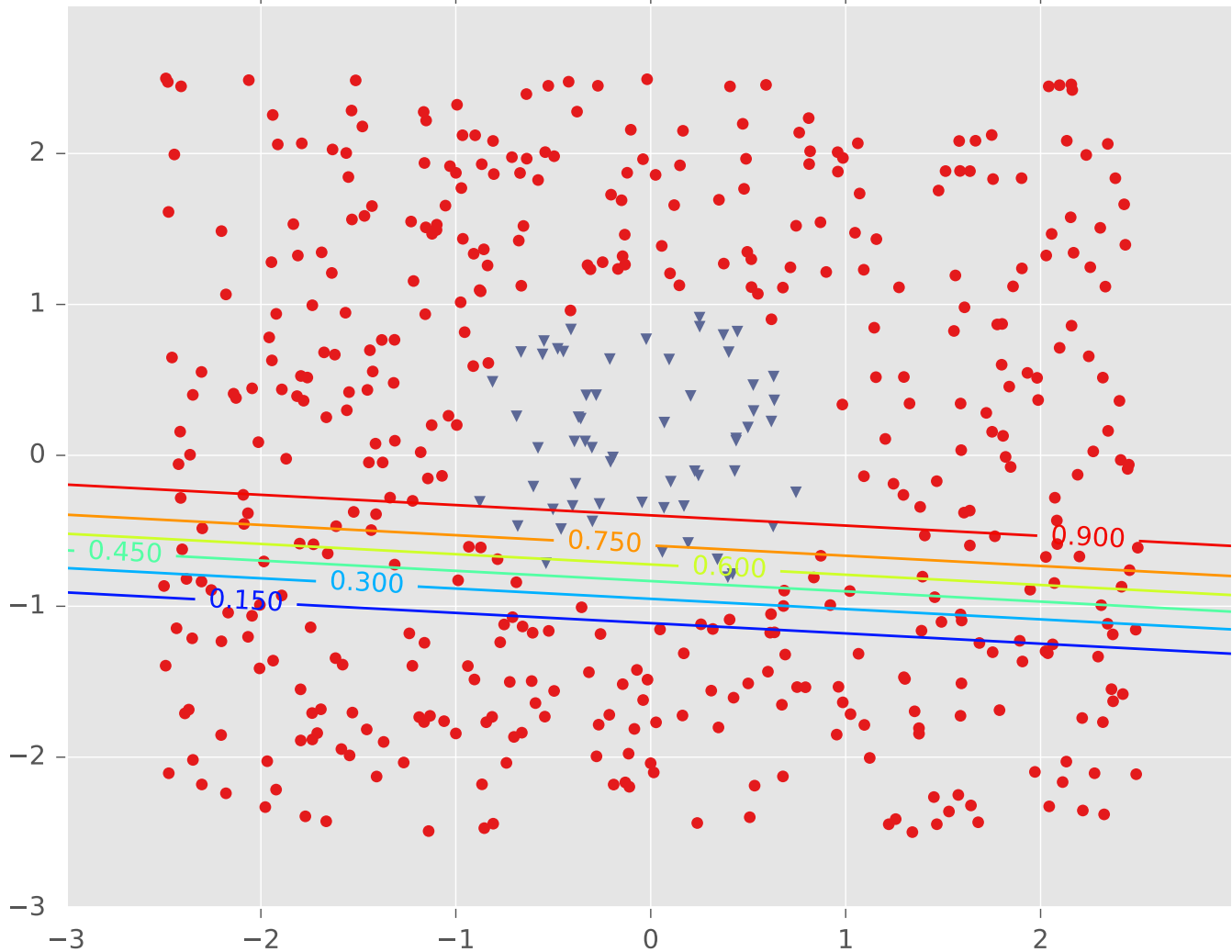LR1 for Tuned Neural Network (hidden=3, activation=logistic)

# 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)
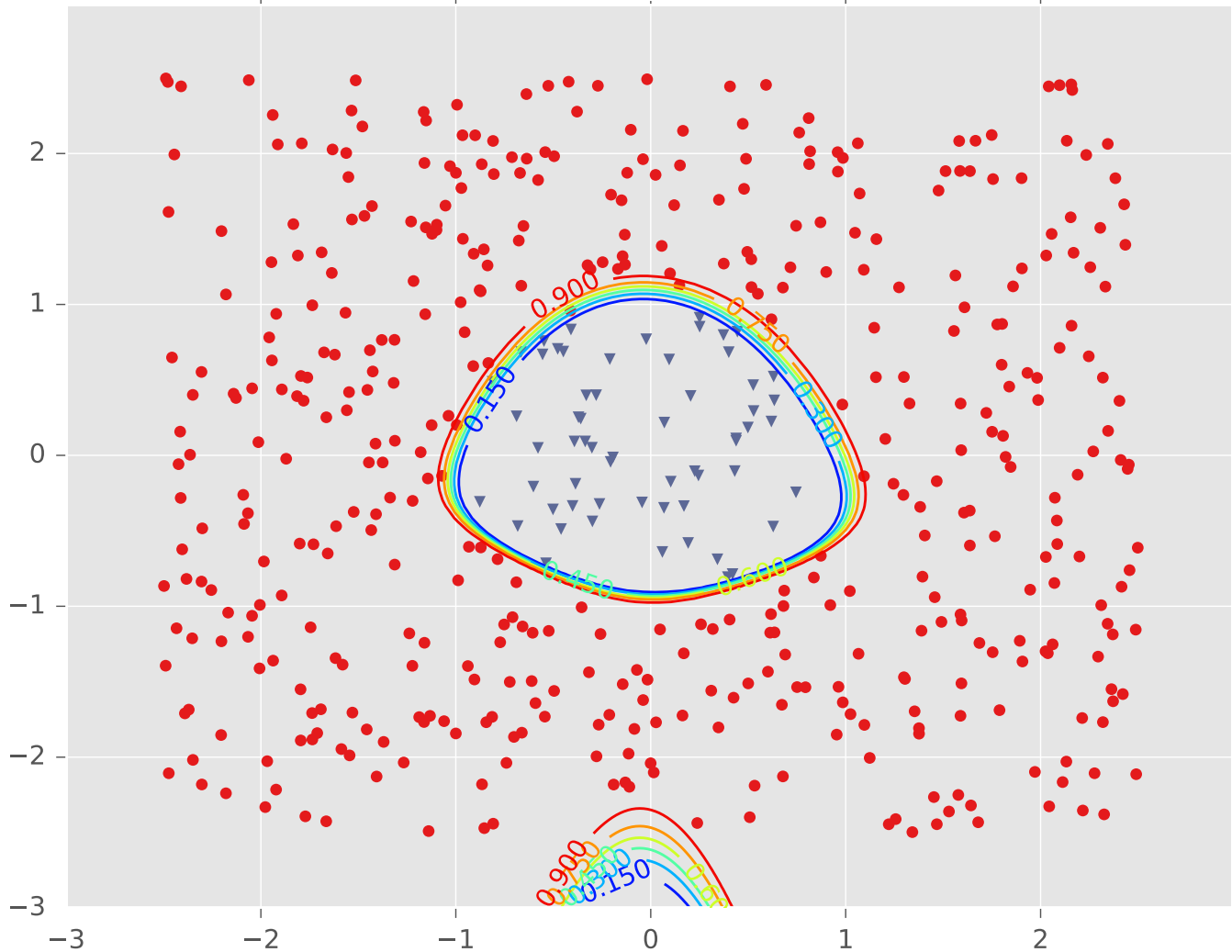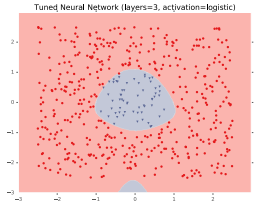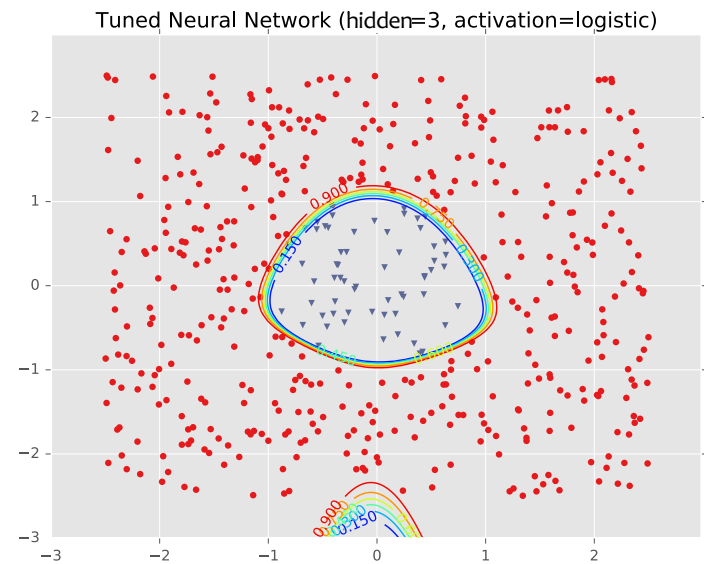
# Example #2: One Pocket


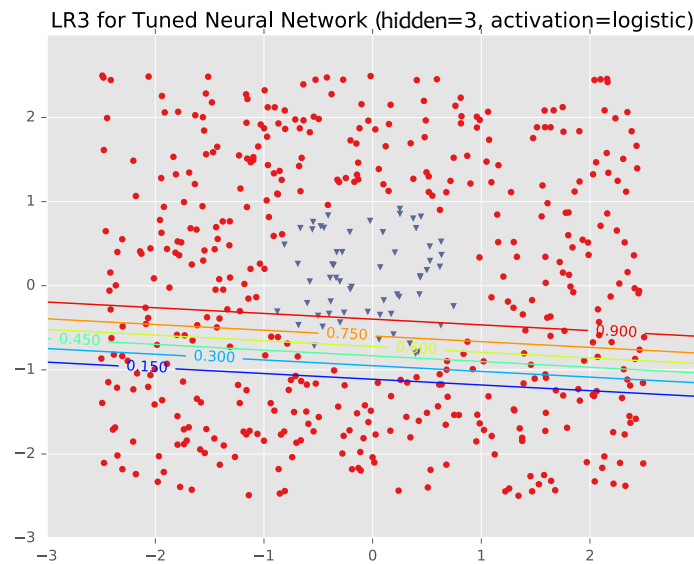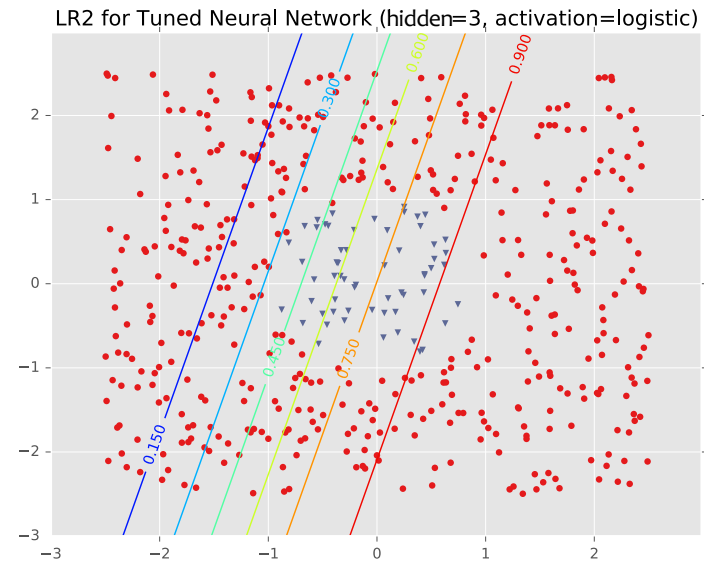LR1 for Tuned Neural Network (hidden=3, activation=logistic)


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



Logistic Regression

# Example #3: Four Gaussians

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

# Example #3: Four Gaussians

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

# Example #3: Four Gaussians



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

# Example #3: Four Gaussians



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

# Example #3: Four Gaussians



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

# Example #4: Two Pockets

# Example #4: Two Pockets



Logistic Regression

# Example #4: Two Pockets

SVM (kernel=linear)

# Example #4: Two Pockets

SVM (kernel=rbf, gamma=80.000000)

# Example #4: Two Pockets

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

# Example #4: Two Pockets



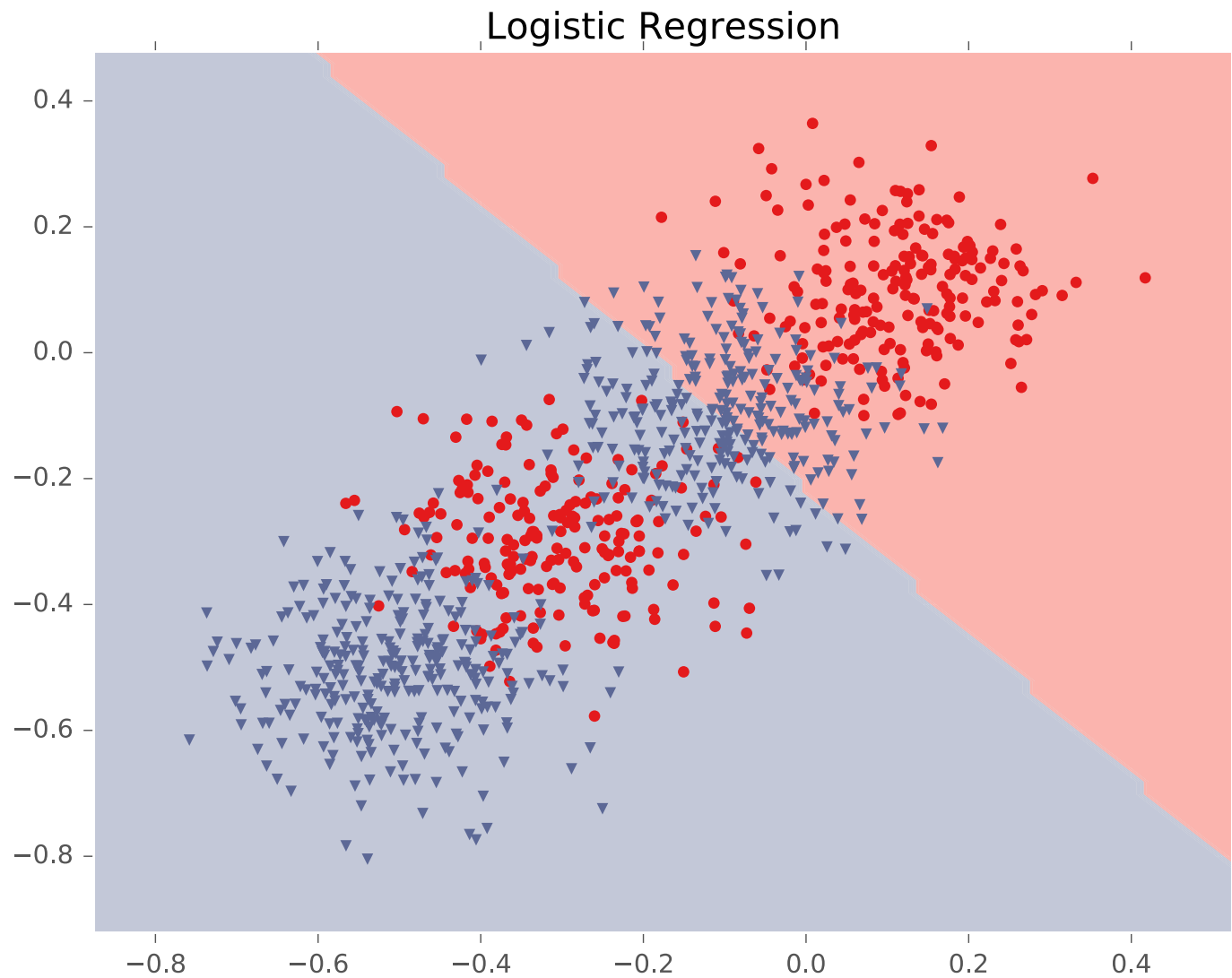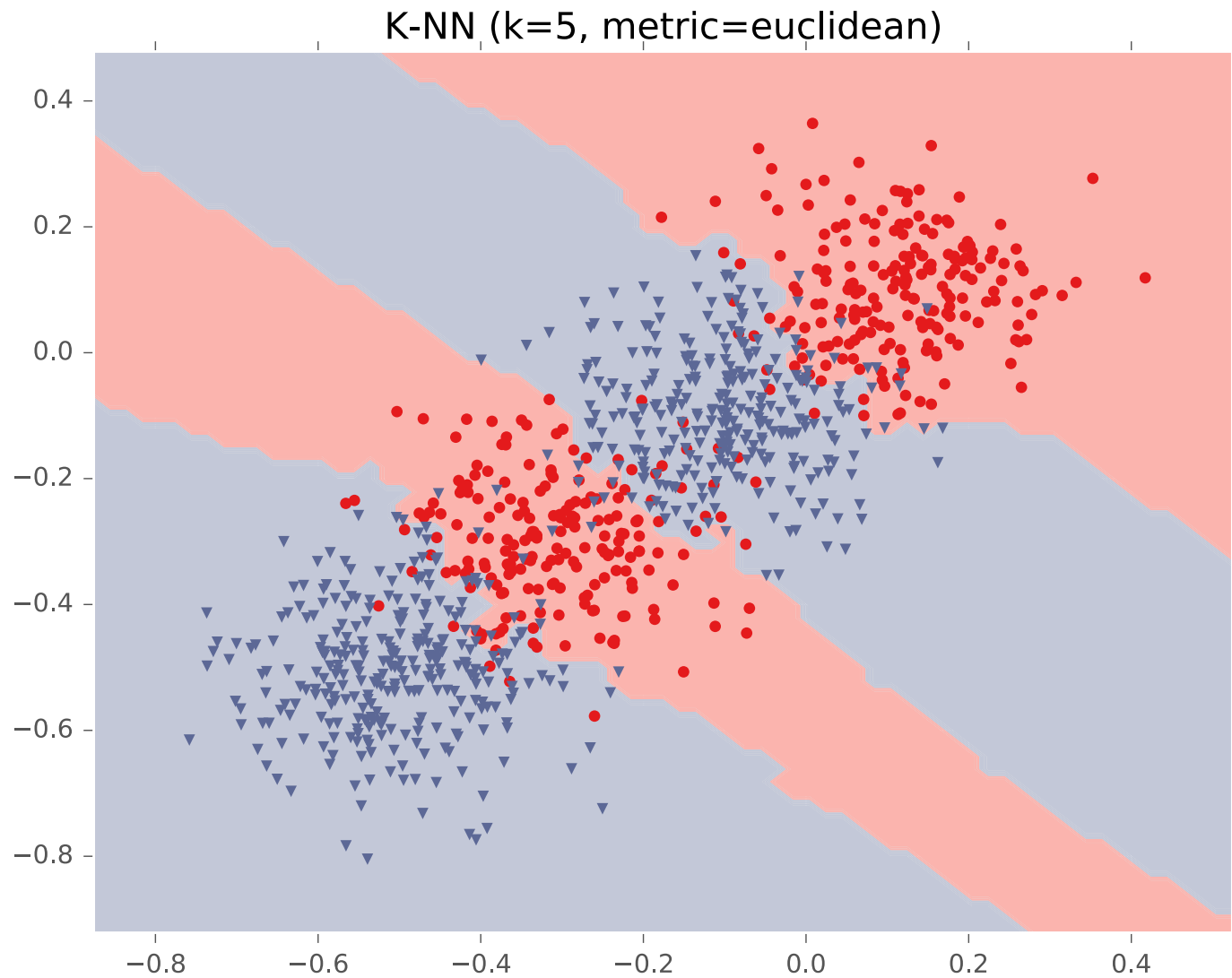Tuned Neural Network (hidden=2, activation=logistic)

# 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
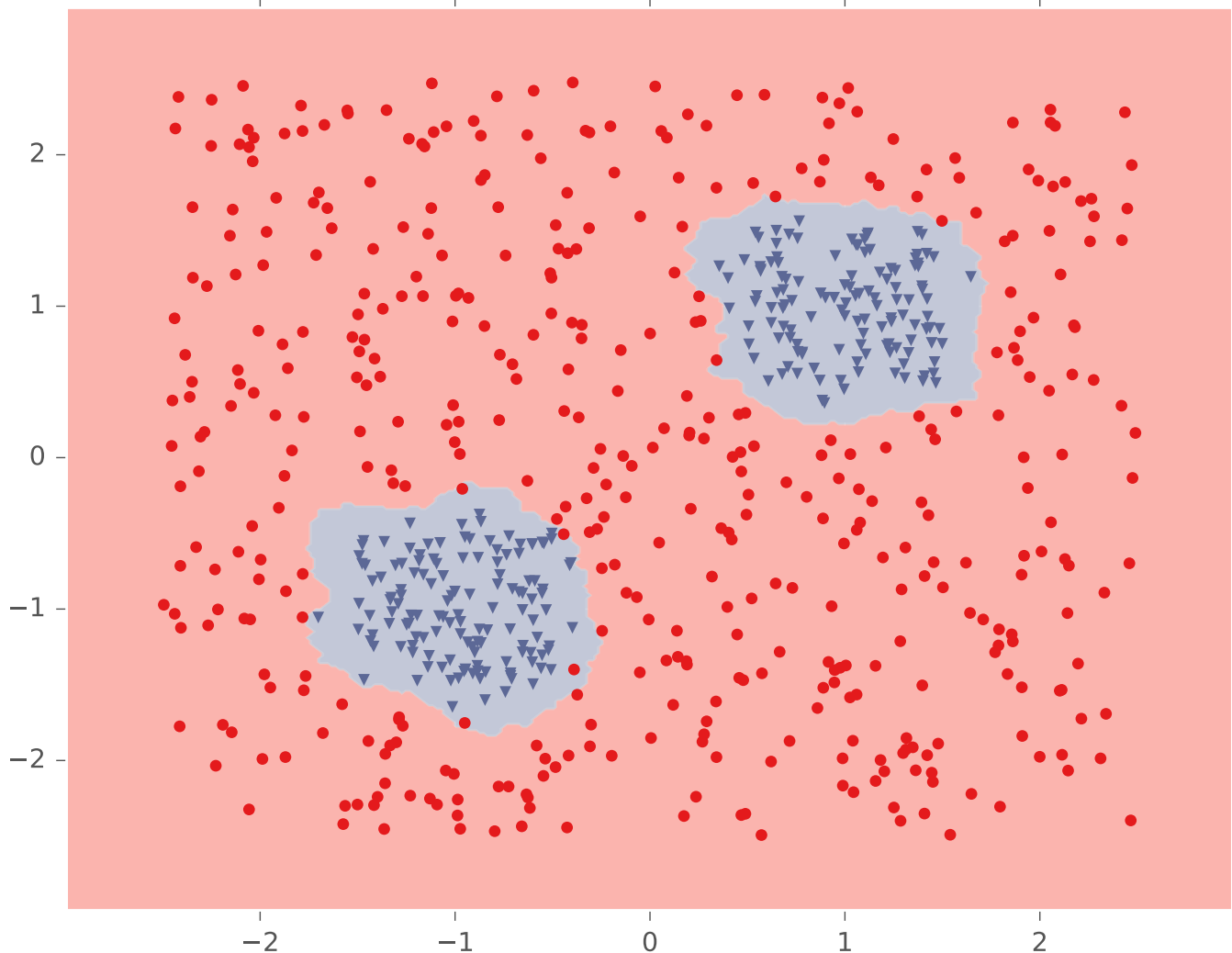- Compare and contrast feature engineering with learning features
- Identify (some of) the options available when designing the architecture of a neural network
- Implement a feed-forward neural network

Computing Gradients

# DIFFERENTIATION

# Background

# A Recipe for Machine Learning

1. Given training data:
$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N}$$

2. Choose each of these:
  - Decision function
$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

  - Loss function
$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

3. Define goal:
$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

4. Train with SGD:

(take small steps opposite the gradient)

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

# Approaches to Differentiation

- **Question 1:**
  When can we compute the gradients for an arbitrary neural network?

- **Question 2:**
  When can we make the gradient computation efficient?

# Approaches to Differentiation

1. **Finite Difference Method**
   - Pro: Great for testing implementations of backpropagation
   - Con: Slow for high dimensional inputs / outputs
   - Required: Ability to call the function f($\mathbf{x}$) on any input $\mathbf{x}$

2. **Symbolic Differentiation**
   - Note: The method you learned in high-school
   - Note: Used by Mathematica / Wolfram Alpha / Maple
   - Pro: Yields easily interpretable derivatives
   - Con: Leads to exponential computation time if not carefully implemented
   - Required: Mathematical expression that defines f($\mathbf{x}$)

3. **Automatic Differentiation - Reverse Mode**
   - Note: Called *Backpropagation* when applied to Neural Nets
   - Pro: Computes partial derivatives of one output f($\mathbf{x}$)$_i$ with respect to all inputs $x_j$ in time proportional to computation of f($\mathbf{x}$)
   - Con: Slow for high dimensional outputs (e.g. vector-valued functions)
   - Required: Algorithm for computing f($\mathbf{x}$)

4. **Automatic Differentiation - Forward Mode**
   - Note: Easy to implement. Uses dual numbers.
   - Pro: Computes partial derivatives of all outputs f($\mathbf{x}$)$_i$ with respect to one input $x_j$ in time proportional to computation of f($\mathbf{x}$)
   - Con: Slow for high dimensional inputs (e.g. vector-valued $\mathbf{x}$)
   - Required: Algorithm for computing f($\mathbf{x}$)

$$\text{Given } f : \mathbb{R}^A \to \mathbb{R}^B, f(\mathbf{x})$$

$$\text{Compute } \frac{\partial f(\mathbf{x})_i}{\partial x_j} \forall i, j$$
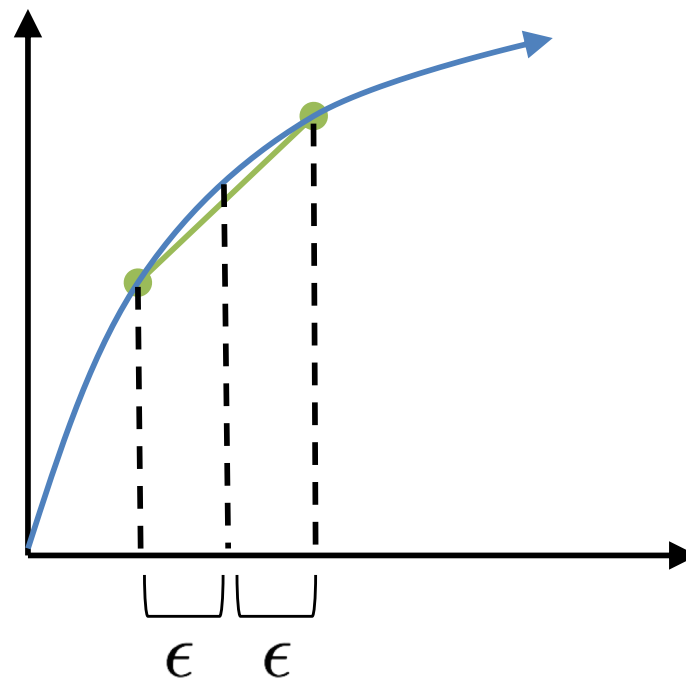
The *centered* finite difference approximation is:

$$\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) \approx \frac{(J(\boldsymbol{\theta} + \epsilon \cdot \boldsymbol{d}_i) - J(\boldsymbol{\theta} - \epsilon \cdot \boldsymbol{d}_i))}{2\epsilon} \qquad (1)$$

where $d_i$ is a 1-hot vector consisting of all zeros except for the $i$th entry of $d_i$, which has value 1.

**Notes:**

- Suffers from issues of floating point precision, in practice
- Typically only appropriate to use on small examples with an appropriately chosen epsilon

$\epsilon$    $\epsilon$

# Symbolic Differentiation

**Differentiation Quiz #1:**

Suppose x = 2 and z = 3, what are dy/dx and dy/dz for the function below? **Round your answer to the nearest integer**.

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

**Answer:** *Answers below are in the form [dy/dx, dy/dz]*

A.  [42, -72]
B.  [72, -42]
C.  [100, 127]
D.  [127, 100]

E.  [1208, 810]
F.  [810, 1208]
G.  [1505, 94]
H.  [94, 1505]

# Symbolic Differentiation

## Differentiation Quiz #2:

A neural network with 2 hidden layers can be written as:

$$y = \sigma(\boldsymbol{\beta}^T \sigma((\boldsymbol{\alpha}^{(2)})^T \sigma((\boldsymbol{\alpha}^{(1)})^T \mathbf{x})))$$

where $y \in \mathbb{R}$, $\mathbf{x} \in \mathbb{R}^{D^{(0)}}$, $\boldsymbol{\beta} \in \mathbb{R}^{D^{(2)}}$ and $\boldsymbol{\alpha}^{(i)}$ is a $D^{(i)} \times D^{(i-1)}$ matrix. Nonlinear functions are applied elementwise:

$$\sigma(\mathbf{a}) = [\sigma(a_1), \ldots, \sigma(a_K)]^T$$

Let $\sigma$ be sigmoid: $\sigma(a) = \frac{1}{1+exp-a}$

What is $\frac{\partial y}{\partial \beta_j}$ and $\frac{\partial y}{\partial \alpha_j^{(i)}}$ for all $i, j$.

# CHAIN RULE

# Training Chain Rule

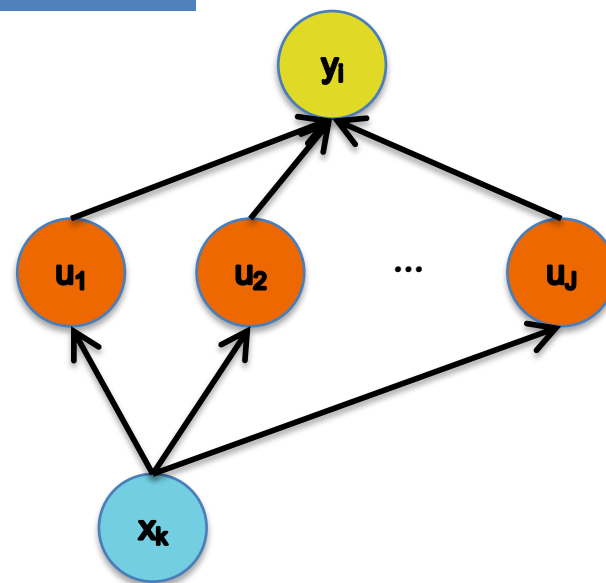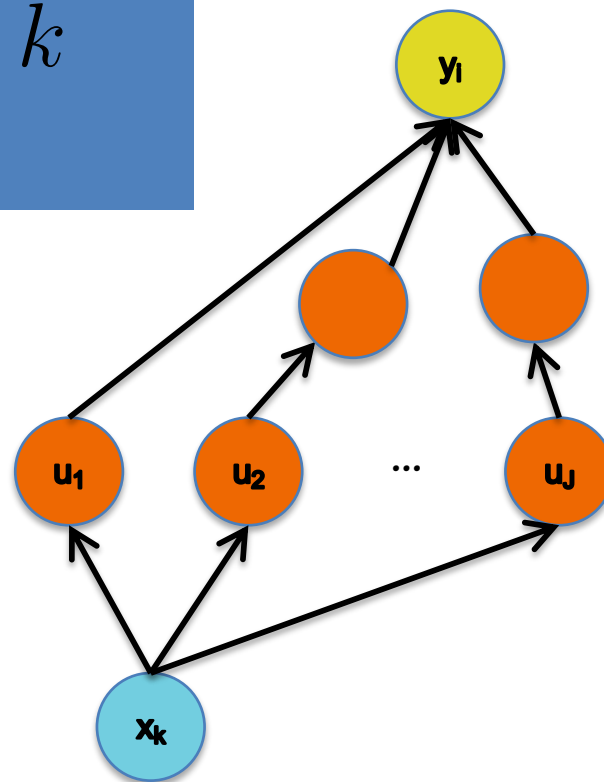*Chalkboard*

– Chain Rule of Calculus

# Chain Rule

**Given:** $y = g(u)$ and $u = h(x)$.

**Chain Rule:**

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j} \frac{du_j}{dx_k}, \quad \forall i, k$$

# Chain Rule

**Given:** $y = g(u)$ and $u = h(x)$.

**Chain Rule:**

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j} \frac{du_j}{dx_k}, \quad \forall i, k$$
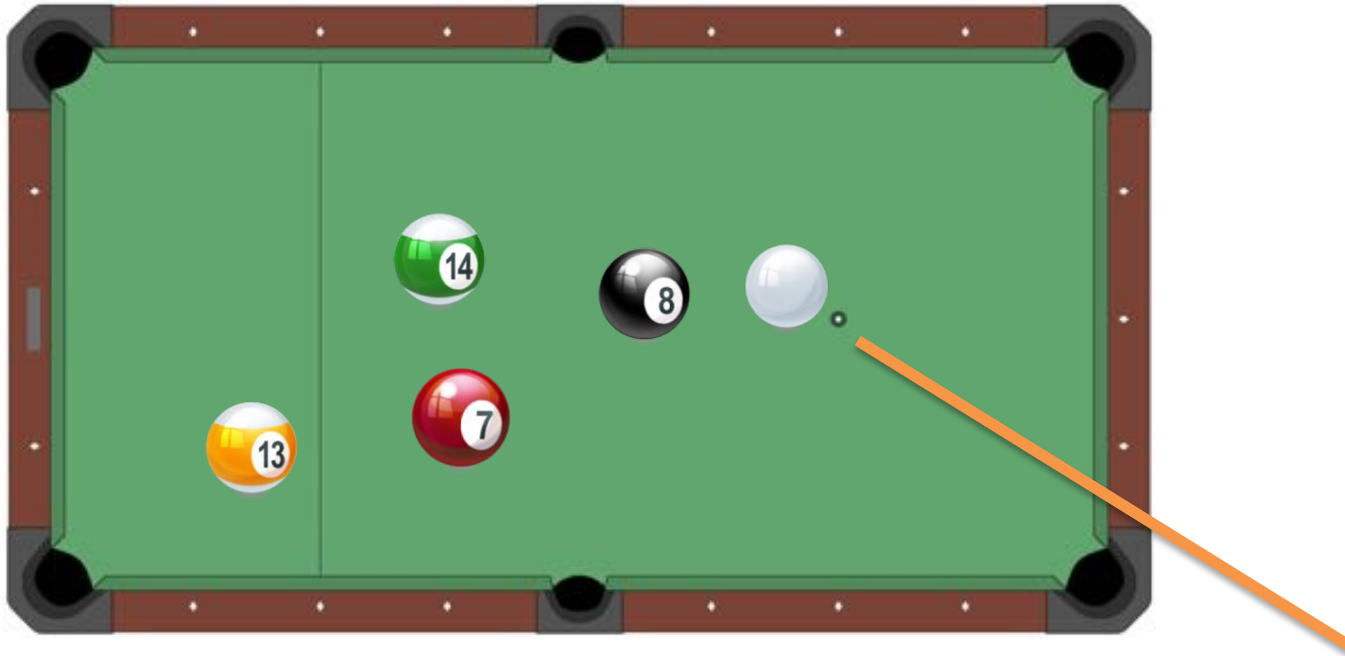
**Backpropagation** is just repeated application of the **chain rule** from Calculus 101.

Intuitions

# BACKPROPAGATION

# Error Back-Propagation

# Error Back-Propagation

Slide from (Stoyanov & Eisner, 2012)

# Error Back-Propagation

# Error Back-Propagation

# Error Back-Propagation

92

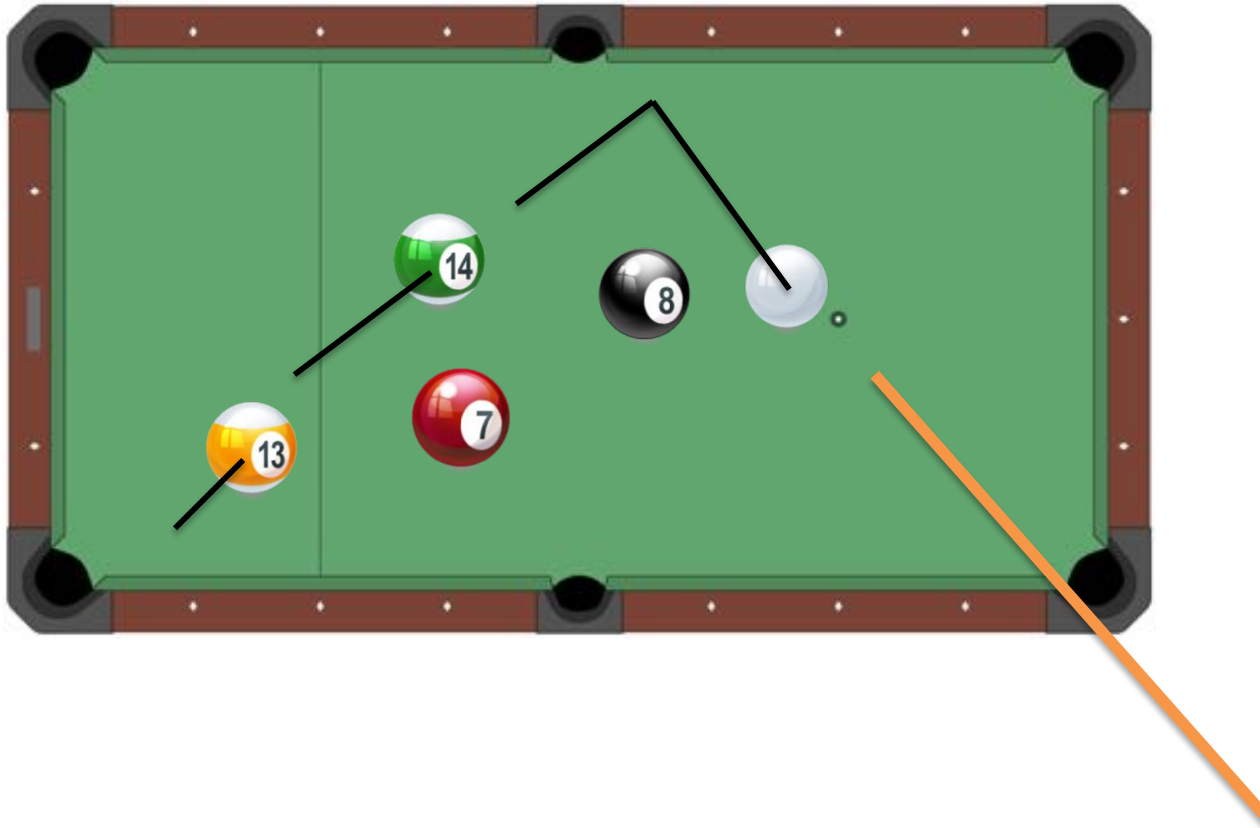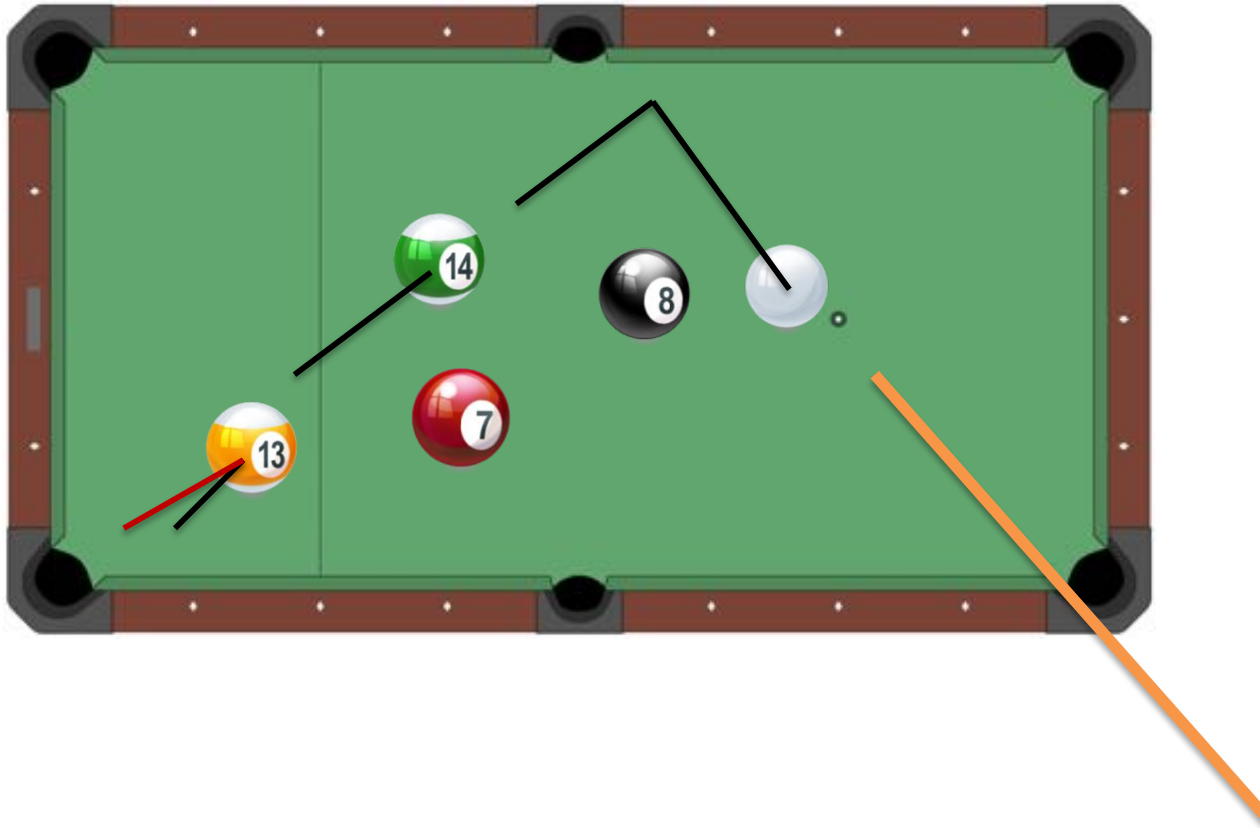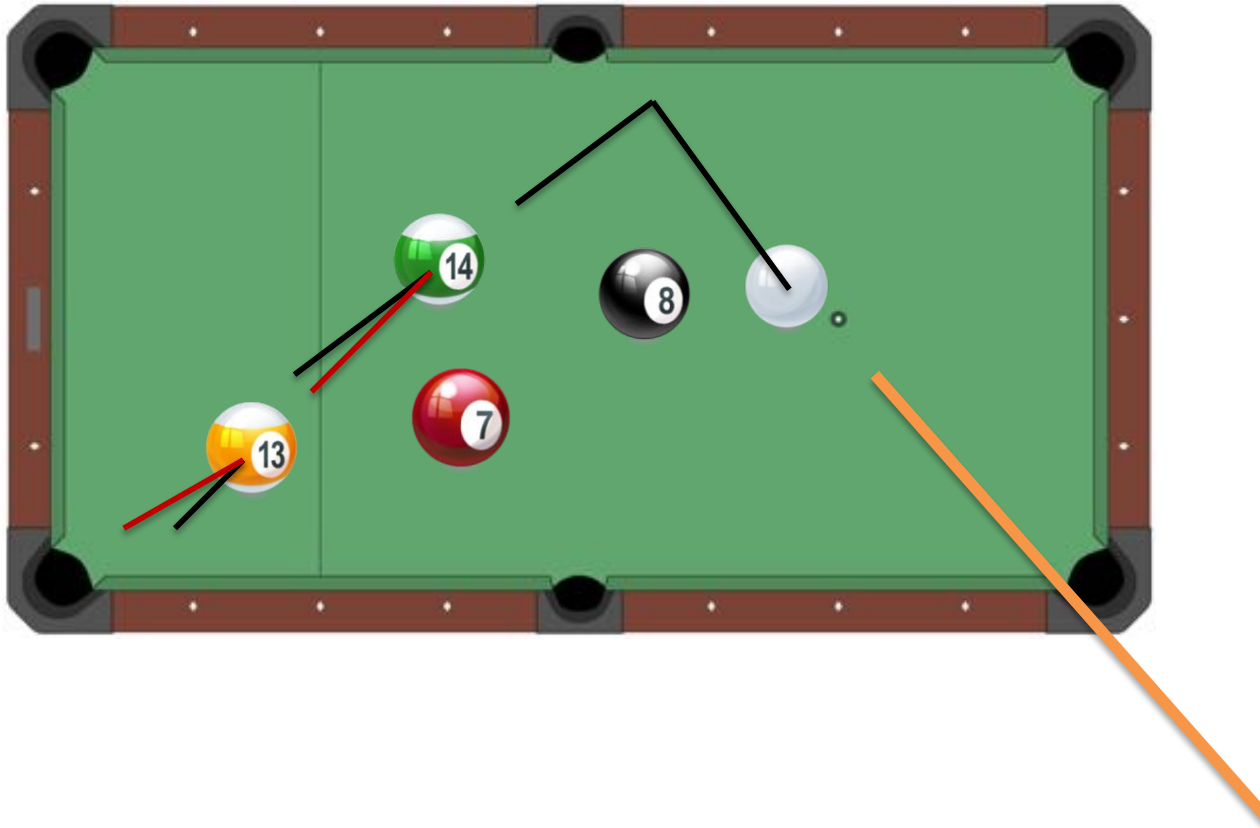# Error Back-Propagation
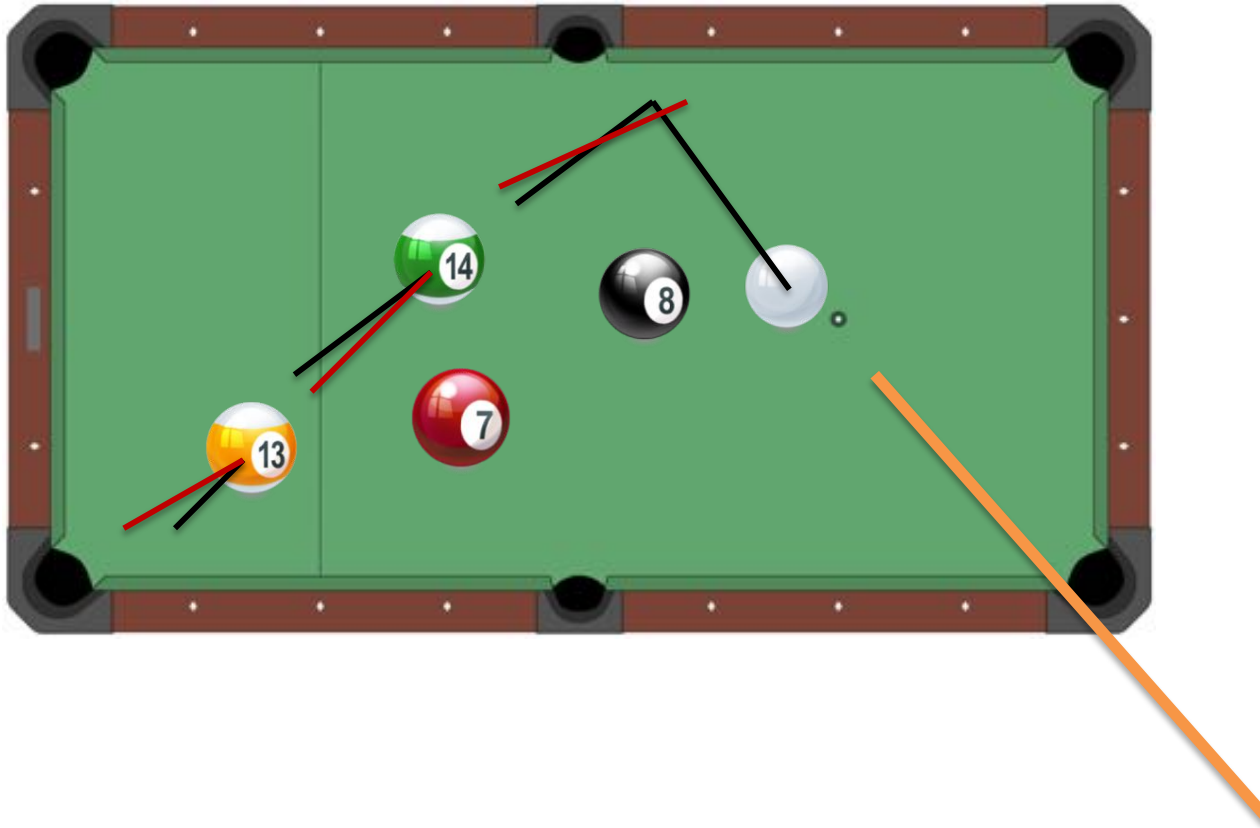
Slide from (Stoyanov & Eisner, 2012)

# Error Back-Propagation

# Error Back-Propagation

# Error Back-Propagation
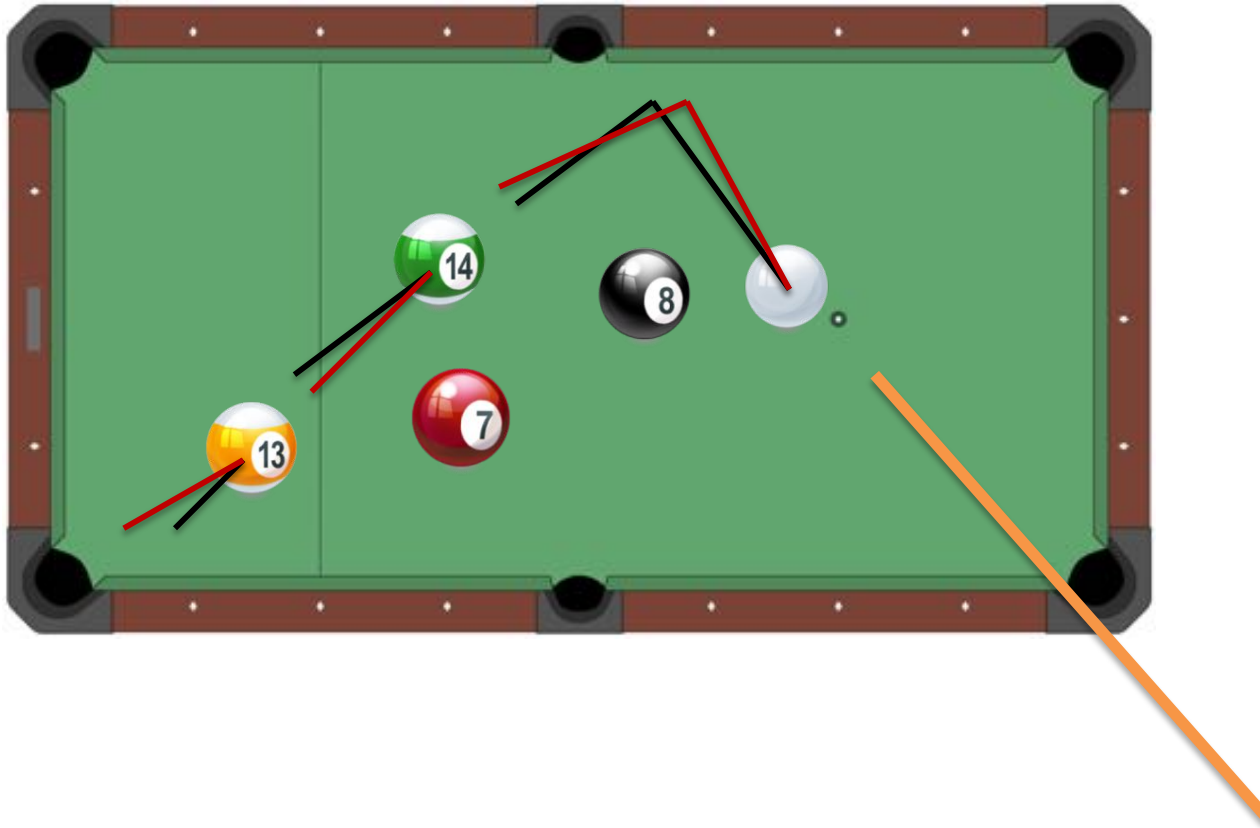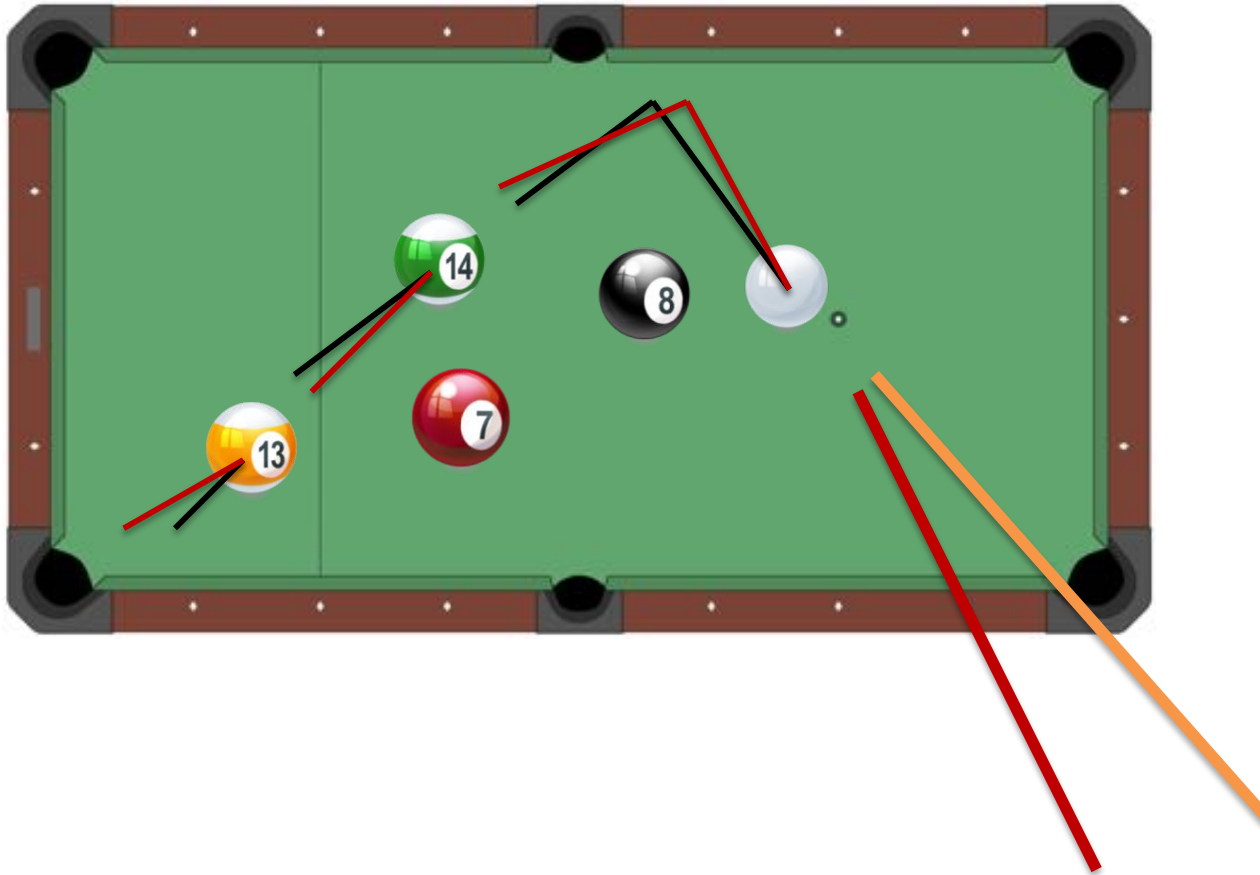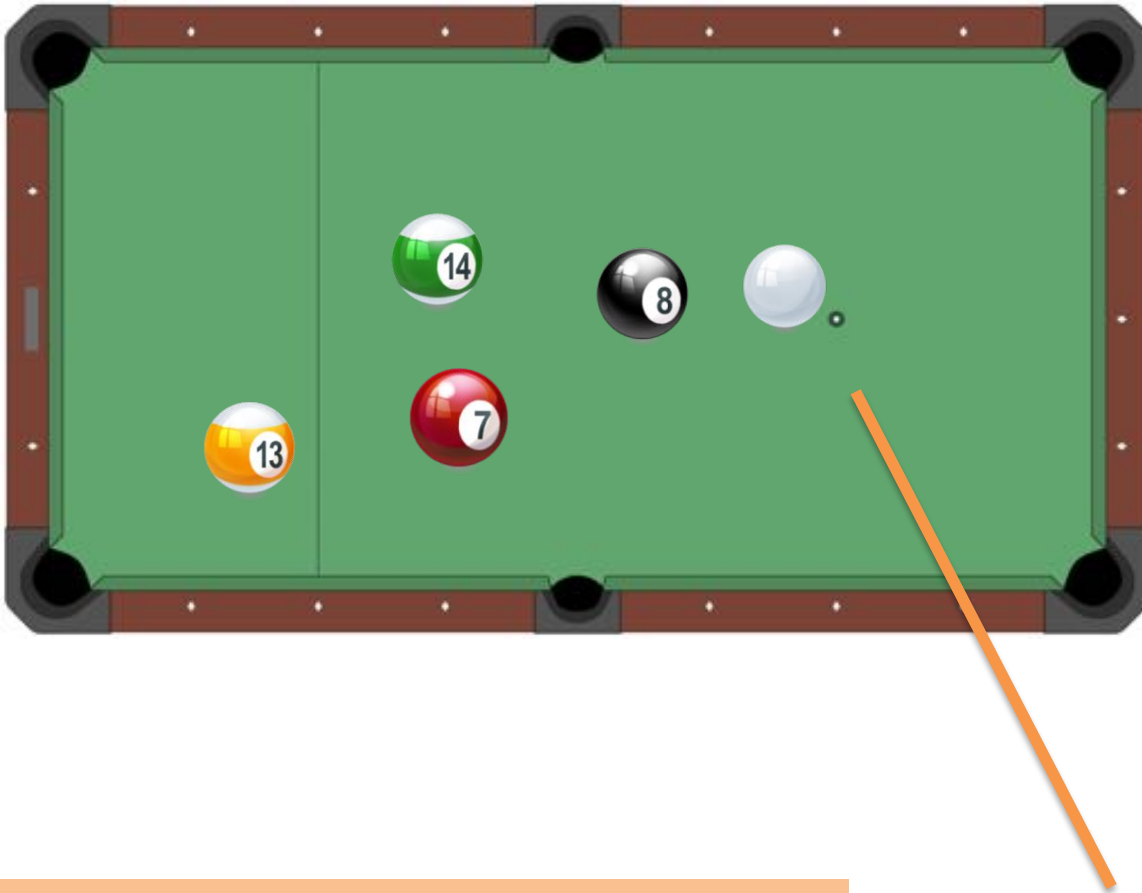
# Error Back-Propagation



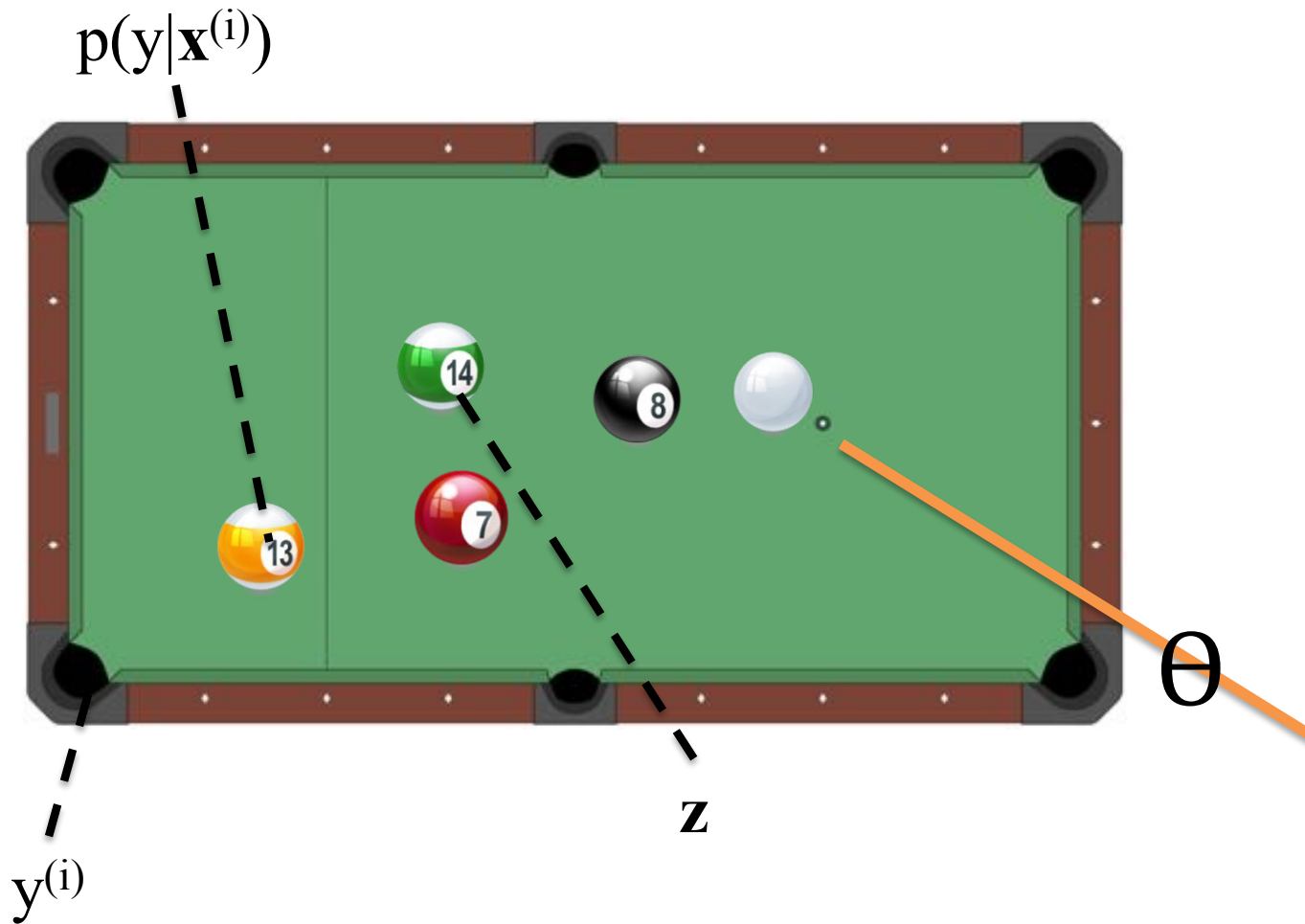$p(y|\mathbf{x}^{(i)})$

$\mathbf{z}$

$\theta$

$y^{(i)}$

Algorithm

# BACKPROPAGATION

# Backpropagation

*Chalkboard*

    – Example: Backpropagation for Chain Rule #1

> **Differentiation Quiz #1:**
>
> Suppose x = 2 and z = 3, what are dy/dx and dy/dz for the function below? **Round your answer to the nearest integer.**
>
> $$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

# Training | Backpropagation

*Chalkboard*

- SGD for Neural Network

- Example: Backpropagation for Neural Network

# Backpropagation

**Automatic Differentiation – Reverse Mode (aka. Backpropagation)**

Forward Computation
1. Write an **algorithm** for evaluating the function y = f(**x**). The algorithm defines a **directed acyclic graph**, where each variable is a node (i.e. the "**computation graph**")
2. Visit each node in **topological order**.
   For variable $u_i$ with inputs $v_1,\ldots, v_N$
   a. Compute $u_i = g_i(v_1,\ldots, v_N)$
   b. Store the result at the node

Backward Computation
1. **Initialize** all partial derivatives $dy/du_j$ to 0 and $dy/dy = 1$.
2. Visit each node in **reverse topological order**.
   For variable $u_i = g_i(v_1,\ldots, v_N)$
   a. We already know $dy/du_i$
   b. Increment $dy/dv_j$ by $(dy/du_i)(du_i/dv_j)$
   (Choice of algorithm ensures computing $(du_i/dv_j)$ is easy)

**Return** partial derivatives $dy/du_i$ for all variables