



10-301/601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

Neural Networks

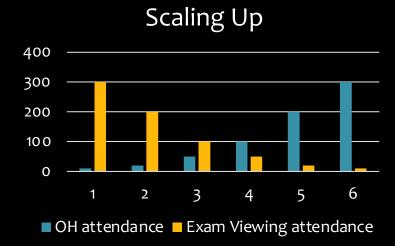


Backpropagation

Matt Gormley Lecture 12 Oct. 4, 2022

Reminders

- Post-Exam Followup:
 - Exam Viewing
 - Exit Poll: Exam 1
 - Grade Summary 1



- Homework 4: Logistic Regression
 - Out: Tue, Oct 4
 - Due: Thu, Oct 13 at 11:59pm
- Homework 5: Neural Networks
 - Out: Thu, Oct 13
 - Due: Thu, Oct 27 at 11:59pm

ARCHITECTURES

Neural Network Architectures

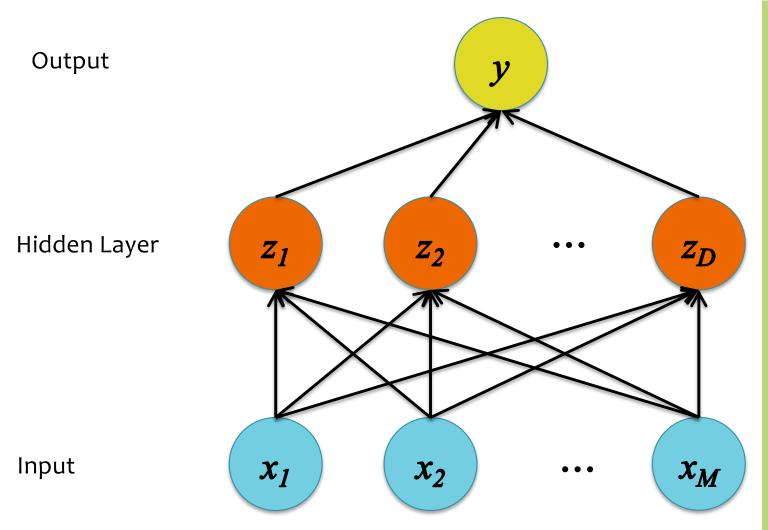
Even for a basic Neural Network, there are many design decisions to make:

- # of hidden layers (depth)
- 2. # of units per hidden layer (width)
- 3. Type of activation function (nonlinearity)
- 4. Form of objective function
- 5. How to initialize the parameters

BUILDING WIDER NETWORKS



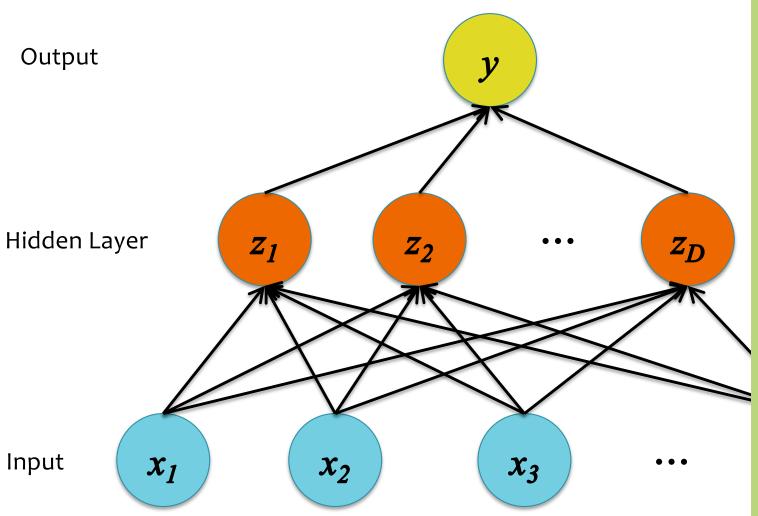
Q: How many hidden units, D, should we use?



- a selection of the most useful features
- nonlinear combinations of the features
- a lower dimensional projection of the features
- a higher dimensional projection of the features
- a copy of the input features
- a mix of the above



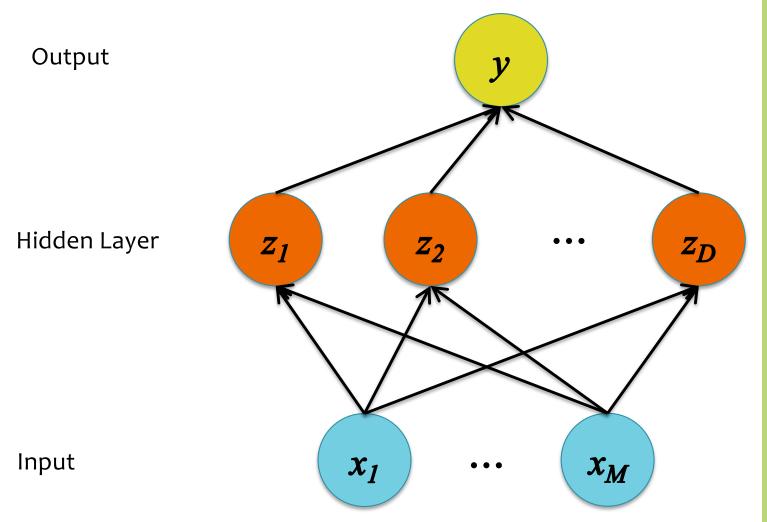
Q: How many hidden units, D, should we use?



- a selection of the most useful features
- nonlinear combinations of the features
- a lower dimensional projection of the features
- a higher dimensional projection of the features
- a copy of the input features
- a mix of the above



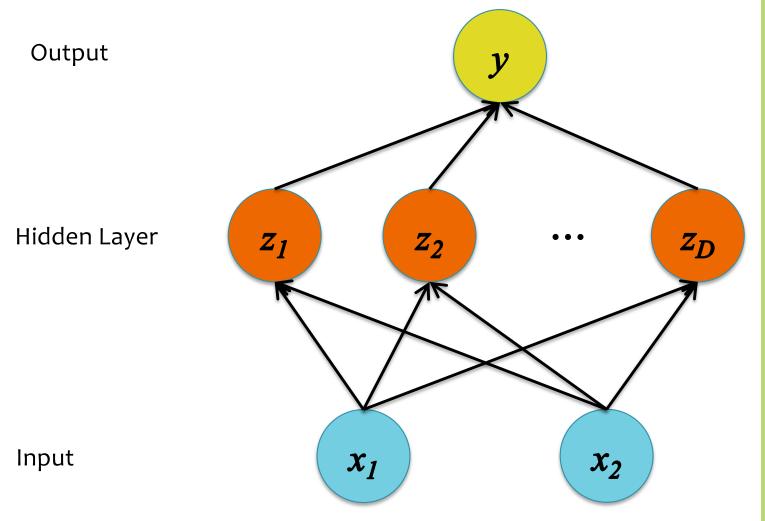
Q: How many hidden units, D, should we use?



- a selection of the most useful features
- nonlinear combinations of the features
- a lower dimensional projection of the features
- a higher dimensional projection of the features
- a copy of the input features
- a mix of the above



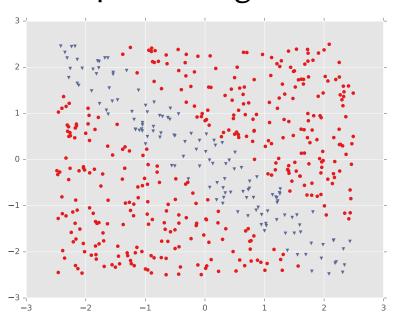
In the following examples, we have two input features, M=2, and we vary the number of hidden units, D.



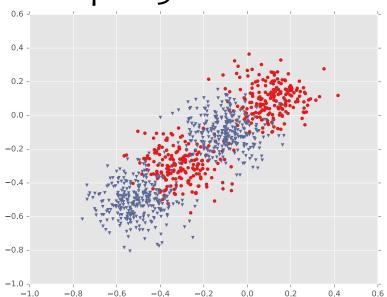
- a selection of the most useful features
- nonlinear combinations of the features
- a lower dimensional projection of the features
- a higher dimensional projection of the features
- a copy of the input features
- a mix of the above

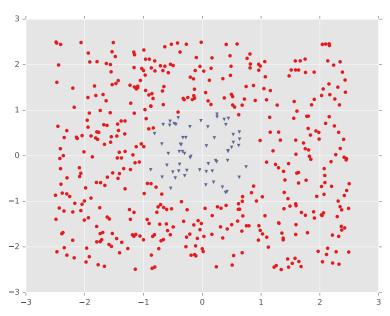
Examples 1 and 2

DECISION BOUNDARY EXAMPLES

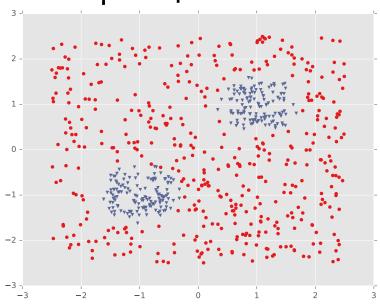


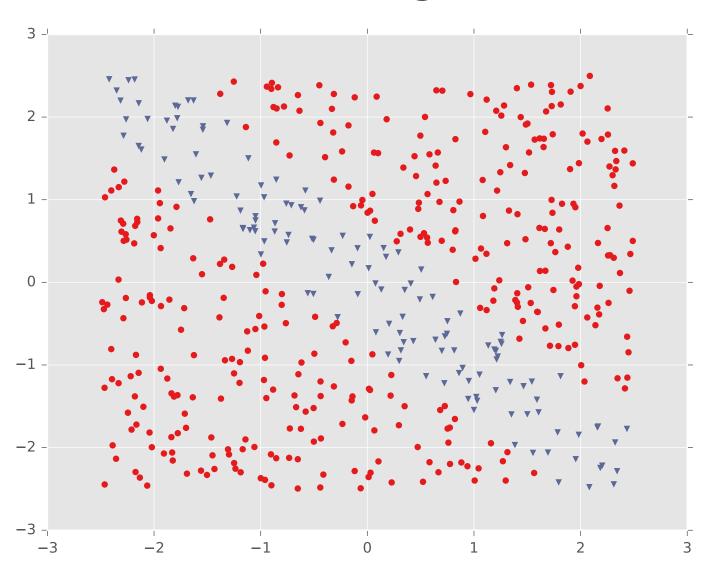
Example #3: Four Gaussians





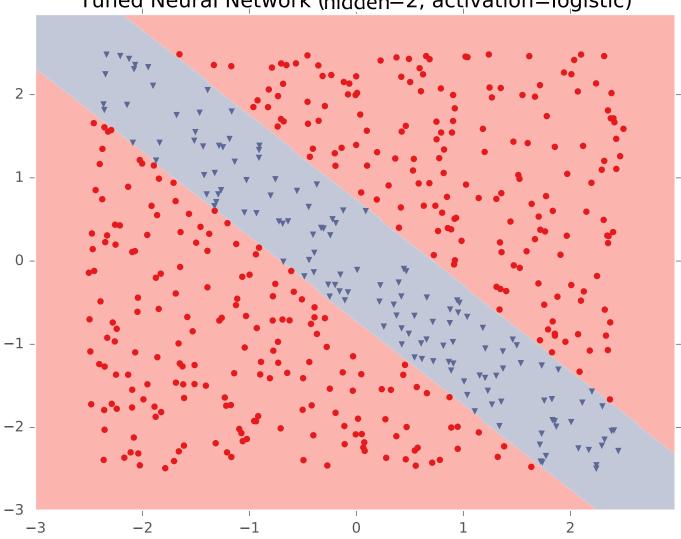
Example #4: Two Pockets

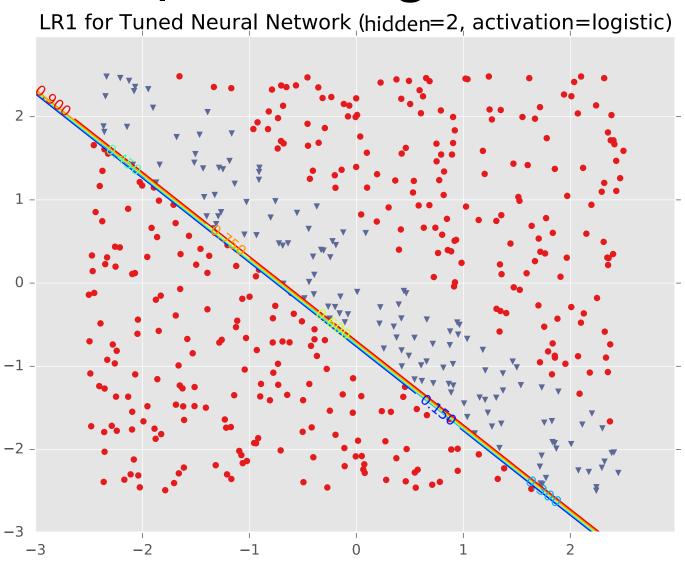


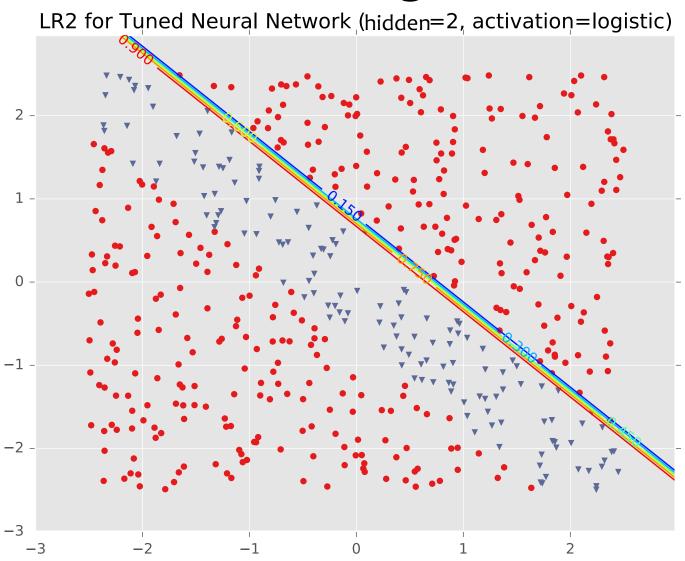


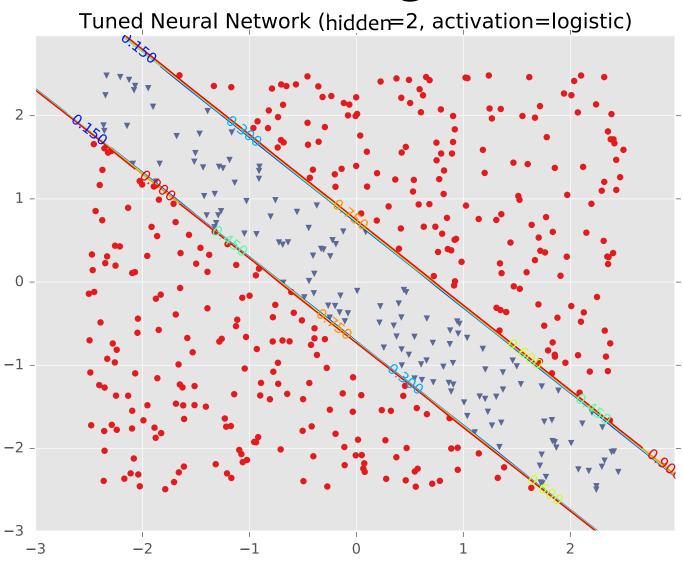


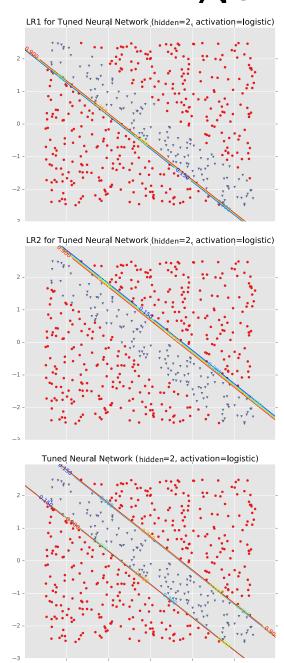
Tuned Neural Network (hidden=2, activation=logistic)



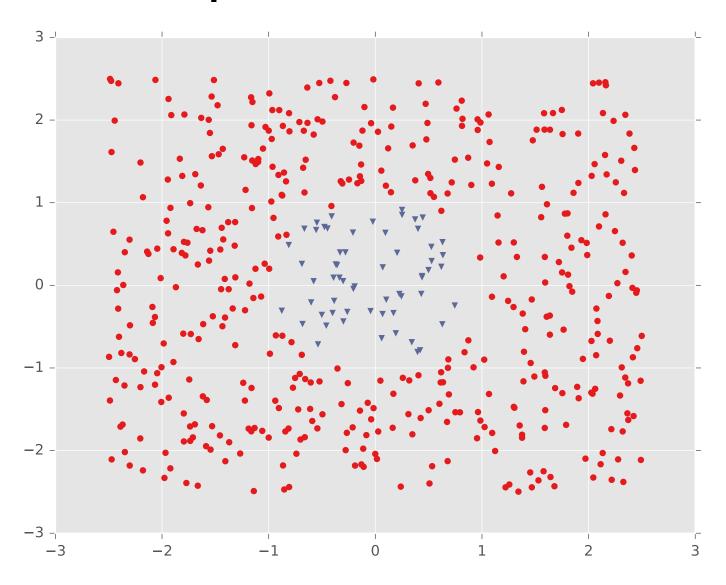


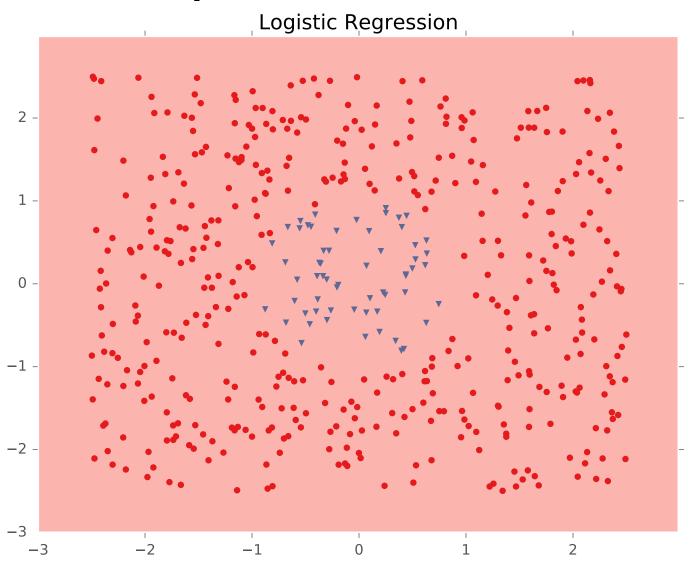




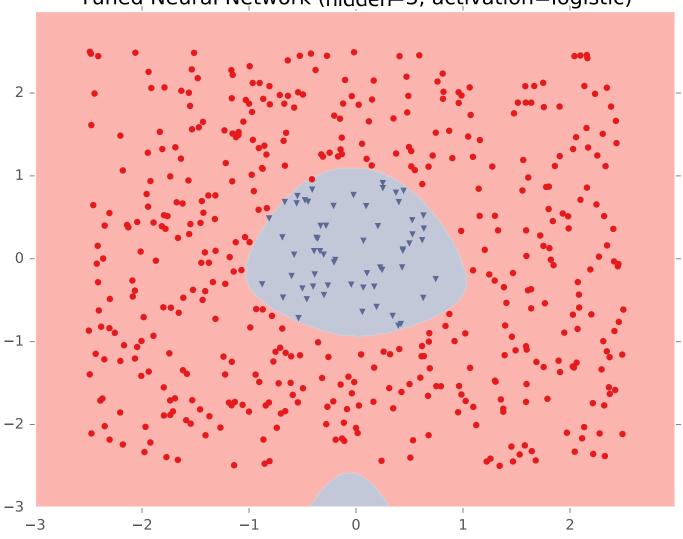


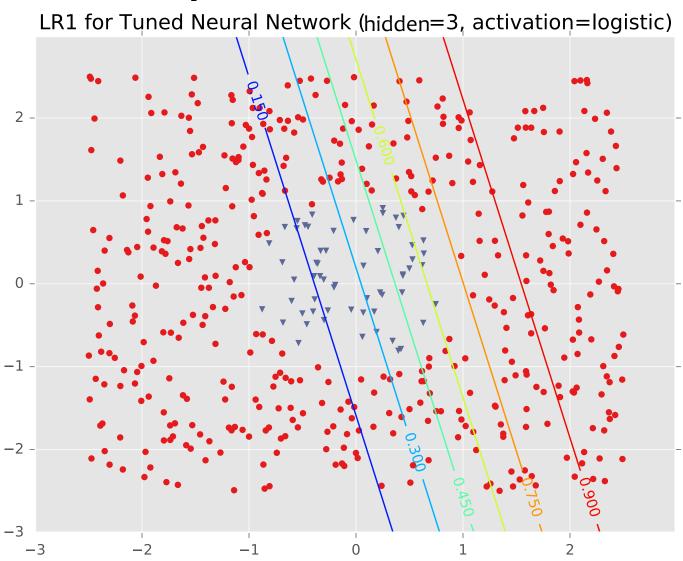
Tuned Neural Network (hidden=2, activation=logistic)

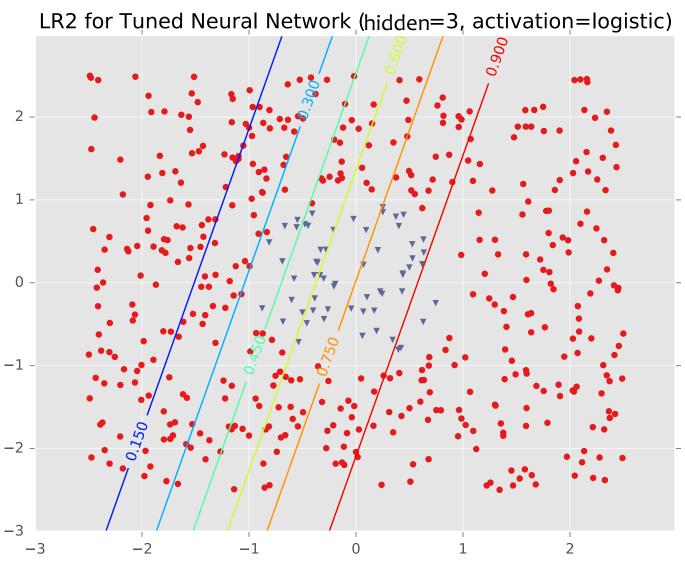


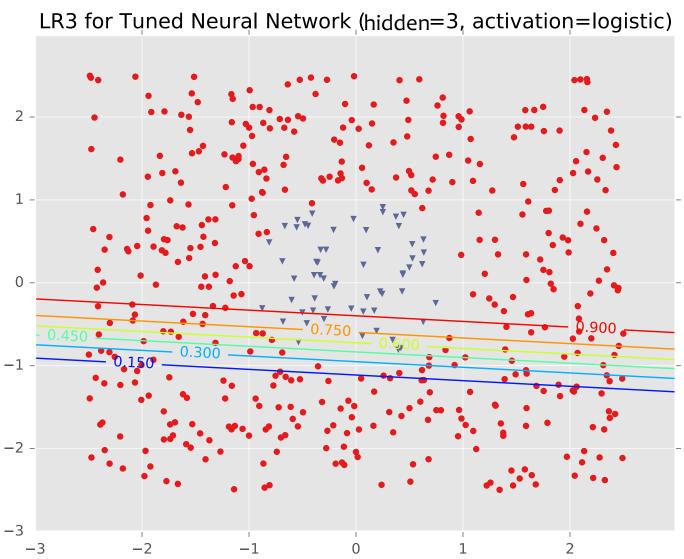


Tuned Neural Network (hidden=3, activation=logistic)

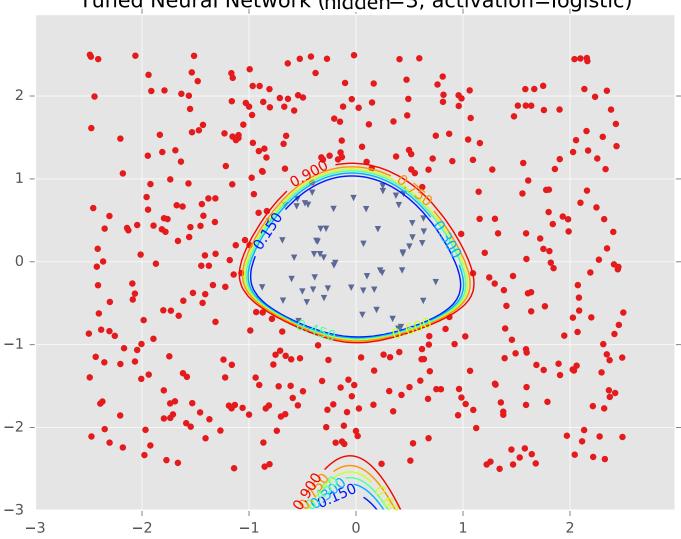


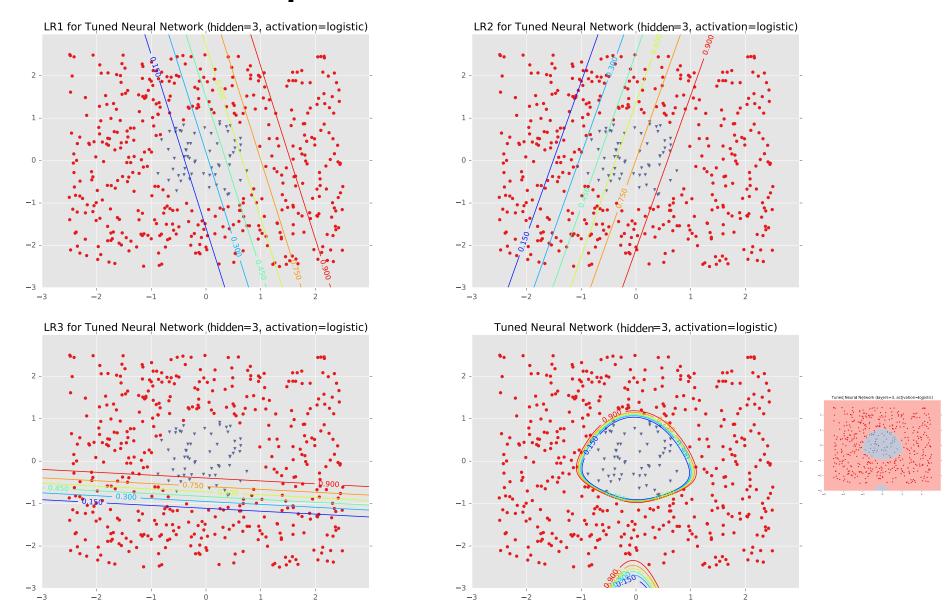






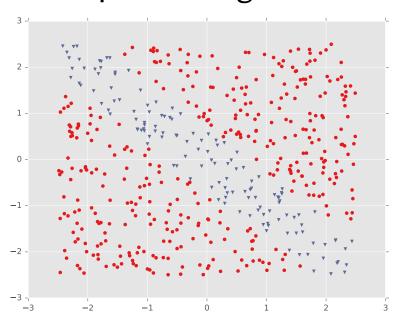
Tuned Neural Network (hidden=3, activation=logistic)



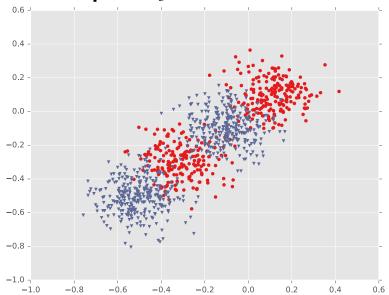


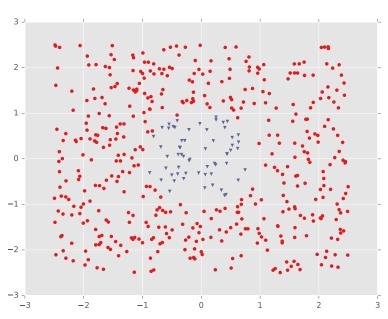
Examples 3 and 4

DECISION BOUNDARY EXAMPLES

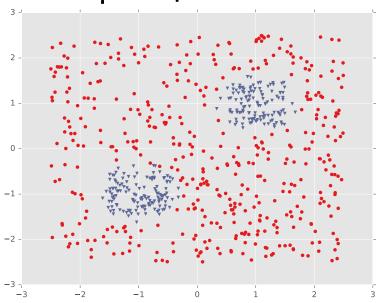


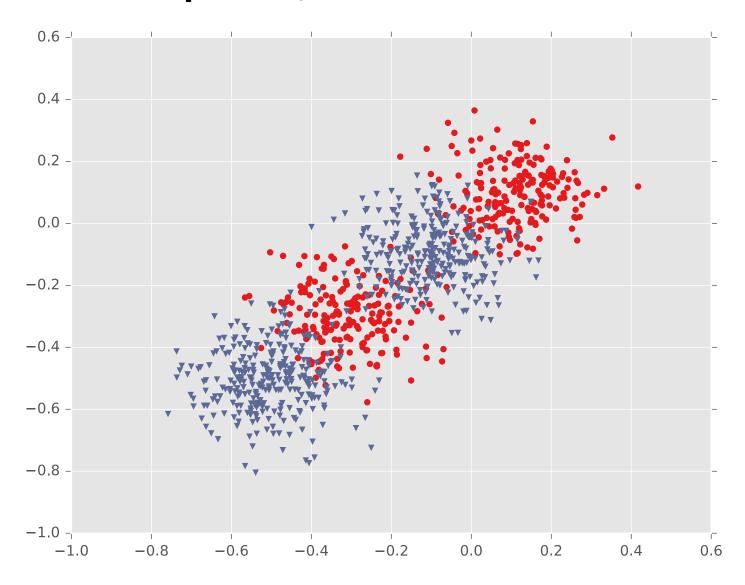
Example #3: Four Gaussians

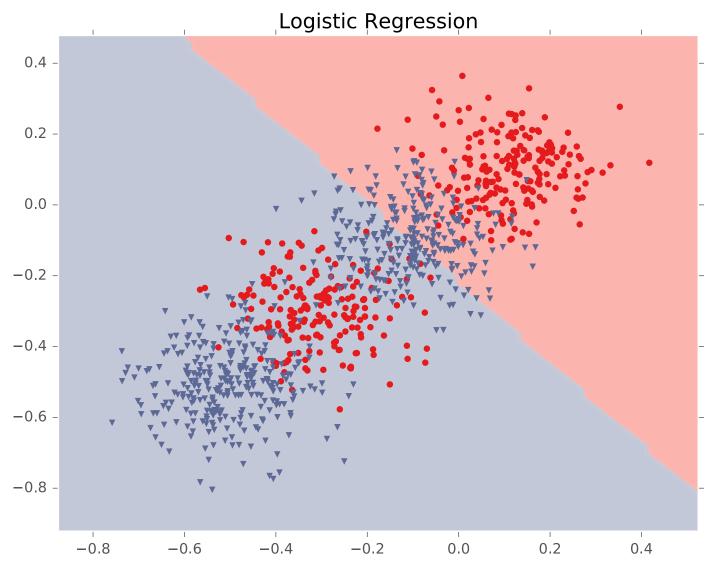


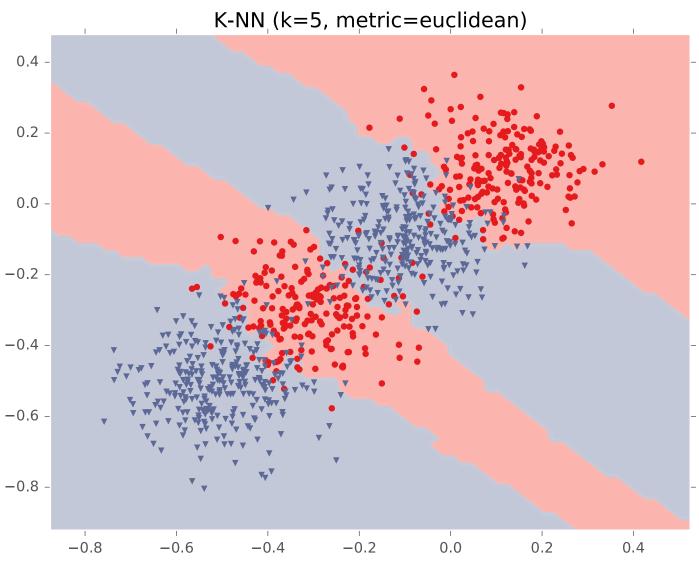


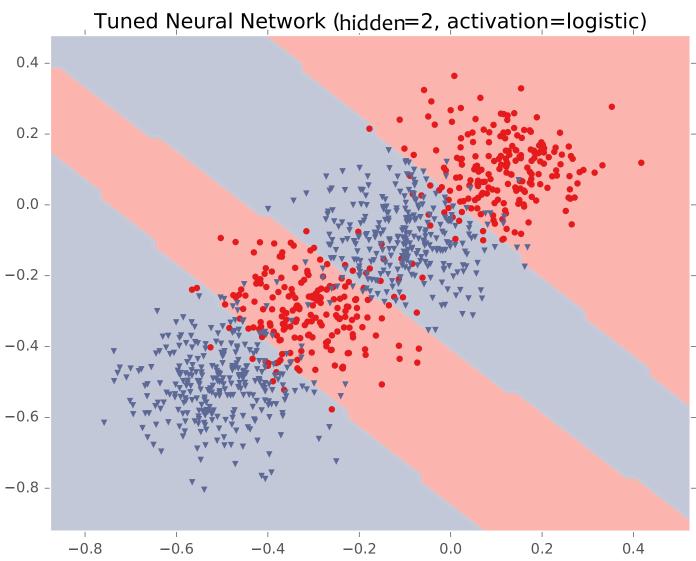
Example #4: Two Pockets

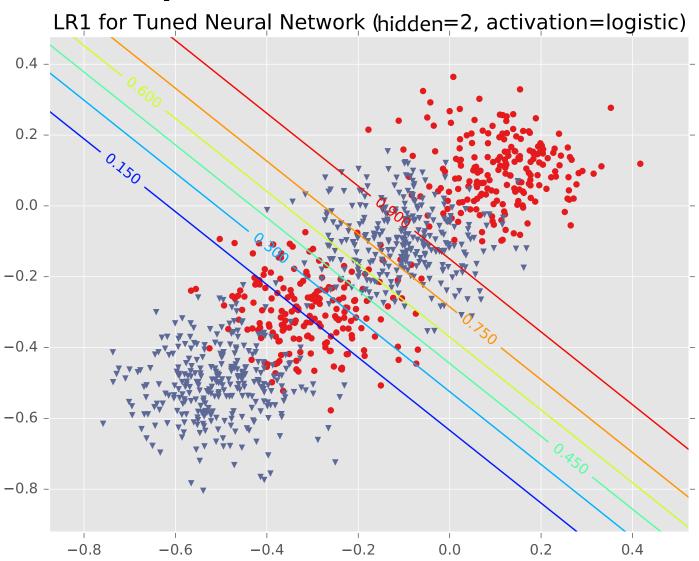


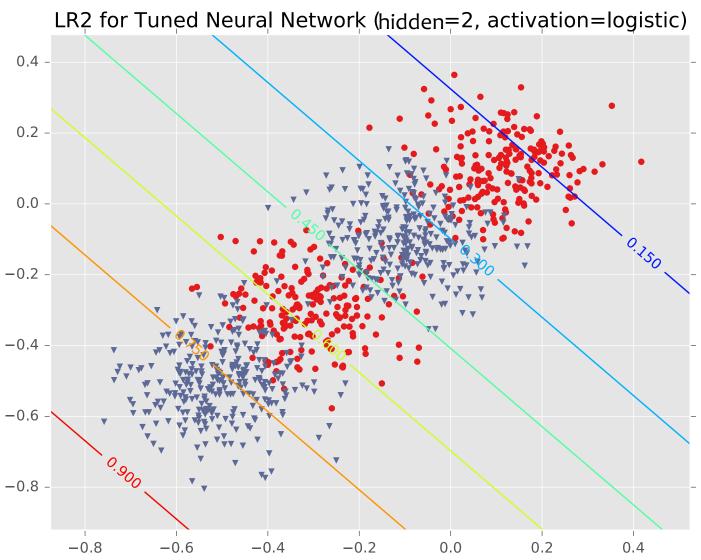


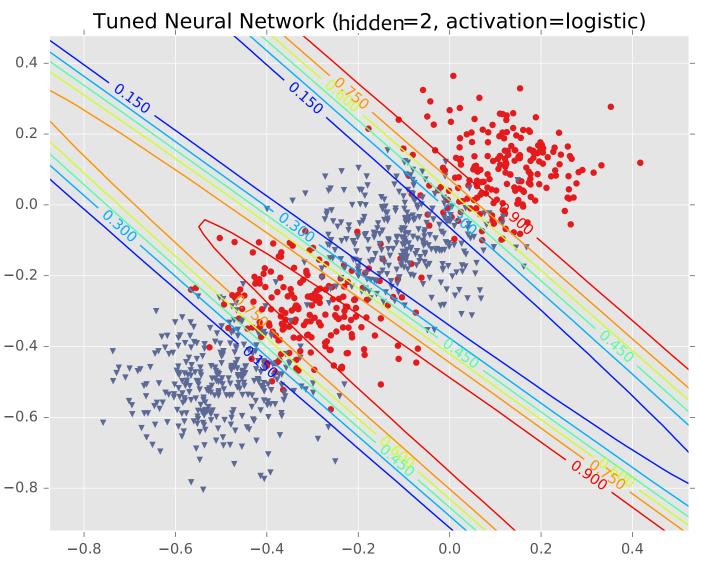




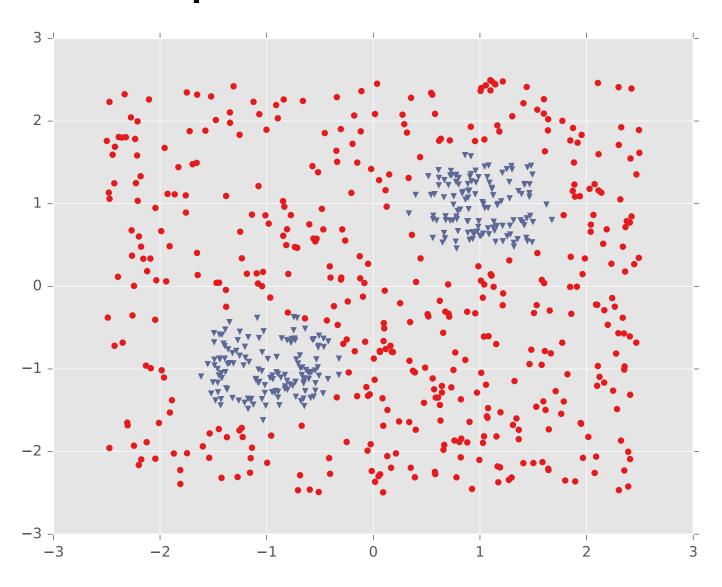


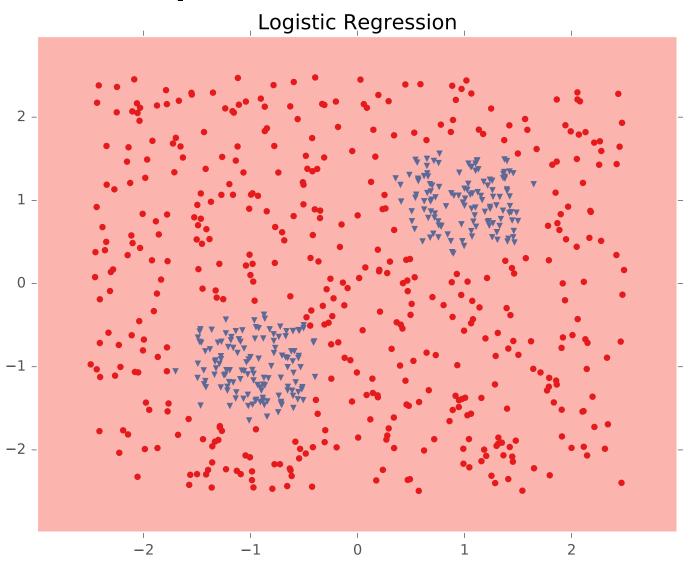


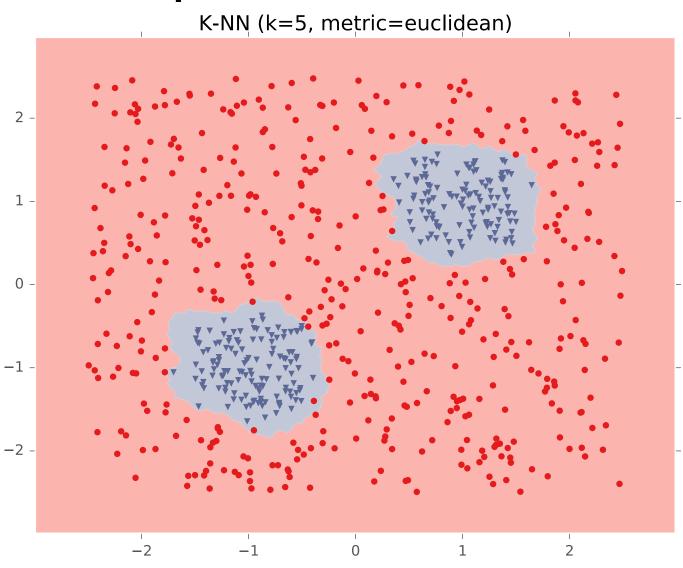




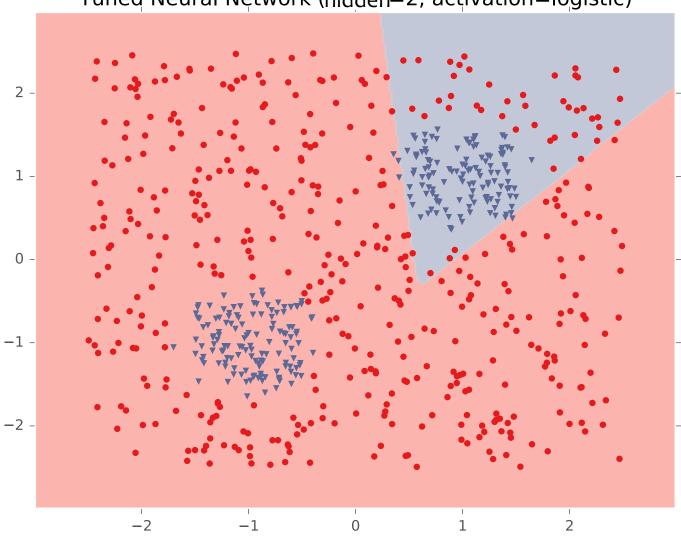
Example #4: Two Pockets



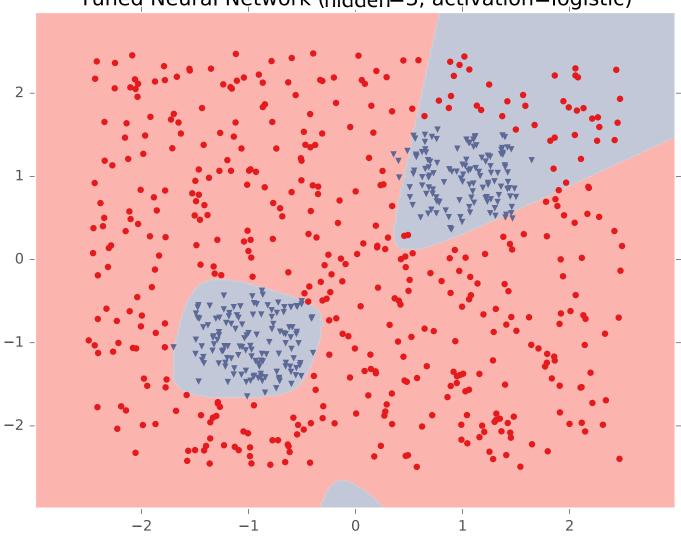




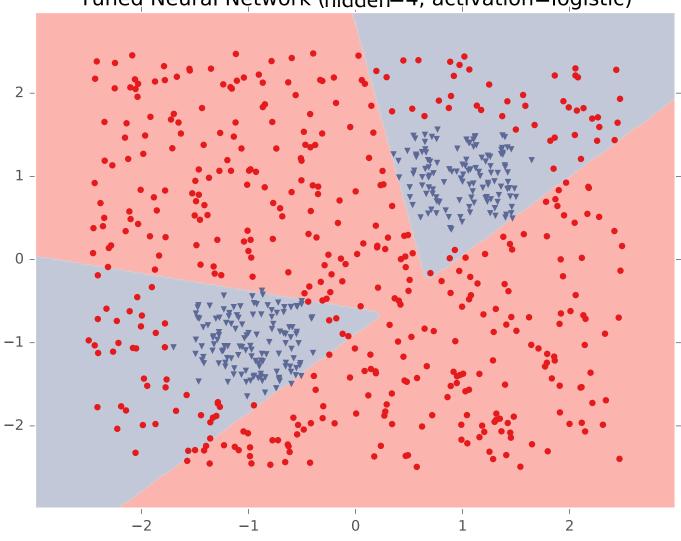
Tuned Neural Network (hidden=2, activation=logistic)



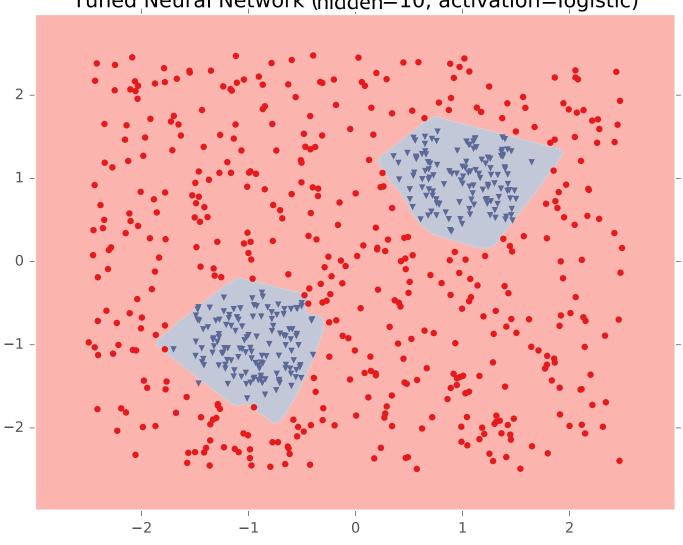
Tuned Neural Network (hidden=3, activation=logistic)



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



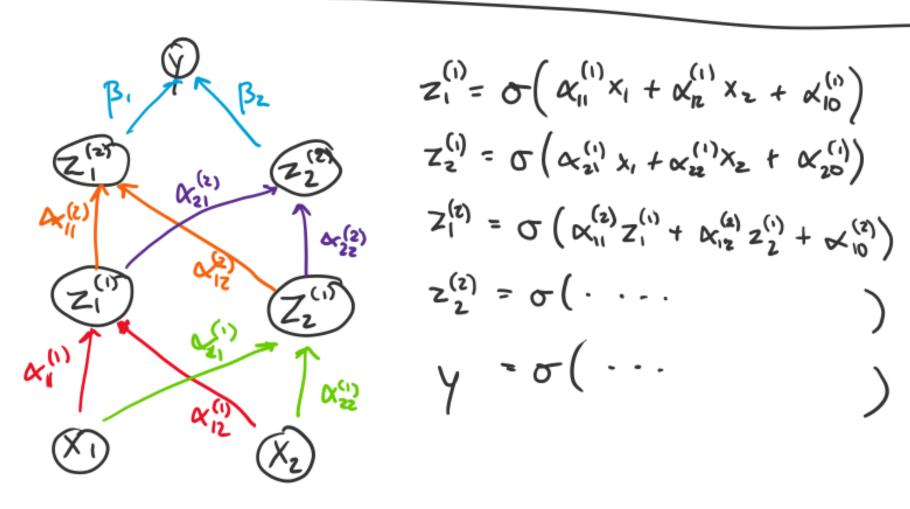
Tuned Neural Network (hidden=10, activation=logistic)



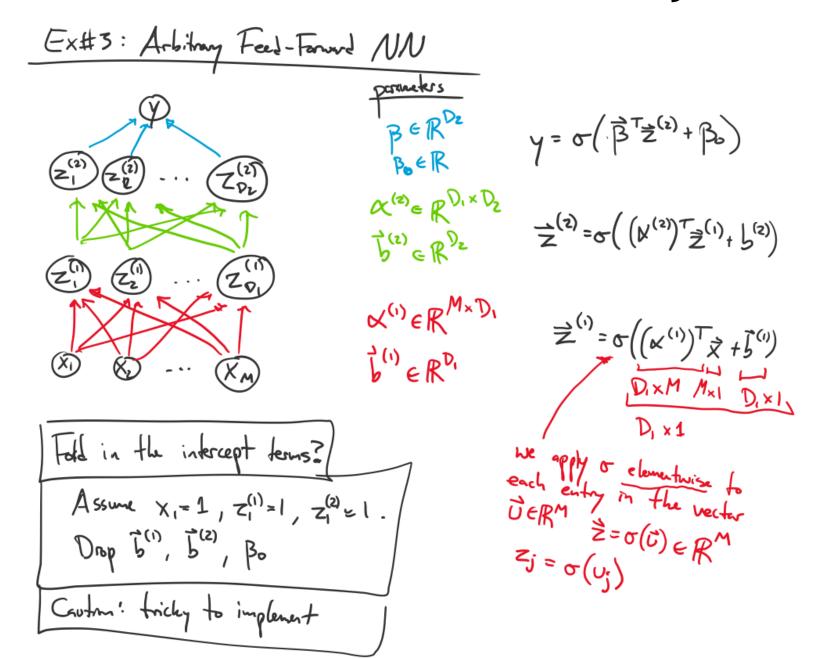
BUILDING DEEPER NETWORKS

Neural Net w/2 Hidden Layers

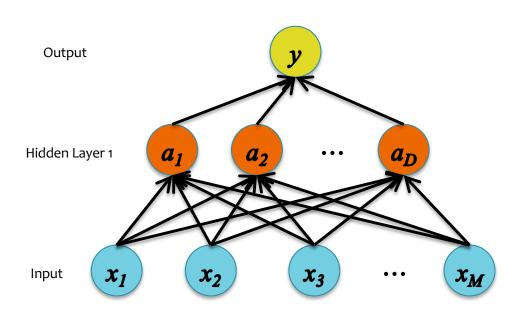
Ex#2: NN w/2 Hidden Layers and 2 units each



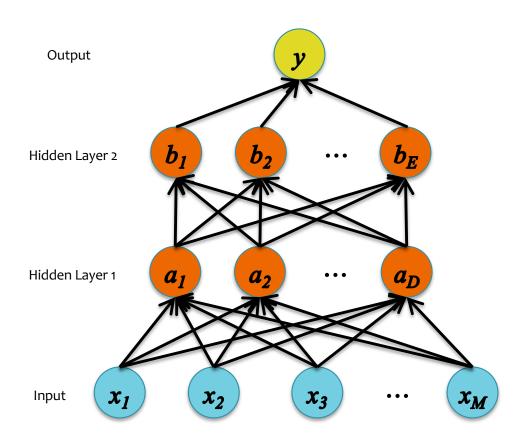
Neural Net w/2 Hidden Layers



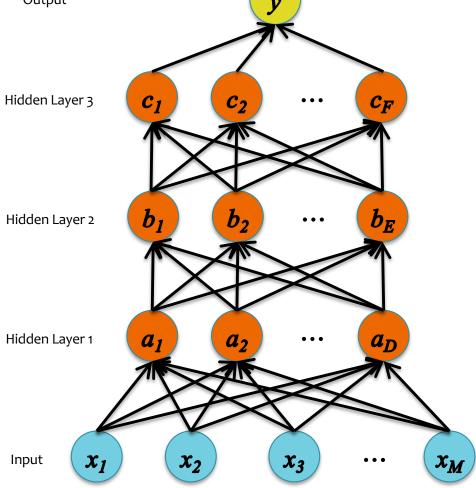
Q: How many layers should we use?



Q: How many layers should we use?



Q: How many layers should we use?



Q: How many layers should we use?

Theoretical answer:

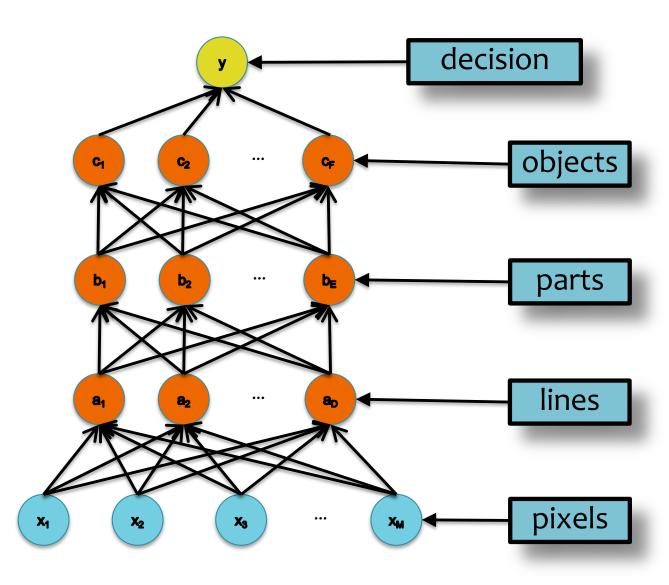
- A neural network with 1 hidden layer is a universal function approximator
- Cybenko (1989): For any continuous function g(x), there exists a 1-hidden-layer neural net $h_{\theta}(x)$ s.t. $|h_{\theta}(x) g(x)| < \epsilon$ for all x, assuming sigmoid activation functions

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"

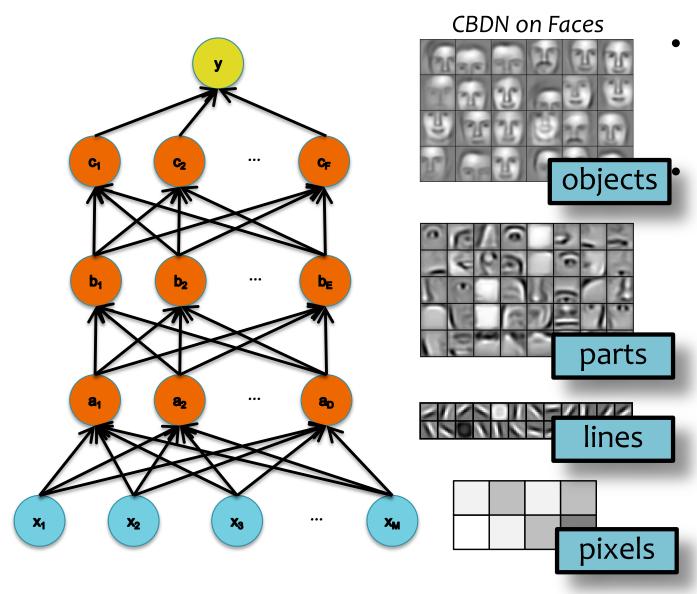
Big caveat: You need to know and use the right tricks.

Feature Learning



- Traditional feature
 engineering: build up
 levels of abstraction
 by hand
- Deep networks (e.g. convolution networks): learn the increasingly higher levels of abstraction from data
 - each layer is a learned feature representation
 - sophistication increases in higher layers

Feature Learning

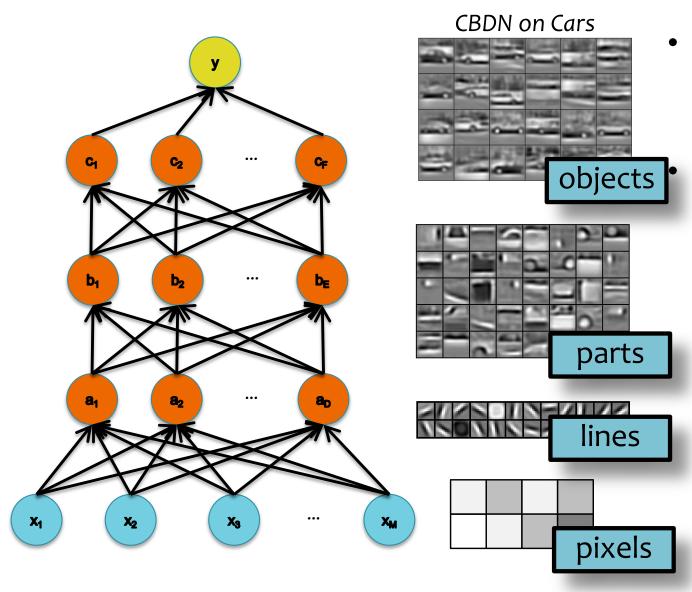


Traditional feature engineering: build up levels of abstraction by hand

Deep networks (e.g. convolution networks): learn the increasingly higher levels of abstraction from data

- each layer is a learned feature representation
- sophistication increases in higher layers

Feature Learning



Traditional feature engineering: build up levels of abstraction by hand

Deep networks (e.g. convolution networks): learn the increasingly higher levels of abstraction from data

- each layer is a learned feature representation
- sophistication increases in higher layers

ARCHITECTURES

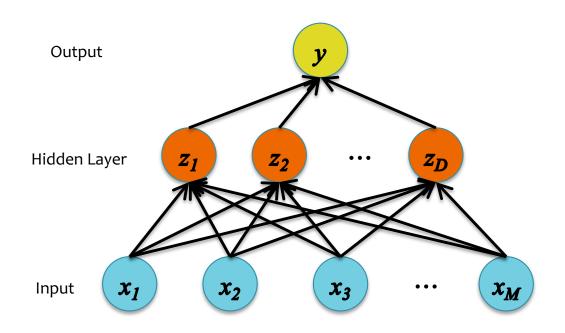
Neural Network Architectures

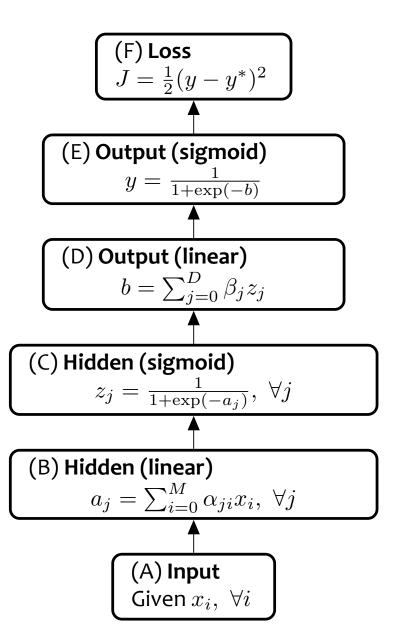
Even for a basic Neural Network, there are many design decisions to make:

- # of hidden layers (depth)
- 2. # of units per hidden layer (width)
- 3. Type of activation function (nonlinearity)
- 4. Form of objective function
- 5. How to initialize the parameters

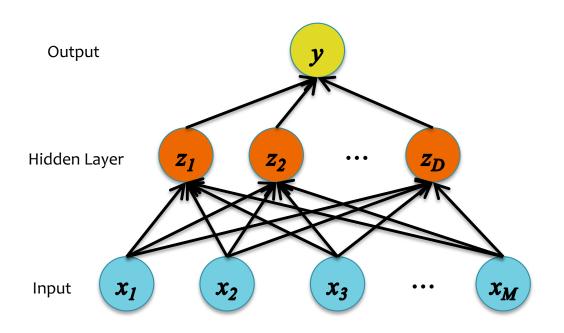
ACTIVATION FUNCTIONS

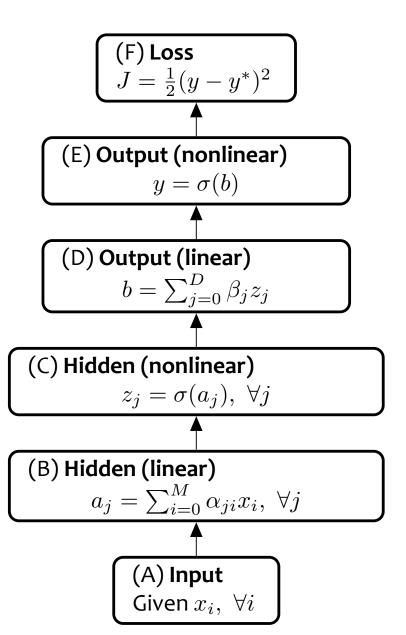
Neural Network with sigmoid activation functions





Neural Network with arbitrary nonlinear activation functions

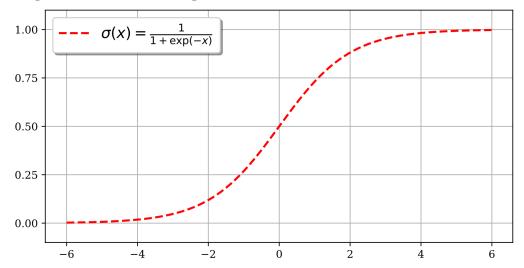




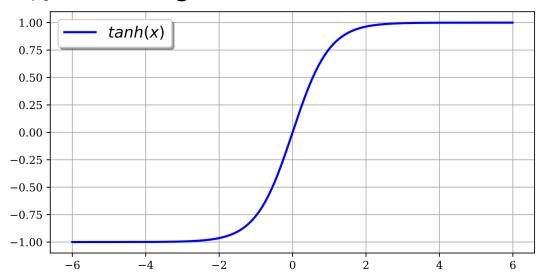
So far, we've assumed that the activation function (nonlinearity) is always the sigmoid function...

... but the sigmoid is not widely used in modern neural networks

Sigmoid (aka. logistic) function

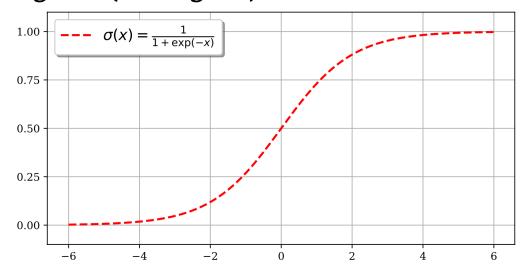


Hyperbolic tangent function

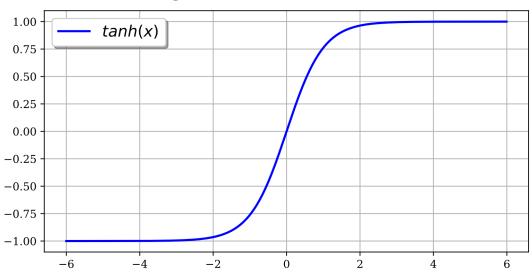


- sigmoid, $\sigma(x)$
 - output in range(0,1)
 - good for probabilistic outputs
- hyperbolic tangent, tanh(x)
 - similar shape to sigmoid, but output in range (-1,+1)

Sigmoid (aka. logistic) function



Hyperbolic tangent function



Understanding the difficulty of training deep feedforward neural networks

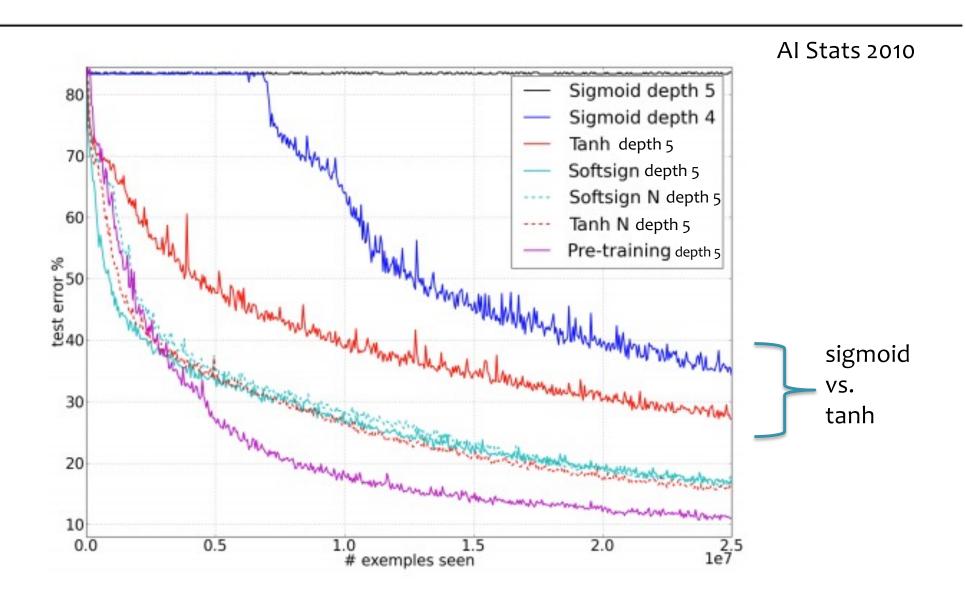
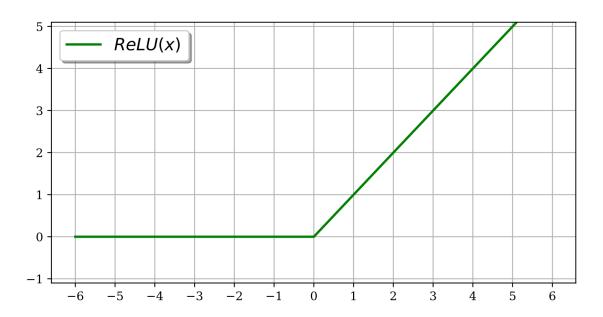


Figure from Glorot & Bentio (2010)

- Rectified Linear Unit (ReLU)
 - avoids the vanishing gradient problem
 - derivative is fast to compute

$$ReLU(x) = max(0, x)$$



- Rectified Linear Unit (ReLU)
 - avoids the vanishing gradient problem
 - derivative is fast to compute

$$ReLU(x) = max(0, x)$$

- Exponential Linear Unit (ELU)
 - same as ReLU on positive inputs
 - unlike ReLU, allows negative outputs and smoothly transitions for x < 0

$$\mathsf{ELU}(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha(\exp(x) - 1), & \text{if } x \le 0 \end{cases}$$

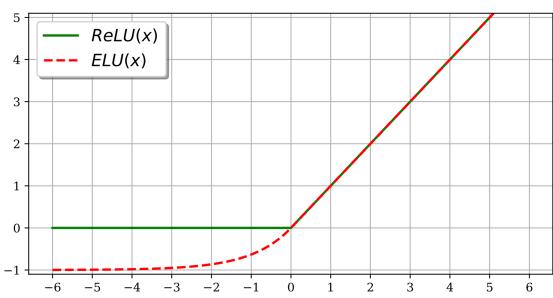
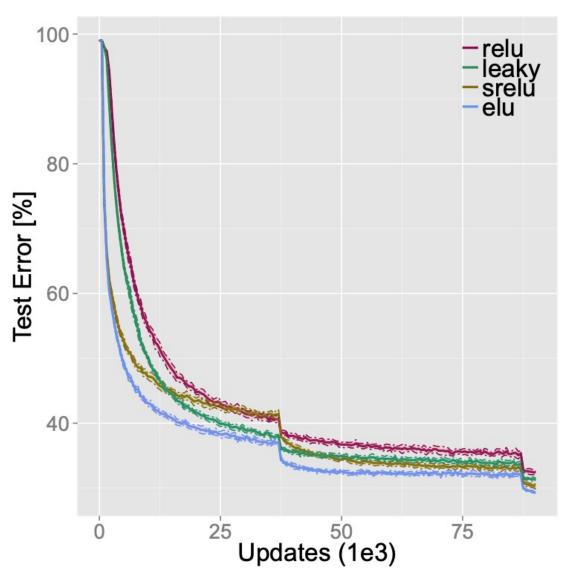


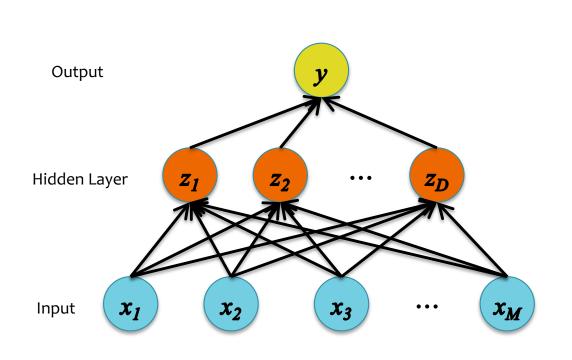
Image Classification Benchmark (CIFAR-10)

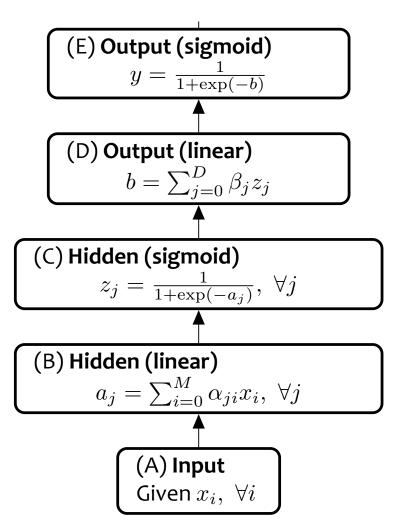


- Training loss converges fastest with ELU
- 2. ELU(x) yields lower test error than ReLU(x) on CIFAR-10

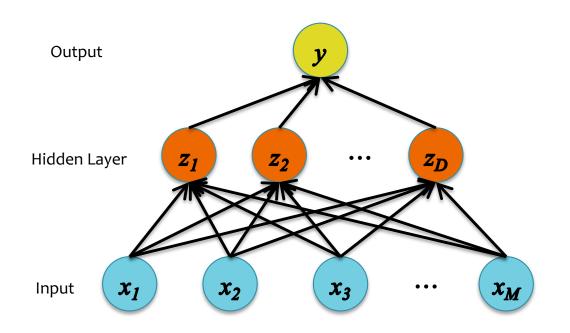
LOSS FUNCTIONS & OUTPUT LAYERS

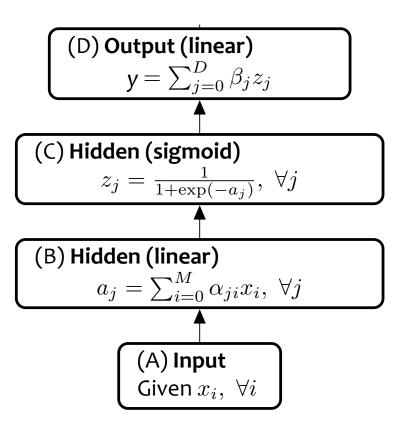
Neural Network for Classification





Neural Network for Regression





Objective Functions for NNs

1. Quadratic Loss:

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

$$J = \ell_Q(y, y^{(i)}) = \frac{1}{2}(y - y^{(i)})^2$$
$$\frac{dJ}{dy} = y - y^{(i)}$$

2. Binary Cross-Entropy:

- the same objective as Binary Logistic Regression
- i.e. negative log likelihood
- This requires our output y to be a probability in [0,1]

$$J = \ell_{CE}(y, y^{(i)}) = -(y^{(i)} \log(y) + (1 - y^{(i)}) \log(1 - y))$$

$$\frac{dJ}{dy} = -\left(y^{(i)} \frac{1}{y} + (1 - y^{(i)}) \frac{1}{y - 1}\right)$$

Objective Functions for NNs

Cross-entropy vs. Quadratic loss

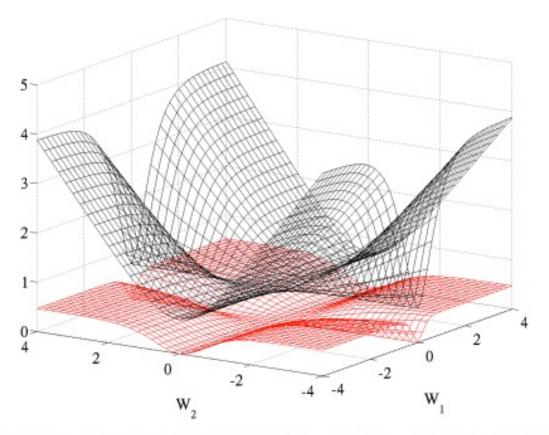
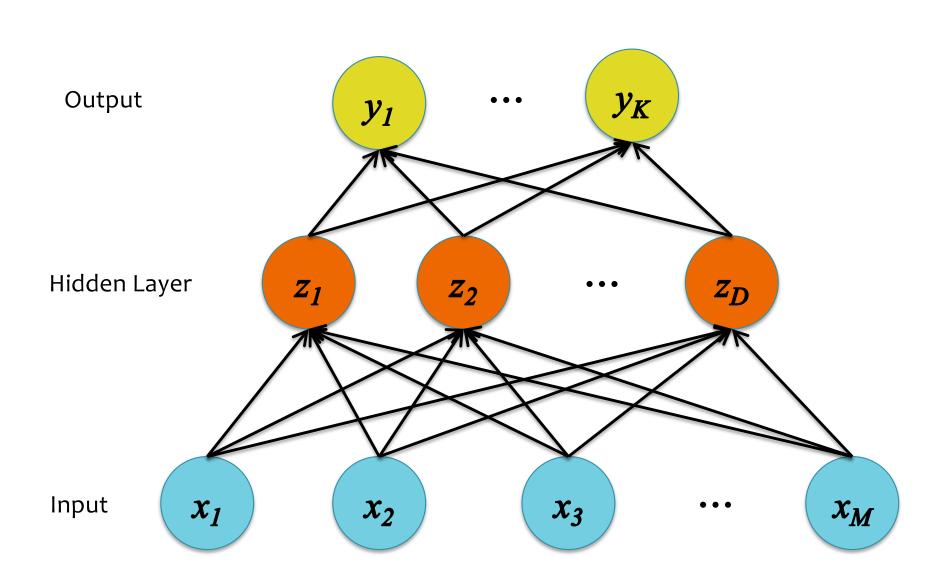


Figure 5: Cross entropy (black, surface on top) and quadratic (red, bottom surface) cost as a function of two weights (one at each layer) of a network with two layers, W_1 respectively on the first layer and W_2 on the second, output layer.

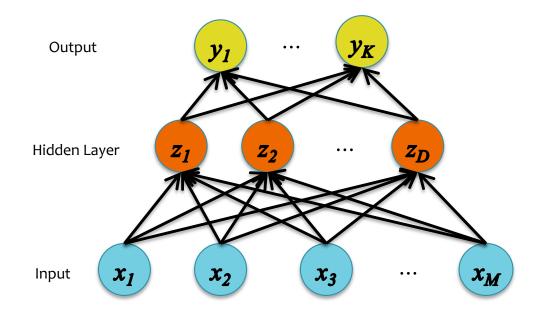
Multiclass Output

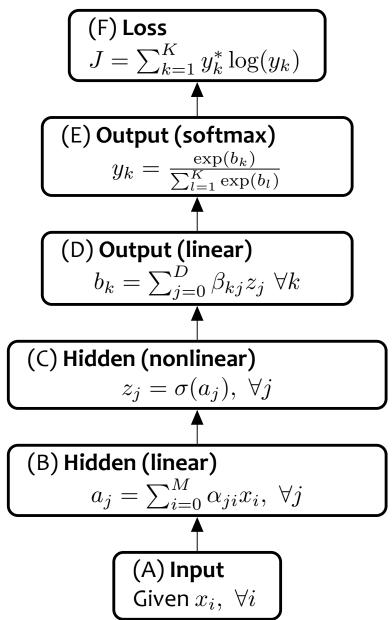


Multiclass Output

Softmax:

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





Objective Functions for NNs

- 3. Cross-Entropy for Multiclass Outputs:
 - i.e. negative log likelihood for multiclass outputs
 - Suppose output is a random variable Y that takes one of K values
 - Let y⁽ⁱ⁾ represent our true label as a one-hot vector:

$$\mathbf{y}^{(i)} = \begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 & \dots & 0 \\ & 1 & 2 & 3 & 4 & 5 & 6 & \dots & K \end{bmatrix}$$

Assume our model outputs a length K vector of probabilities:

$$y = softmax(f_{scores}(x, \theta))$$

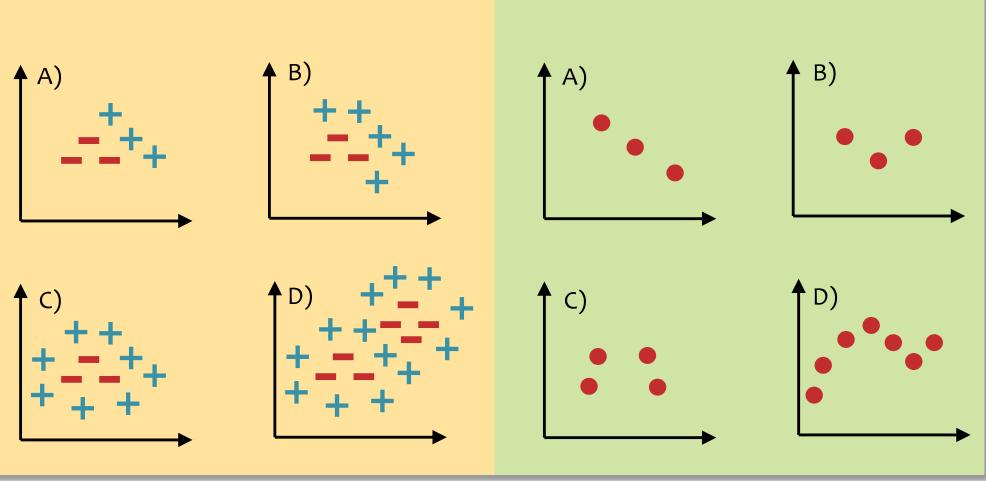
- Then we can write the log-likelihood of a single training example $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ as:

$$J = \ell_{CE}(\mathbf{y}, \mathbf{y}^{(i)}) = -\sum_{k=1}^{K} y_k^{(i)} \log(y_k)$$

Neural Network Errors

Question X: For which of the datasets below does there exist a one-hidden layer neural network that achieves zero *classification* error? **Select all that apply.**

Question Y: For which of the datasets below does there exist a one-hidden layer neural network for *regression* that achieves *nearly* zero MSE? **Select all that apply.**



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

APPROACHES TO DIFFERENTIATION

Background

A Recipe for Machine Learning

1. Given training data:

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

3. Define goal:

$$oldsymbol{ heta}^* = rg\min_{oldsymbol{ heta}} \sum_{i=1}^N \ell(f_{oldsymbol{ heta}}(oldsymbol{x}_i), oldsymbol{y}_i)$$

- 2. Choose each of these:
 - Decision function

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

Loss function

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

4. Train with SGD:

(take small steps opposite the gradient)

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

Background

A Recipe for Gradients

1. Given training dat

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

- 2. Choose each of the
 - Decision function

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

Loss function

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

Backpropagation can compute this gradient!

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

opposite the gradient)
$$oldsymbol{ heta}^{(t)} - \eta_t
abla \ell(f_{oldsymbol{ heta}}(oldsymbol{x}_i), oldsymbol{y}_i)$$

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(x) on any input 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(x)

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

Compute $\dfrac{\partial f(\mathbf{x})_i}{\partial x_j} orall i, j$

Approaches to Differentiation

- 3. Automatic Differentiation Reverse Mode
 - Note: Called Backpropagation when applied to Neural Nets
 - Pro: Computes partial derivatives of one output f(x)_i with respect to all inputs x_i in time proportional to computation of f(x)
 - Con: Slow for high dimensional outputs (e.g. vector-valued functions)
 - Required: Algorithm for computing f(x)
- 4. Automatic Differentiation Forward Mode
 - Note: Easy to implement. Uses dual numbers.
 - Pro: Computes partial derivatives of all outputs f(x)_i with respect to one input x_i in time proportional to computation of f(x)
 - Con: Slow for high dimensional inputs (e.g. vector-valued x)
 - Required: Algorithm for computing f(x)

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

Compute $\frac{\partial f(\mathbf{x})_i}{\partial x_j} orall i, j$

THE FINITE DIFFERENCE METHOD

Finite Difference Method

The centered finite difference approximation is:

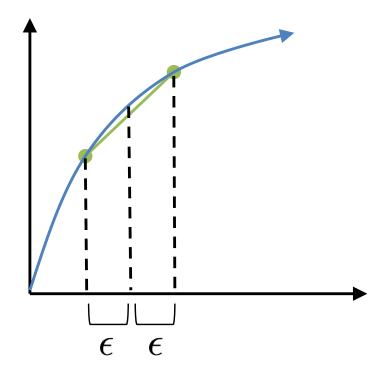
$$\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) \approx \frac{(J(\boldsymbol{\theta} + \epsilon \cdot \boldsymbol{d}_i) - J(\boldsymbol{\theta} - \epsilon \cdot \boldsymbol{d}_i))}{2\epsilon} \tag{1}$$

where d_i is a 1-hot vector consisting of all zeros except for the ith

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



Differentiation Quiz

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]

Differentiation Quiz

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:

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

Let σ 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.

