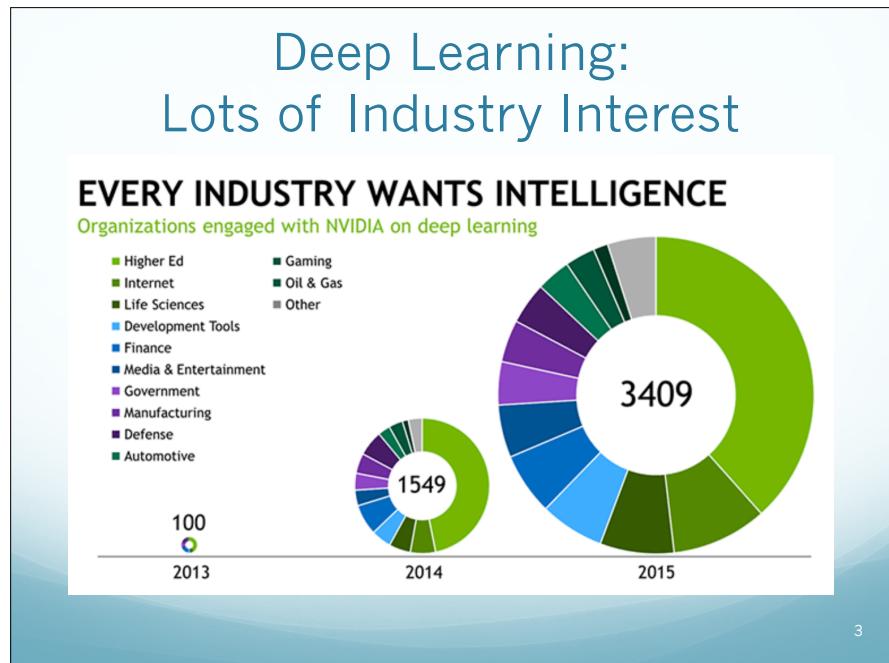
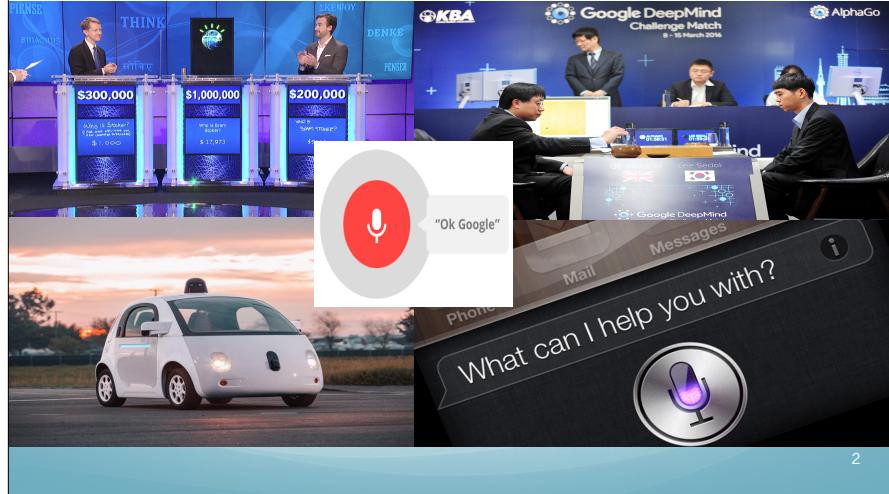


Artificial Neural Networks

Mark Dredze
Machine Learning
CS 601.475

1

Deep Learning: Lots of Results



Deep Learning: Also Some Backlash

The Case Against Deep-Learning Hype

Is there more to AI than neural networks? Gary Marcus, professor of psychology at NYU and ex-director of Uber's AI lab, thinks so. He's [published a critique of deep-learning systems](#) that use neural nets, and it skewers some of the current AI hype.

Ferenc Huszár @fhuszar Following Intelligent Machines

Is AI Riding a One-Trick Pony?

Just about every AI advance you've heard of depends on a breakthrough that's three decades old. Keeping up the pace of progress will require confronting AI's serious limitations.

by James Somers September 29, 2017

4

Outline

- Lecture 1: Neural Networks
 - Nonlinearities
 - Objective functions
 - Training
 - Gradient Computations
- Lecture 2: Deep Learning 1
 - Deep networks
 - Backpropagation
 - Training options
- Lecture 3: Deep Learning 2
 - Activation functions
 - Regularization
 - Dropout
 - Architecture examples

5

Handling Non-Linear Data

- Option 1: Add features by hand that make the data separable
 - Requires feature engineering
- Option 2: Learn a small number of additional features that will suffice
 - Today
- Option 3: Kernel trick

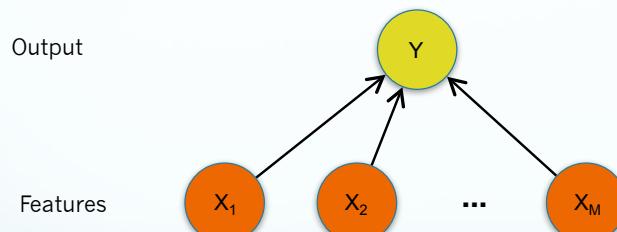
6

Motivation

- Where do features come from?
 - We build them by hand
- What if we wanted to *learn* features?
 - Goal: learn features that give linearly separable data
- After learning features apply usual linear classifier

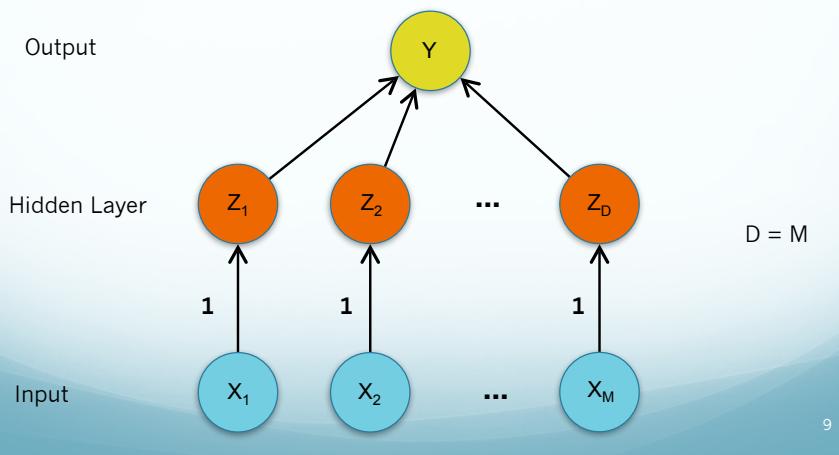
7

Perceptron: Graphical Representation

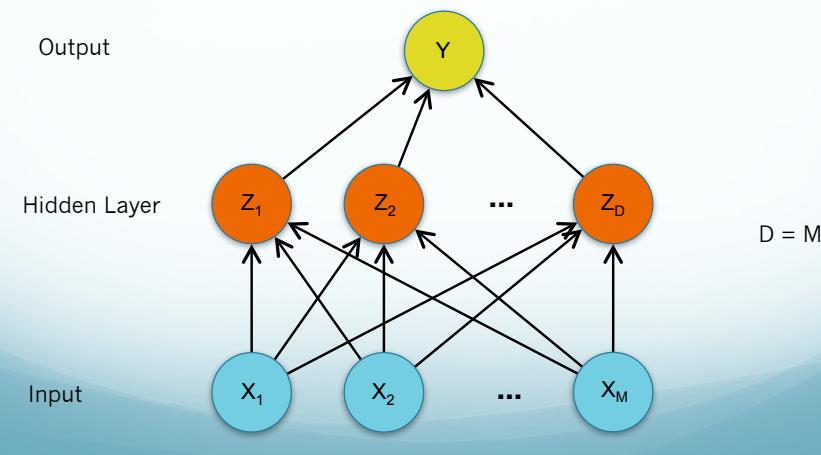


8

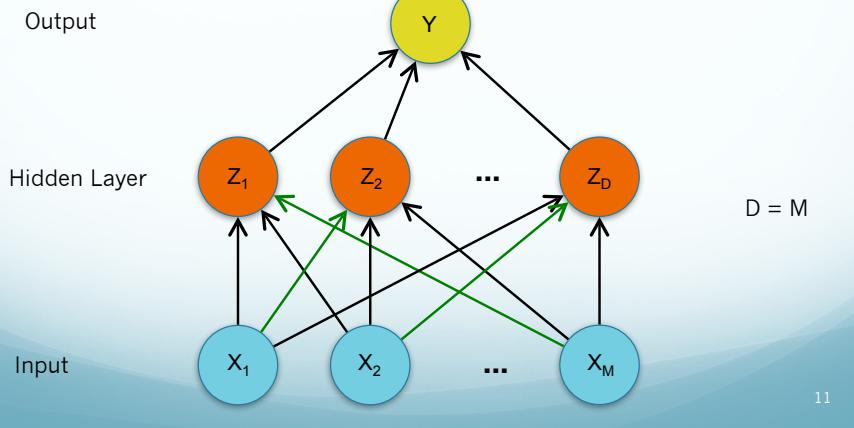
Perceptron: Graphical Representation



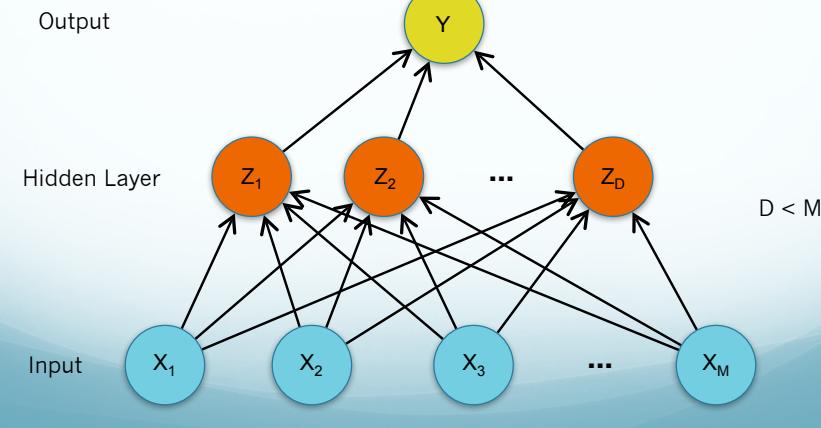
Perceptron: Graphical Representation



Perceptron: Graphical Representation



Perceptron: Graphical Representation



Why?

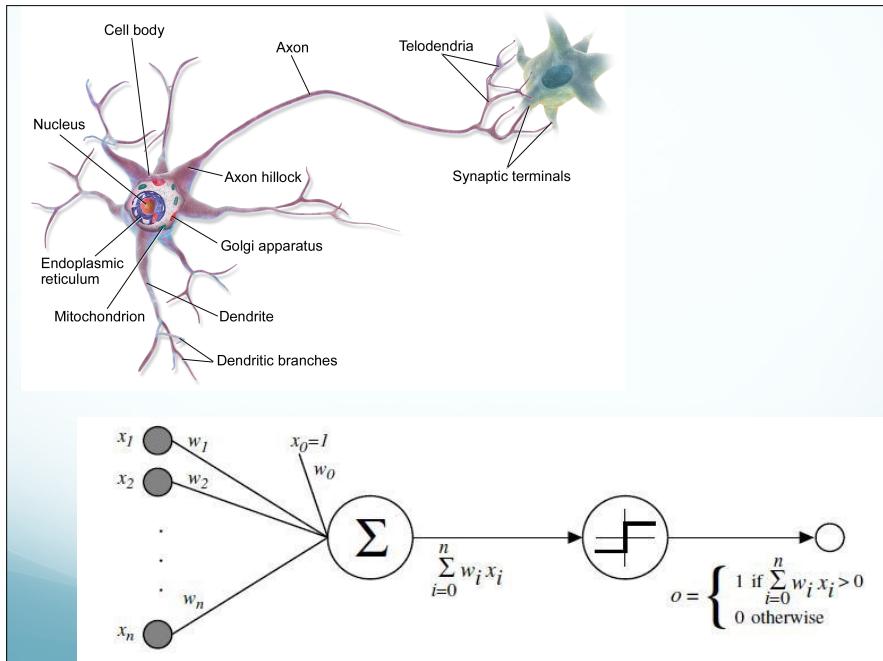
- Constraints from X
 - When $D = M$, likely to copy the features from X to Z
 - When $D < M$, cannot make an exact copy of X
 - Must come up with a representation that is more efficient
- Constraints from Y
 - Z should be a representation that helps learn Y
 - Forces the low-dimensional representation to capture properties of X useful in predicting Y

13

Why Non-Linear

- Generalized linear classifiers!
 - Start with linear function
 - $w \cdot x$
 - Pass the output through a non-linear function
 - $\hat{y} = h(w \cdot x)$
 - What is h?
 - Non-linear function
 - Logistic
 - Tanh
 - Sign
 - ReLu
- Each Z is the output of a non-linear function
 - Combinations of Z are now non-linear in X

14



Multi-Layer Perceptrons Fitting a function to data

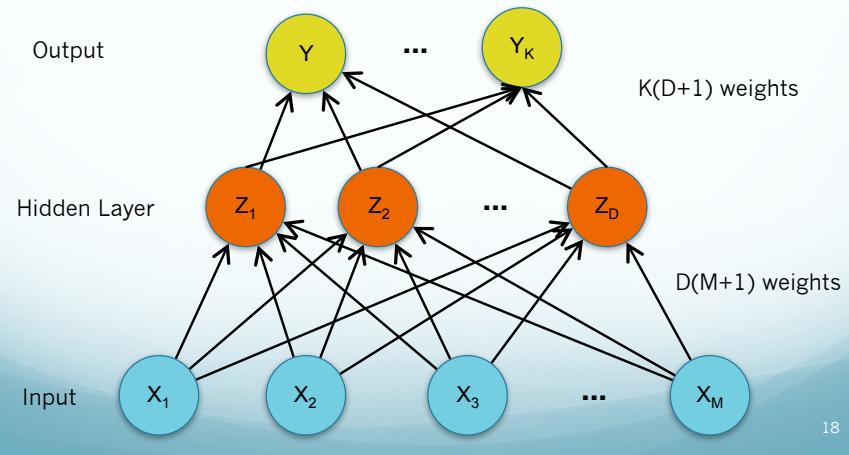
- **Fitting: what type of optimization algorithm?**
- Function: non-linear: linear combination of generalized linear functions
- Data: Data/model assumptions? How we use data?

How Will We Learn?

- Perceptron: a training method for generalized linear classifiers
 - Training method for linear classifiers
 - Minimize the error of the training data
 - Chain multiple Perceptrons together
 - Update rule:
$$w^{j+1} = w^j + \nabla f(x, y)$$
 - The real work will be in computing the gradient

17

Example: Multi-Class



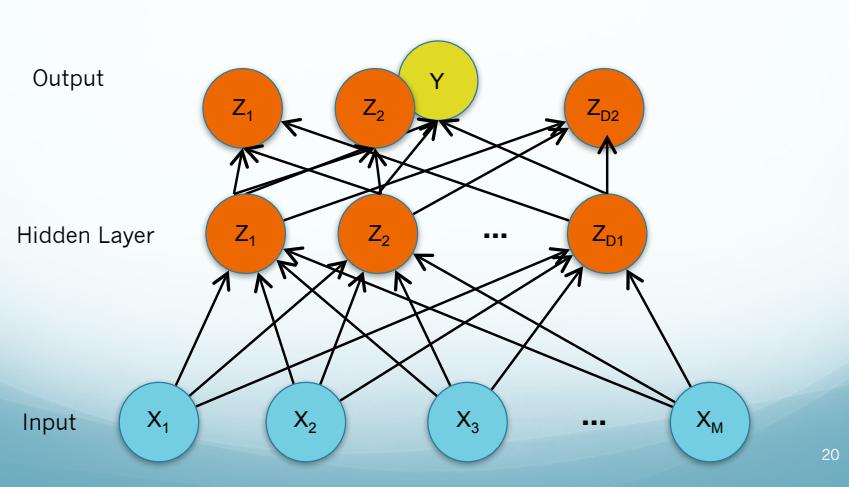
18

Network Terminology

- Input nodes: x
- Output node: y
- Hidden nodes: z
 - This network has 1 hidden layer
 - 2 layer network (two layers to learn)
- h for hidden nodes are called activation functions
- h for output depends on task
 - Identity for regression
 - Logistic for classification

19

Deep Networks



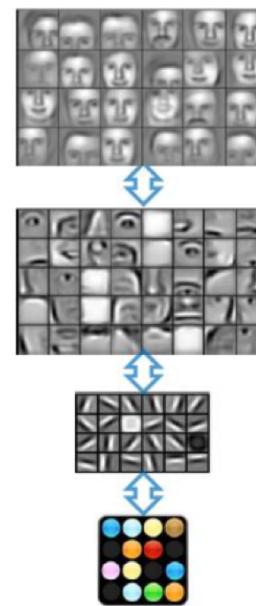
20

Deep Networks

- Learn multiple levels of features at higher and higher abstractions
- Same learning techniques
 - Just more complex gradients

21

Feature representation



3rd layer
"Objects"

2nd layer
"Object parts"

1st layer
"Edges"

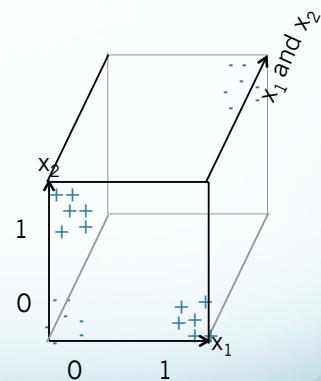
Pixels

Example:
Image
Processing

22

An Non-linear Example

- Consider the xor function
 - $y(x) = 1$ iff $x_1 \text{ xor } x_2$
- Clearly non-linear
 - No values for w will produce desired output
- We could solve this by adding a new feature
 - $x_3 = x_1 \text{ xor } x_2$

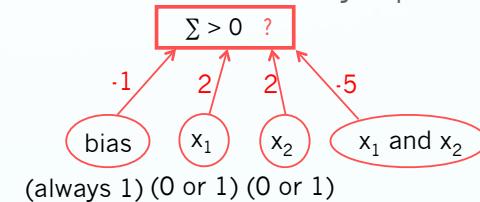


Example from Jason Eisner

23

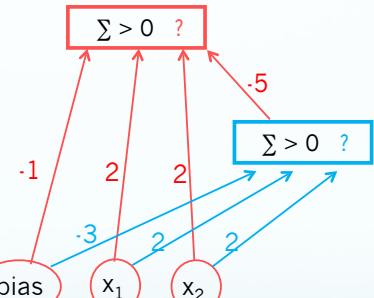
The Neural Network Solution

- Learn new features that are linearly separable
 - $\Sigma > 0$?
 - bias
 - x_1
 - x_2
 - $x_1 \text{ and } x_2$
- We now have a linear classifier for XOR
- How do we learn these features?



24

The Neural Network Solution



- The new features are learned by linear classifiers
 - All other hidden nodes (not shown) just replicate input
- The activation function makes the feature 1 or 0

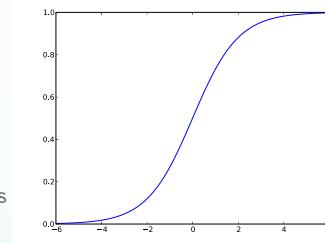
25

Non-linear Activation Functions

- What non-linear function should we use for activation function h ?
 - Typically use sigmoid functions
 - Logistic function

$$g_a(x) = \frac{1}{1 + e^{-ax}}$$

- Each hidden node has a threshold for activation
 - Will be 0 and then quickly transition to 1
- This is what we use when we stack Perceptrons
- This is why we think of hidden nodes as features
 - They are off and then when enough input they turn on
 - Learning input weights turns on the feature!



26

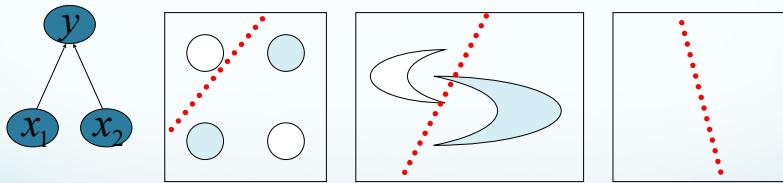
Hypothesis Class

- What can a neural network learn?
 - Obviously highly non-linear outputs
- Universal approximators
 - With enough hidden layers and hidden nodes a neural network can model any continuous function on compact input domain (some number of inputs)
 - The power of the networks depends on its structure
 - General result independent of activation functions

27

Decision Boundary

- 0 hidden layers: linear classifier
 - Hyperplanes

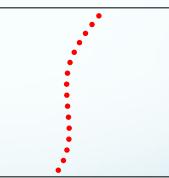
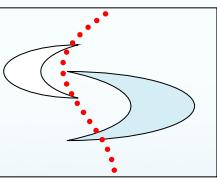
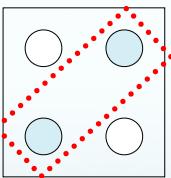
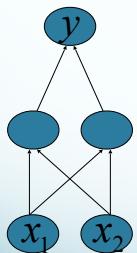


Example from Eric Postma via Jason Eisner

28

Decision Boundary

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

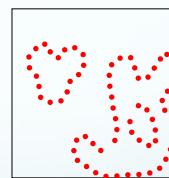
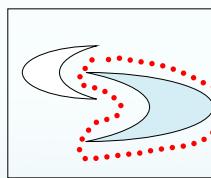
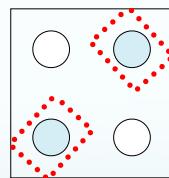
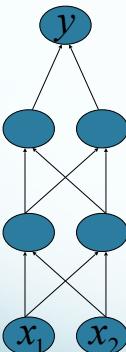


Example from to Eric Postma via Jason Eisner

29

Decision Boundary

- 2 hidden layers
 - Combinations of convex regions



Example from to Eric Postma via Jason Eisner

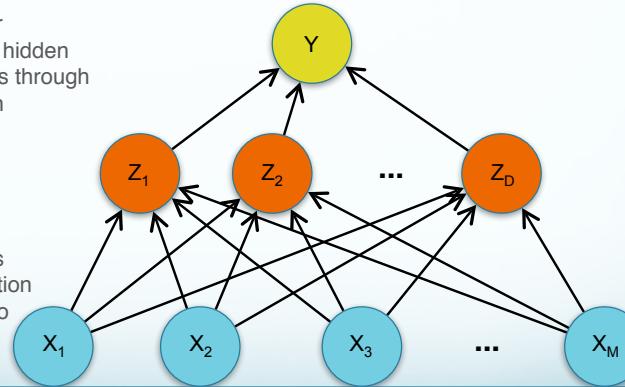
30

Prediction

Compute linear combination of hidden nodes and pass through logistic function

Hidden nodes are now new features

Compute each hidden node as linear combination of x and pass to logistic



- Forward propagation through the network

31

Classification Objective

- Define an error function and minimize
- Cross entropy error function

$$E(\mathbf{w}) = - \sum_{i=1}^N \{y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)\}$$

- This arises naturally when we consider a logistic probability model and take the negative log likelihood
 - Details in the book

32

Regression Objective

- For regression we use the sum of squares error

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

- If we assume a Gaussian model for y , the error function arises from maximizing the likelihood function
 - We saw the same thing for linear regression

33

Combined Model

- Writing the generalized linear classifier with generalized non-linear basis functions

$$y(\mathbf{x}; \mathbf{w}) = h^{(2)} \left(\sum_{j=1}^M w_j^{(2)} h^{(1)} \left(\sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_0^{(2)} \right)$$

- $h^{(1)}$ is the non-linear function for the basis function
- $h^{(2)}$ is the non-linear function for the output
- $w^{(1)}$ are the parameters for the basis function
- $w^{(2)}$ are the parameters for the linear model
- w_0 are the bias parameters (shown here for clarity)

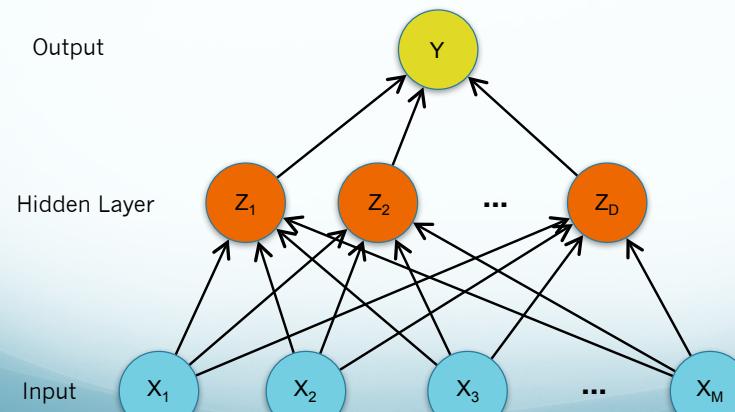
34

Training

- Prediction is relatively easy
- Learning is where the magic happens
- Strategy: compute the gradient of the objective function
 - Similar to perceptron
 - Gradient based update

35

Graphical Representation



36

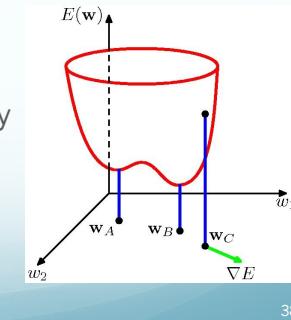
Training

- Prediction is relatively easy
- Learning is where the magic happens
- Strategy: compute the gradient of the objective function
 - Similar to perceptron
 - Gradient based update
- For the moment: assume black box computes gradient

37

Gradient Based Optimization

- The objective function is now non-convex
- Gradient based optimization NOT guaranteed to find global optimum
- For now: use gradient stochastic gradient and hope for the best
- Next time: tricks for non-convexity key to learning good networks



38

Computing the Gradient

- Next time: discuss methods for computing gradients

39

Algorithm: Neural Network

- Train: Given examples X and Y
 - Y can be multiple outputs
 - Define a network structure
 - ex. 2 layer feed forward, D nodes in hidden layer
 - Learn parameters w
- Predict: given example x
 - For 2 layer feed forward, compute output as

$$y(\mathbf{x}; \mathbf{w}) = h^{(2)} \left(\sum_{j=1}^D w_j^{(2)} h^{(1)} \left(\sum_{i=1}^M w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_0^{(2)} \right)$$

40

Multi-Layer Perceptrons

Fitting a function to data

- Fitting: gradient based optimization
- Function: non-linear: linear combination of generalized linear functions
 - Universal approximations
 - can model any continuous function on compact input domain (some number of inputs)
- Data: Batch training using stochastic methods

Next Time
Deep(er) Networks