Slides adapted from Prof Carpuat and Duraiswami



Furong Huang / furongh@cs.umd.edu



UNIVERSITY OF MARYLAND

#### **Neural Networks**

# Last time

What are Neural Networks? (key components)

How to make a prediction given an input?

Why are neural networks powerful?



### **Neural Networks**

### Last time

What are Neural Networks?
Multilayer perceptron
How to make a prediction given an input?
Simple matrix operations + nonlinearities
Why are neural networks powerful?
Universal function approximators!

### Today

how to train neural networks?



# Example: a neural network to solve the XOR problem



#### Example

In new space, the examples are linearly separable!



#### Example

.The final net





### **Activation Functions**



### **Activation Functions**

• Sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{1+e^x}$ , logistic function

• Derivative: 
$$\frac{d}{dx}\sigma(x) = \frac{e^x}{1+e^x} = \sigma(x) (1 - \sigma(x))$$



# **Activation Functions**

- Relu:  $relu(x) = x^+ = max(0, x)$  rectified linear unit
  - A smooth approximation: softplus  $f(x) = \log(1 + \exp x)$
  - The derivative of softplus is logistic function  $f'(x) = \frac{e^x}{1+e^x}$



#### **Exercise**

 Learn softmax function <u>https://en.wikipedia.org/wiki/Softmax\_fu</u> <u>nction</u>



## Forward Propagation: given input x, compute network output

#### **Algorithm 24** TwoLayerNetworkPredict( $\mathbf{W}, v, \hat{x}$ )

- 1: for i = 1 to number of hidden units do
- $_{2:}$   $h_i \leftarrow \tanh(\boldsymbol{w}_i \cdot \hat{\boldsymbol{x}})$

// compute activation of hidden unit i

- 3: end for
- 4: return  $v \cdot h$

// compute output unit



### **Neural Network Training**

#### **Backpropagation algorithm**

### Gradient descent + Chain rule



# Recall: gradient descent for linear classifiers





# What's our Training Objective?

We'll consider the following objective

$$\min_{\mathbf{W},v} \quad \sum_{n} \frac{1}{2} \left( y_n - \sum_{i} v_i f(\boldsymbol{w}_i \cdot \boldsymbol{x}_n) \right)^2$$

i.e. our goal is to find parameters  $\boldsymbol{W},$  v that minimize squared error

Other objectives are possible (e.g., other loss functions, add regularizer)



#### **Backprop in a 2-layer network**

**Algorithm 25** TwoLayerNetworkTrain( $\mathbf{D}$ ,  $\eta$ , K, *MaxIter*)  $W \leftarrow D \times K$  matrix of small random values // initialize input layer weights  $_{2}$ :  $v \leftarrow K$ -vector of small random values // initialize output layer weights  $_{3:}$  for *iter* = 1 ... *MaxIter* do  $\mathbf{G} \leftarrow D \times K$  matrix of zeros // initialize input layer gradient  $g \leftarrow K$ -vector of zeros // initialize output layer gradient 5: for all  $(x,y) \in \mathbf{D}$  do 6: for  $i = \tau$  to K do 7:  $a_i \leftarrow w_i \cdot \hat{x}$ 8:  $h_i \leftarrow tanh(a_i)$ // compute activation of hidden unit i9: end for 10:  $\hat{y} \leftarrow \boldsymbol{v} \cdot \boldsymbol{h}$ // compute output unit 11:  $e \leftarrow y - \hat{y}$ // compute error 12:  $g \leftarrow g - eh$ // update gradient for output layer 13: for i = 1 to K do 14:  $\mathbf{G}_i \leftarrow \mathbf{G}_i - ev_i(1 - \tanh^2(a_i))\mathbf{x}$ // update gradient for input layer 15: end for 16: end for 17:  $\mathbf{W} \leftarrow \mathbf{W} - \eta \mathbf{G}$ // update input layer weights 18:  $v \leftarrow v - \eta g$ // update output layer weights 19: 20: end for 21: return W, v



#### **Backprop in a 2-layer network**





# What's our Training Objective?

We'll consider the following objective

$$\min_{\mathbf{W},v} \quad \sum_{n} \frac{1}{2} \left( y_n - \sum_{i} v_i f(\boldsymbol{w}_i \cdot \boldsymbol{x}_n) \right)^2$$

i.e. our goal is to find parameters  $\boldsymbol{W},$  v that minimize squared error

Other objectives are possible (e.g., other loss functions, add regularizer)



# Gradient of objective w.r.t. output layer weights v



# **Gradient of objective** w.r.t. hidden unit weights $w_i$

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2} \left( y - \sum_{i} v_{i} f(\boldsymbol{w}_{i} \cdot \boldsymbol{x}) \right)^{2}$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_{i}} = \frac{\partial \mathcal{L}}{\partial f_{i}} \frac{\partial f_{i}}{\partial \boldsymbol{w}_{i}}$$
Chain rule
$$\frac{\partial \mathcal{L}}{\partial f_{i}} = -\left( y - \sum_{i} v_{i} f(\boldsymbol{w}_{i} \cdot \boldsymbol{x}) \right) v_{i} = -ev_{i}$$

$$\frac{\partial f_{i}}{\partial \boldsymbol{w}_{i}} = f'(\boldsymbol{w}_{i} \cdot \boldsymbol{x})\boldsymbol{x}$$

 $\nabla_{w_i} = -ev_i f'(w_i \cdot x) x$ 

**)EAS** (This is on one example only)

// // // // //

#### **Backprop in a 2-layer network**

<b>ιν</b> ( <b>D</b> , η, Κ, MaxIter)
es // initialize input layer weights
// initialize output layer weights
// initialize input layer gradient
// initialize output layer gradient
// compute activation of hidden unit $i$
// compute output unit
// compute error
// update gradient for output layer
// update gradient for input layer
<pre>// update input layer weights</pre>
// update output layer weights

Forward

Update

Update

gradients

parameters

propagation



# Tricky issues with neural network training

# Sensitive to initialization

Objective is non-convex, many local optima In practice: start with random values rather than zeros

#### Many other hyperparameters

Number of hidden units (and potentially hidden layers) Gradient descent learning rate

Stopping criterion



# Neural networks vs. linear classifiers

Advantages of Neural Networks: More expressive Less feature engineering

Inconvenients of Neural Networks:

Harder to train

Harder to interpret



### **Neural Network Architectures**

# We focused on a **2-layer feedforward** network

### Other architectures are possible

- More than 2 layers (aka deep learning)
- Recurrent network (i.e. network has cycles)
- Can still be trained with backpropagation
  - But more issues arise when networks get more complex (e.g., vanishing gradients)



# Try different architectures and training parameters here:

# http://playground.tensorflow.org



### What you should know

What are Neural Networks?

Multilayer perceptron

How to make a prediction given an input?

Forward propagation: Simple matrix operations + nolinearities

Why are neural networks powerful?

Universal function approximators!

How to train neural networks?

The backpropagation algorithm





Furong Huang 3251 A.V. Williams, College Park, MD 20740 301.405.8010 / furongh@cs.umd.edu