

Non-linear hypotheses

### **Non-linear Classification**



#### **Computer Vision: Car detection**



Testing:



### What is this?





Andrew Ng









```
Andrew Ng
```



### Model representation I

### **Neural Networks**

Origins: Algorithms that try to mimic the brain. Was very widely used in 80s and early 90s; popularity diminished in late 90s.

Recent resurgence: State-of-the-art technique for many applications

#### Neuron in the brain



#### Neurons in the brain



[Credit: US National Institutes of Health, National Institute on Aging]

#### Neuron model: Logistic unit



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

Sigmoid (logistic) activation function.

#### **Neural Network**



 $a_i^{(j)} =$  "activation" of unit i in layer j

 $\Theta^{(j)} = \text{matrix of weights controlling} \\ \text{function mapping from layer } j \text{ to} \\ \text{layer } j+1$ 

$$\begin{aligned} a_1^{(2)} &= g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3) \\ a_2^{(2)} &= g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3) \\ a_3^{(2)} &= g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3) \\ h_{\Theta}(x) &= a_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)}) \end{aligned}$$

If network has  $s_j$  units in layer j,  $s_{j+1}$  units in layer j + 1, then  $\Theta^{(j)}$  will be of dimension  $s_{j+1} \times (s_j + 1)$ .



## Model

representation II

### Forward propagation: Vectorized implementation



$$a_{1}^{(2)} = g(\Theta_{10}^{(1)}x_{0} + \Theta_{11}^{(1)}x_{1} + \Theta_{12}^{(1)}x_{2} + \Theta_{13}^{(1)}x_{3})$$

$$a_{2}^{(2)} = g(\Theta_{20}^{(1)}x_{0} + \Theta_{21}^{(1)}x_{1} + \Theta_{22}^{(1)}x_{2} + \Theta_{23}^{(1)}x_{3})$$

$$a_{3}^{(2)} = g(\Theta_{30}^{(1)}x_{0} + \Theta_{31}^{(1)}x_{1} + \Theta_{32}^{(1)}x_{2} + \Theta_{33}^{(1)}x_{3})$$

$$h_{\Theta}(x) = g(\Theta_{10}^{(2)}a_{0}^{(2)} + \Theta_{11}^{(2)}a_{1}^{(2)} + \Theta_{12}^{(2)}a_{2}^{(2)} + \Theta_{13}^{(2)}a_{3}^{(2)})$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \qquad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$
$$z^{(2)} = \Theta^{(1)} x$$
$$a^{(2)} = g(z^{(2)})$$
$$Add \ a_0^{(2)} = 1.$$
$$z^{(3)} = \Theta^{(2)} a^{(2)}$$
$$h_{\Theta}(x) = a^{(3)} = g(z^{(3)})$$

Andrew Ng

### **Neural Network learning its own features**



#### **Other network architectures**





# Examples and intuitions I

### Non-linear classification example: XOR/XNOR

 $x_1$ ,  $x_2$  are binary (0 or 1).



### Simple example: AND

 $x_1, x_2 \in \{0, 1\}$  $y = x_1 \text{ AND } x_2$ 



$x_1$	$x_2$	$h_{\Theta}(x)$
0	0	
0	1	
1	0	
1	1	

1.0 g(z)

### **Example: OR function**





## Examples and intuitions II

$$x_1 \text{ AND } x_2$$

$$x_1 \text{ OR } x_2$$

### Negation:



$$h_{\Theta}(x) = g(10 - 20x_1)$$

(NOT  $x_1$ ) AND (NOT  $x_2$ )

Putting it together:  $x_1$  XNOR  $x_2$ 











$x_1$	$x_2$	$a_1^{(2)}$	$a_2^{(2)}$	$h_{\Theta}(x)$
0	0			
0	1			
1	0			
1	1			

#### **Neural Network intuition**



### Handwritten digit classification



### Handwritten digit classification



Andrew Ng



### Multi-class classification

### Multiple output units: One-vs-all.





Andrew Ng


Neural Networks: Learning

## **Cost function**

#### Neural Network (Classification)



#### **Binary classification**

y = 0 or 1

1 output unit

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

L = total no. of layers in network

 $s_l = \max l$  no. of units (not counting bias unit) in layer l

$$\underbrace{ \textit{Multi-class classification}}_{y \in \mathbb{R}^{K} \textit{ E.g. } \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \\ \textit{pedestrian car motorcycle truck} }$$

#### K output units

#### **Cost function**

#### Logistic regression:

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^{m} y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_{j}^{2}$$

#### Neural network:

$$h_{\Theta}(x) \in \mathbb{R}^{K} \quad (h_{\Theta}(x))_{i} = i^{th} \text{ output}$$

$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^{m} \sum_{k=1}^{K} y_{k}^{(i)} \log(h_{\Theta}(x^{(i)}))_{k} + (1 - y_{k}^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))_{k}) \right]$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_{l}} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^{2}$$



Neural Networks: Learning

## Gradient descent

## Have some function $J(\theta_0, \theta_1)$ Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

#### **Outline:**

- Start with some  $heta_0, heta_1$
- Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$ until we hopefully end up at a minimum





#### **Gradient descent algorithm**

repeat until convergence {

$$\theta_{j} := \theta_{j} - \alpha \frac{\partial}{\partial \theta_{j}} J(\theta_{0}, \theta_{1}) \quad \text{(for } j = 0 \text{ and } j = 1\text{)}$$
}



## Neural Networks: Learning

Backpropagation algorithm

#### **Gradient computation**

$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^{m} \sum_{k=1}^{K} y_k^{(i)} \log h_\theta(x^{(i)})_k + (1 - y_k^{(i)}) \log(1 - h_\theta(x^{(i)})_k) \right] \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_j^{(l)})^2$$

 $\min_{\Theta} J(\Theta)$ 

Need code to compute:

$$\begin{array}{l} - J(\Theta) \\ - \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) \end{array} \end{array}$$

#### **Backpropagation algorithm**

Training set  $\{(x^{(1)}, y^{(1)}), ..., (x^{(m)}, y^{(m)})\}$ Set  $\triangle_{ij}^{(l)} = 0$  (for all l, i, j). For i = 1 to m

Set  $a^{(1)} = x^{(i)}$ 

Perform forward propagation to compute  $a^{(l)}$  for l = 2, 3, ..., LUsing  $y^{(i)}$ , compute  $\delta^{(L)} = a^{(L)} - y^{(i)}$ Compute  $\delta^{(L-1)}, \delta^{(L-2)}, ..., \delta^{(2)}$   $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$   $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}$  if  $j \neq 0$   $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$  if j = 0 $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$ 



## Neural Networks: Learning

## Random

## initialization

#### Initial value of $\Theta$

For gradient descent and advanced optimization method, need initial value for ⊖. optTheta = fminunc(@costFunction, initialTheta, options)

Consider gradient descent

Set initialTheta = zeros(n,1) ?

#### **Zero initialization**



$$\Theta_{ij}^{(l)} = 0$$
 for all  $i, j, l$ .

After each update, parameters corresponding to inputs going into each of two hidden units are identical.

#### **Random initialization: Symmetry breaking**

Initialize each  $\Theta_{ij}^{(l)}$  to a random value in  $[-\epsilon, \epsilon]$  (i.e.  $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$ )

E.g.



## Neural Networks: Learning

# Putting it together

#### Training a neural network

Pick a network architecture (connectivity pattern between neurons)







No. of input units: Dimension of features  $x^{(i)}$ 

No. output units: Number of classes

Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

#### Training a neural network

- 1. Randomly initialize weights
- 2. Implement forward propagation to get  $h_{\Theta}(x^{(i)})$  for any  $x^{(i)}$
- 3. Implement code to compute cost function  $J(\Theta)$
- 4. Implement backprop to compute partial derivatives  $\frac{\partial}{\partial \Theta_{ik}^{(l)}} J(\Theta)$

#### for i = 1:m

Perform forward propagation and backpropagation using example  $(x^{(i)}, y^{(i)})$ 

(Get activations  $a^{(l)}$  and delta terms  $\delta^{(l)}$  for l = 2, ..., L).



#### Training a neural network

5. Use gradient checking to compare  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$  computed using backpropagation vs. using numerical estimate of gradient of  $J(\Theta)$ .

Then disable gradient checking code.

6. Use gradient descent or advanced optimization method with backpropagation to try to minimize  $J(\Theta)$  as a function of parameters  $\Theta$ 





Machine Learning

## Neural Networks: Learning

Backpropagation example: Autonomous driving







### Neural Networks: Feature Learning

## Autoencoder

#### **Autoencoders (and sparsity)**







Input

Features

#### Sparse autoencoders





Input

Features

#### **Unsupervised Feature Learning/Self-taught learning**





#### **Unsupervised pre-training + Fine-tuning**





## Neural Networks: Feature Learning

# Deep Learning
## **Unsupervised pre-training + Fine-tuning**







Input Softmax (Features II) classifier

Input

Features I

Output

Input Features II (Features I) Output

Andrew Ng







Softmax Input (Features II) classifier

Input

Features I Output

Input Features II Output (Features I)



Andrew Ng