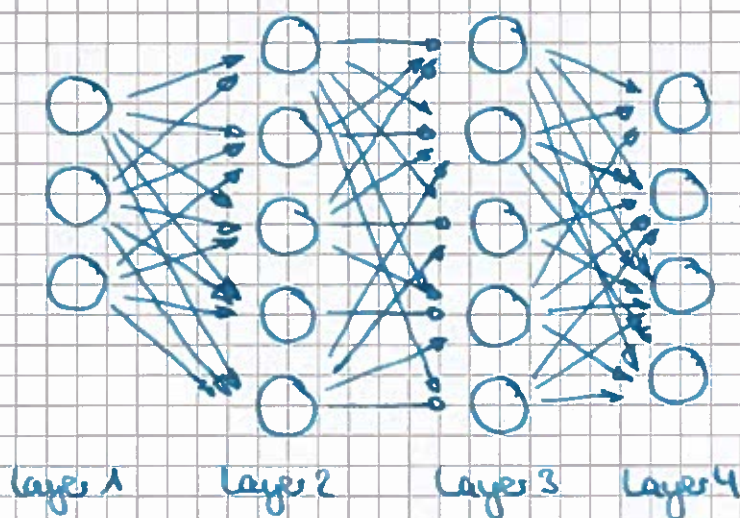


9. Neural Networks

- Cost function & Back Propagation

Neural Network - Cost function & Backpropagation Chapter 3 [1]

Neural Network (Classification)



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

L = total # of layers in Network $\Rightarrow L=4$

S_L = # of units (not counting bias unit) in layer L

$$S_1=3, S_2=5, S_3=5, S_4=4=S_L$$

Binary Classification

$$y = 0 \text{ or } 1 \Rightarrow \begin{cases} 0 \\ 1 \end{cases} \rightarrow h\theta(x)$$

1 Output unit

$$h\theta(x) \Rightarrow \mathbb{R}$$

$$S_L = 1 \Rightarrow K=1$$

of output units in Binary classification = 1.

Multi-Class classification (K classes)

$$y \in \mathbb{R}^K \quad \text{eg.} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

K output units

$$h\theta(x) \Rightarrow \mathbb{R}^K$$

$$S_L = K \quad (K \geq 3)$$

If only 2 output classes use Binary classification

logistic regression:

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

Neural Network:

$h_{\Theta}(x) \in \mathbb{R}^k \Rightarrow (h_{\Theta}(x))_i = i^{\text{th}}$ output

$$J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\Theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \cdot \log(1 - h_{\Theta}(x^{(i)}))_k \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{3l} \sum_{j=1}^{3l+1} (\Theta_{ji}^{(l)})^2$$

called weight decay

K sum of output unit

for neural networks, cost function is a generalization. It generates K outputs $\Rightarrow K$ dimensional vector

$$\begin{bmatrix} 0 \\ x \\ 0 \\ x \\ 0 \end{bmatrix} \Rightarrow y_k \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

4 output units in final layer

$h_{\Theta}(x)$ is K dim. vector $\Rightarrow h_{\Theta}(x)_i$ refers to i^{th} value in vector

Back propagation Function

3

Gradient Computation:

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

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

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

Need to compute

$-J(\Theta)$

$$-\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = 0 \Rightarrow \Theta_{ij}^{(l)} \in \mathbb{R}$$

Given one training example (x, y)

Forward propagation

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

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

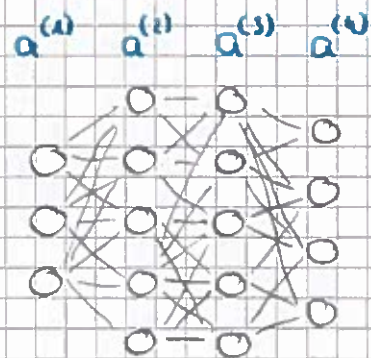
$$a^{(2)} = g(z^{(2)}) \text{ (add } a_0^{(2)} \text{) bias term}$$

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)}) \text{ (add } a_0^{(3)} \text{)}$$

$$z^{(4)} = \Theta^{(3)} a^{(3)}$$

$$a^{(4)} = g(z^{(4)}) = h_{\Theta}(x) = g(z^{(4)})$$



Back propagation algorithm

[4]

Gradient computation: Backpropagation algorithm

Intuition: $\delta_j^{(L)} = \text{"error" of node } j \text{ in layer } L$.

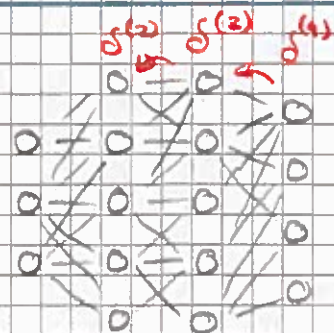
For each output unit (layer $L=4$)

$$\delta_j^{(L)} = \underbrace{a_j^{(L)}}_{(h \circ w)_j} - y_j$$

→ all vectors

$$\delta^{(L)} = a^{(L)} - y$$

dimension of vector is equal to number of output units



$$\delta^{(2)} = (\Theta^{(3)})^T \delta^{(4)} * g'_i(z^{(3)})$$

→ $a^{(3)} * (1 - a^{(3)})$ vectorized

$$\delta^{(2)} = (\Theta^{(3)})^T \delta^{(4)} * g'_i(z^{(3)})$$

(No $\delta^{(1)}$)

$a^{(2)} * (1 - a^{(2)})$ vectorized.

$$\hookrightarrow \frac{\partial}{\partial \Theta_{ij}} J(\Theta) = a_j^{(L)} \cdot \delta_i^{(L+1)} \quad (\text{ignoring } \lambda; \text{ if } \lambda = \emptyset)$$

Back propagation algorithm:

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

Set $\Delta_{ij}^{(L)} = \emptyset$ (for all L, i, j).

For $i = 1:m$

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

perform forward propagation $a^{(L)} - y^{(i)}$

L = layer

Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

j = node in the layer

i = error of affected node in layer

$$\Delta_{ij}^{(L)} := \Delta_{ij}^{(L)} + a_j^{(L)} \cdot \delta_i^{(L+1)}$$

=> accumulate partial derivative terms

end

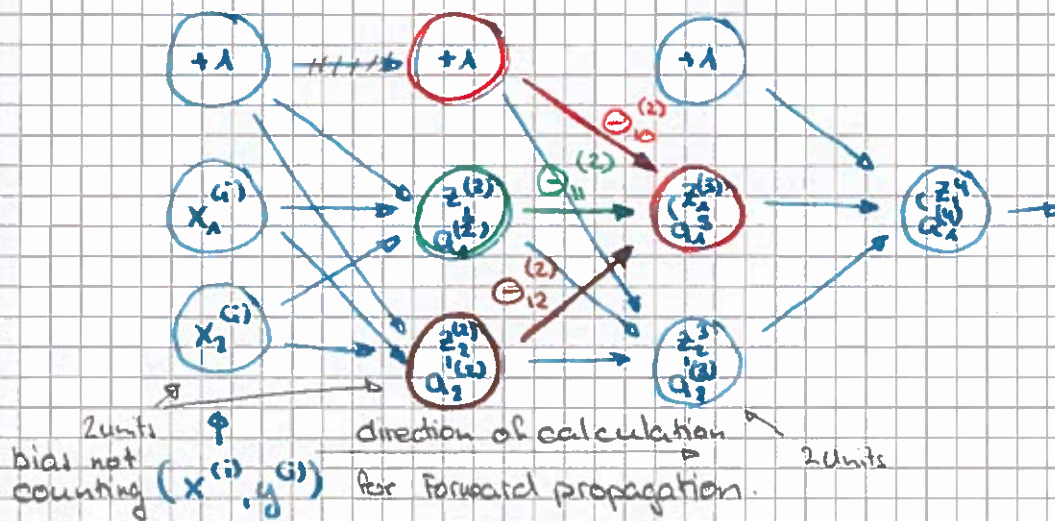
$$\hookrightarrow \text{vectorized } \Delta^{(L)} := \Delta^{(L)} + \delta^{(L+1)} \cdot (a^{(L)})^T$$

$$D_{ij}^{(L)} := \frac{1}{m} \Delta_{ij}^{(L)} + \lambda \Theta_{ij}^{(L)} \quad \text{if } j \neq \emptyset$$

$$D_{ij}^{(L)} := \frac{1}{m} \Delta_{ij}^{(L)} \quad \text{if } j = \emptyset$$

$$\frac{\partial}{\partial \Theta_{ij}} J(\Theta) = D_{ij}^{(L)}$$

Forward propagation:



calculate: $z_1^{(3)} = \underbrace{\theta_{10}^{(2)} \cdot 1}_{\text{red}} + \underbrace{\theta_{11}^{(2)} a_1^{(2)}}_{\text{green}} + \underbrace{\theta_{12}^{(2)} a_2^{(2)}}_{\text{red}}$

What is backpropagation doing?

$$J(\Theta) = \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(\text{h}\Theta(x^{(i)}) - y^{(i)}) + (1 - y^{(i)}) \log(1 - (\text{h}\Theta(x^{(i)}))) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\Theta_{ji}^{(l)})^2$$

Focus on a single example $x^{(i)}, y^{(i)}$, the case of 1 output unit, and ignoring regularization ($\lambda = 0$),

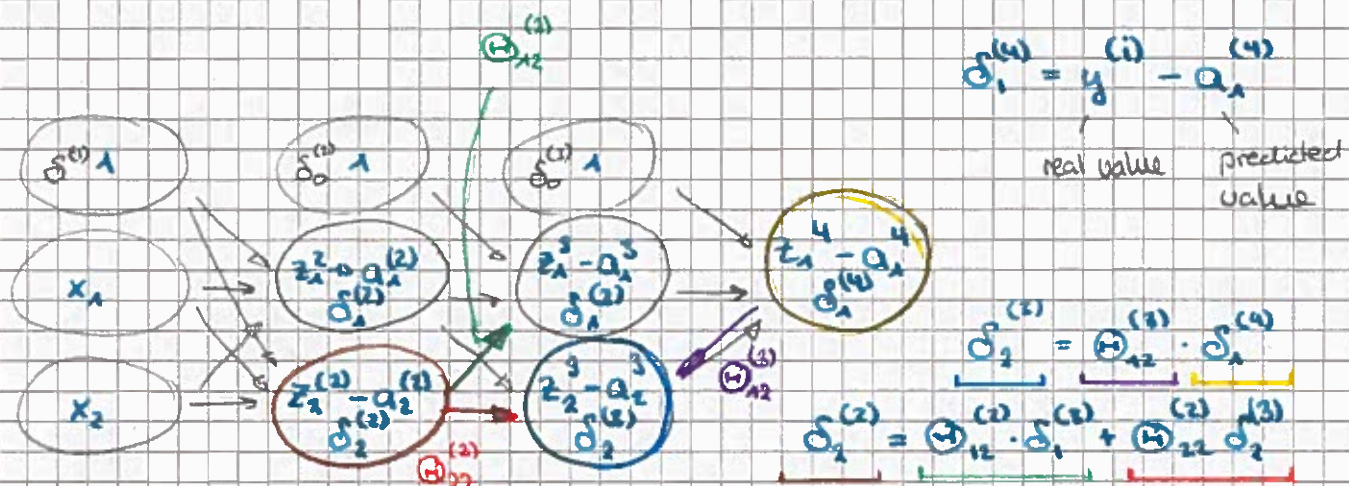
$$\text{cost}(i) = y^{(i)} \log(\text{h}\Theta(x^{(i)}) - y^{(i)}) + (1 - y^{(i)}) \log(1 - (\text{h}\Theta(x^{(i)})))$$

\hookrightarrow think of $\text{cost}(i) \approx (\text{h}\Theta(x^{(i)}) - y^{(i)})^2$

how well is the network doing on example i ?

Backward propagation intuition

[6]



$\delta_j^{(l)}$ = "error" or cost for $a_j^{(l)}$ (unit j in layer l)

Formally, $\delta_j^{(l)} = \frac{\partial}{\partial a_j^{(l)}} \cdot \text{cost}(l)$ for $j \geq 0$, where

$$\text{cost}(i) = y^{(i)} \log(h(x^{(i)})) + (1 - y^{(i)}) \log(h(x^{(i)}))$$

Advanced optimization:

Function [jval gradient] = costfunction(theta)

... $\rightarrow \mathbb{R}^{n+1}$ vector \mathbb{R}^{n+1} vectors for Linear Regression

optTheta = fminunc(@costfunction, initialTheta, options)

Neural Network (L=4)

$\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ - matrices (Theta1, Theta2, Theta3)

$D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (D1, D2, D3)

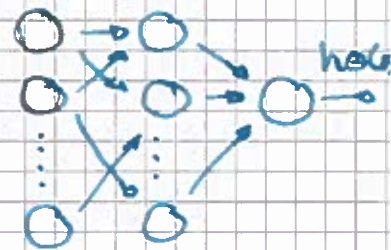
"Unroll" into Vectors

Example:

Input layer 10 units
 $s_1 = 10, s_2 = 10, s_3 = 1$ - output layer 1 unit
 hidden layer 10 units

$\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$

$D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$



in Octave:

create a vector $\begin{cases} \text{thetaVec} = [\text{Theta1}(:); \text{Theta2}(:); \text{Theta3}(:)] \\ \text{DVec} = [D1(:); D2(:); D3(:)] \end{cases}$

Theta1 = reshape(thetaVec(1:110), 10, 11);

Theta2 = reshape(thetaVec(111:220), 10, 11);

Theta3 = reshape(thetaVec(221:231), 1, 11);

Learning Algorithm

- Have initial Parameters $\Theta^{(1)}$, $\Theta^{(2)}$, $\Theta^{(3)}$
- Unroll to get initialTheta to pass to `fminunc (@costFunction, initialTheta, options)`

Function [jval gradientVec] = costFunction(thetaVec)

From thetaVec get $\Theta^{(1)}$, $\Theta^{(2)}$, $\Theta^{(3)}$ \rightarrow reshape to get matrices

Use forward propagation/backward prop. to compute

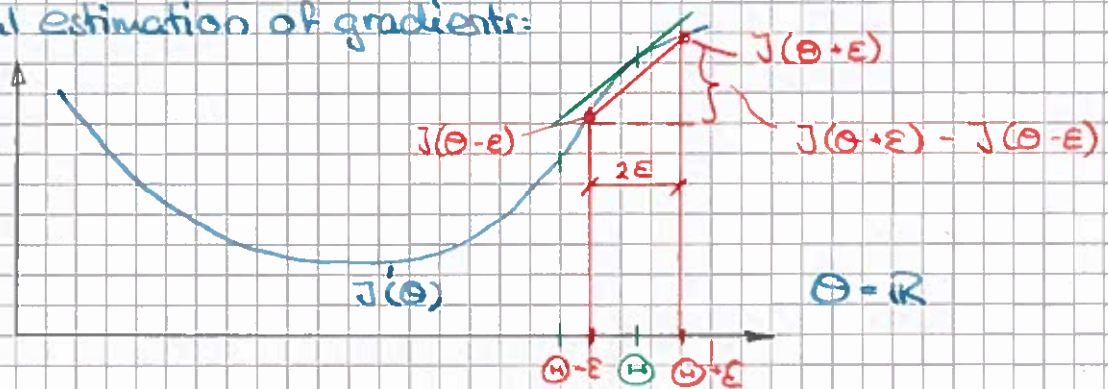
$D^{(1)}$, $D^{(2)}$, $D^{(3)}$ and $J(\Theta)$

Unroll $D^{(1)}$, $D^{(2)}$, $D^{(3)}$ to get gradientVec

Gradient checking

9

Numerical estimation of gradients:



$$\frac{\partial}{\partial \theta} \cdot J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

$$\epsilon = 10^{-4}$$

2 sided difference

more accurate - use this

$$\frac{J(\theta + \epsilon) - J(\theta)}{\epsilon}$$

1 sided difference

Implementation:

$$\text{gradApprox} = \frac{J(\text{theta} + \text{EPSILON}) - J(\text{theta} - \text{EPSILON})}{(2 \cdot \text{EPSILON})}$$

Parameter Vector:

$\theta \in \mathbb{R}^n$ (e.g. θ is an "unrolled" version of $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$)

$$\theta = \theta_1, \theta_2, \dots, \theta_n$$

$$\frac{\partial}{\partial \theta_1} \cdot J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \theta_2} \cdot J(\theta) = \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$$

$$\vdots$$
$$\frac{\partial}{\partial \theta_n} \cdot J(\theta) = \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_n - \epsilon)}{2\epsilon}$$

Octave implementation:

for i = 1:n

thetaPlus = theta;

↳ unrolled version

thetaPlus(i) = thetaPlus(i) + EPSILON;

thetaMinus = theta;

thetaMinus(i) = thetaMinus(i) - EPSILON

gradApprox(i) = (J(thetaPlus) - J(thetaMinus)) / (2 * EPSILON)

end; $\frac{\partial}{\partial \theta_i} (J(\theta))$

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{i+\epsilon} \\ \vdots \\ \theta_n \end{bmatrix} \rightarrow \theta_{i+\epsilon}$$

Check that **gradApprox** \approx **DVec**

↳ numerically computed

↳ backwards propagation

Implementation Note:

- Implement back prop to compute DVec (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$)
- Implement numerical gradient check to compute "gradApprox"
- Make sure they give similar values.
- Turn off gradient checking. Use back prop code for learning
↳ DVec computationally efficient

Important:

- Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient checking computation on every iteration, or gradient descent (or in the inner loop of cost function) your code will be very slow.

Initial Value of Θ

For gradient descent and advanced optimization method, need initial value for Θ

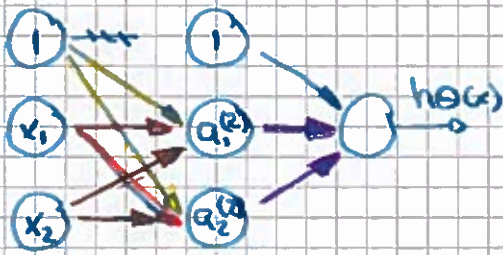
$$\text{optTheta} = \text{fminunc}(@\text{costfunction}, \text{initialTheta}, \text{options})$$

Consider gradient descent

$$\text{set initialTheta} = \text{zeros}(n, 1)$$

Zero initialization

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



Elements in same color will have the same value
we end up for all our training samples

$$\Rightarrow a_1^{(2)} = a_2^{(2)}, \text{ also } \delta_1^{(2)} = \delta_2^{(2)}$$

\Rightarrow same applies for the partial derivatives.

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

Even after the first iteration of gradient descent update

$$a_1^{(2)} = a_2^{(2)}$$

Random initialization: Symmetry breaking

Initialize each $\Theta_{ij}^{(l)}$ to a random value $[-\epsilon, \epsilon]$ (e.g. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)
 \rightarrow rand 10x11 Matrix (between 0 and 1)

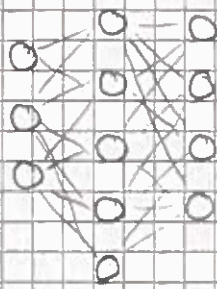
$$\text{Theta1} = \text{rand}(10, 11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$$

$$\text{Theta2} = \text{rand}(1, 11) * (2 * \text{INIT_EPSILON}) - \text{INIT_EPSILON};$$

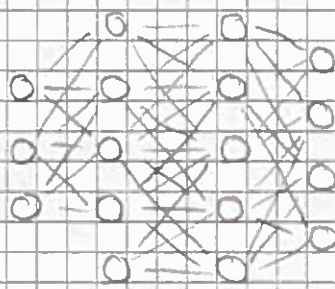
e.g. 10^{-4}

Training a neural network:

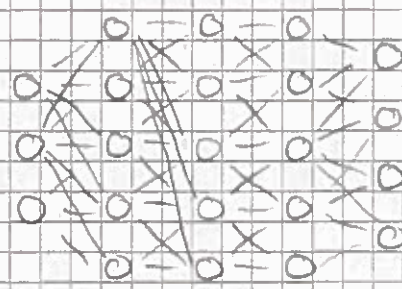
Pick the network architecture (connectivity patterns between neurons)



3 5 4



3 5 5 4



3 5 5 5 4

No. of input units:

Dimension of features $x^{(i)}$

No. of output units:

Number of classes

Reasonable default:

1 hidden layer, or if > 1 hidden layer, have same number of hidden layers units in every layer (usually the more the better)

=> multi class classification (e.g. $y \in \{1, 2, 3, 4, \dots, 10\}$)

$$y_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$y_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$y_5 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

eg. # features = 120 / output 0-9

input units = 120

hidden layer 2-3 * 120 units

output units 10

Training a neural network:

- 1) Randomly initialize weights small values near zero
- 2) Implement forward propagation $h(\theta(x^{(i)}))$ for any $x^{(i)}$
- 3) Implement code to compute cost function $J(\theta)$
- 4) Implement backpropagation to compute partial derivatives $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$

for $i = 1:m$ {

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

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

$$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l)} (a^{(l)})^T$$

}

...

compute $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$ \Rightarrow must also include regularization terms

- 5) Use gradient checking to compute $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$ 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 backprop to try to minimize $J(\theta)$ as a function of parameters θ

$$\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$$

$J(\theta)$ - non convex \rightarrow can get to local minimum which is not a problem.

