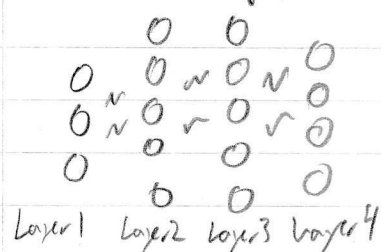


Machine Learning - Week 5 - Neural Networks: Learning

TL - Cost Function and Backpropagation

LI - Cost Function

Neural Network (Classification)



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

L = total # of layers in network

S_l = number of units (not counting bias unit) in layer l

Recall:

Binary Classification

$y = 0$ or 1

1 output unit

$$h_{\theta}(x) \in \mathbb{R}$$

$$S_L = 1$$

$$K = 1$$

Multi-class Classification (K classes)

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

person car motorcycle truck

$$K \text{ output units} \rightarrow h_{\theta}(x) \in \mathbb{R}^K$$

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

Cost Function

Recall: Logistic Regression:

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

Neural Network:

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

$$J(\theta) = -\frac{1}{n} \left[\sum_{i=1}^n \sum_{k=1}^K y_k^{(i)} \log (h_{\theta}(x^{(i)})_k) + (1 - y_k^{(i)}) \log (1 - h_{\theta}(x^{(i)})_k) \right] - \frac{\lambda}{2n} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

L2-Backpropagation Algorithm

Recall what $J(\theta)$ is (on previous page)

want $\min_{\theta} J(\theta)$

Need code to compute:

$$\rightarrow J(\theta)$$

$$\rightarrow \frac{d}{d\theta_{ij}} J(\theta)$$

Gradient computation:

Given one training example (x, y)

Forward Propagation:

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

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

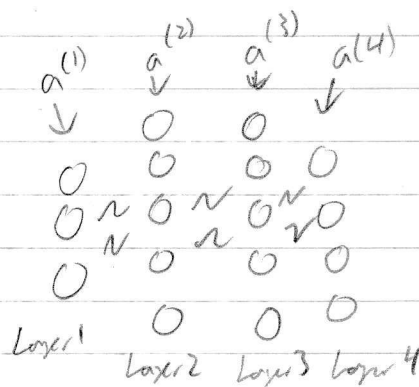
$$a^{(2)} = g(z^{(2)}) \text{ (add } a_0^{(2)})$$

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

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

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

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



Backpropagation Algorithm

Intuition: $d_j^{(l)}$ = "error" of node j in layer l .

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

$$d_j^{(4)} = a_j^{(4)} - y_j \quad (\text{vectorized: } d^{(4)} = a^{(4)} - y)$$

$$d^{(3)} = (\theta^{(3)})^T d^{(4)} \cdot g'(z^{(3)})$$

$$d^{(2)} = (\theta^{(2)})^T d^{(3)} \cdot g'(z^{(2)})$$

No $d^{(1)}$, since the first layer is the input layer (feature layer) and we don't want to change those values

$$\frac{d}{d\theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} d_i^{(l+1)} \quad (\text{ignoring } \lambda; \text{ if } \lambda=0)$$

2

Backpropagation algorithm.

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

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j). (Used to compute $\frac{d}{d\theta_{ij}^{(l)}} J(\theta)$)

For $i=1$ to m

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

Perform forward propagation to compute $a^{(l)}$ for $l=2, 3, \dots, L$

Using $y^{(i)}$ compute $d^{(L)} = a^{(L)} - y^{(i)}$

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

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} d_i^{(l+1)}$

Vectorized: $\Delta^{(l)} := \Delta^{(l)} + d^{(l+1)} (a^{(l)})^T$

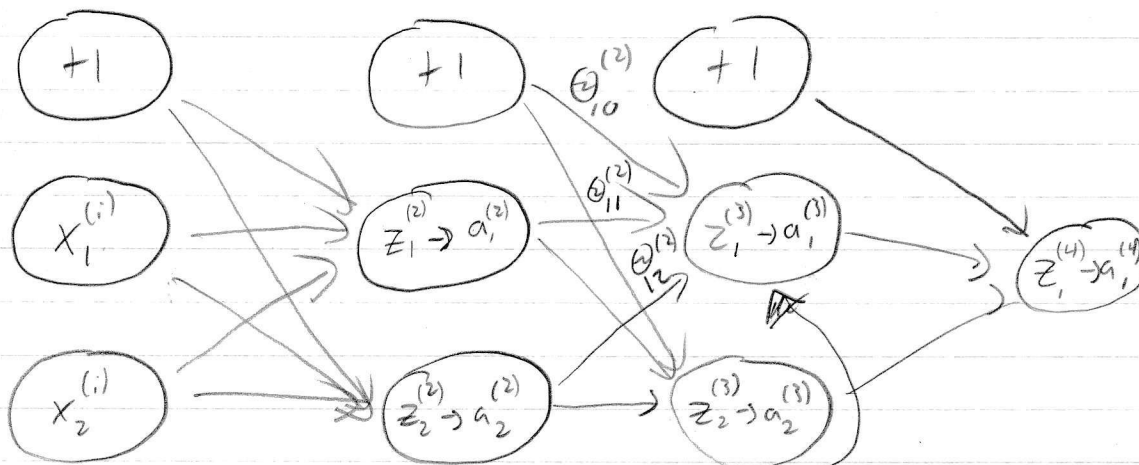
$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}$ if $j \neq 0$

$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$ if $j = 0$

$\frac{d}{d\theta_{ij}^{(l)}} J(\theta) = D_{ij}^{(l)}$

L3- Backpropagation Intuition

Recall Forward Propagation



$$z_1^{(3)} = (\theta_{10}^{(2)} \cdot 1) + (\theta_{11}^{(2)} \cdot a_1^{(2)}) + (\theta_{12}^{(2)} \cdot a_2^{(2)})$$

What is Backpropagation doing?

Recall:
$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ij}^{(l)})^2$$
 $(x^{(i)}, y^{(i)})$

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

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

(Think of $\text{cost}(i) \approx (h_{\theta}(x^{(i)}) - y^{(i)})^2$)

i.e. how well is the network doing on example i ?

watch rest of video for example

(3)

T2 - Backpropagation in Practice

L1 - Implementation Note: Unrolling Parameters

Advanced Optimization

function [jval, gradient] = costFunction(theta)
 $\hookrightarrow \mathbb{R}^{n+1}$ $\hookrightarrow \mathbb{R}^{n+1}$ (vectors)
 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

Ex. $S_1=10, S_2=10, S_3=1$
 $\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$
 $D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$
 $\begin{bmatrix} 0 & 0 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix} \approx 0 \rightarrow h_{\theta}(x)$

thetaVec = [Theta1(:); Theta2(:); Theta3(:)];

DVec = [D1(:); D2(:); D3(:)];

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

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

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

Learning Algorithm

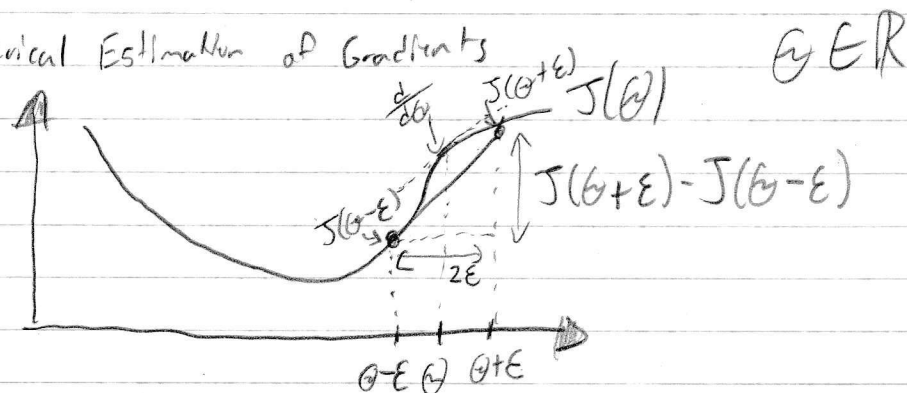
Have initial parameter $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.Unroll to get initialTheta to pass to fminunc(@costFunction, initialTheta, options)

function [jval, gradientVec] = costFunction(thetaVec)

→ From thetaVec, get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$. (reshape)→ Use forward Prop/back prop to compute $D^{(1)}, D^{(2)}, D^{(3)}$ and $J(\theta)$.→ Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradientVec.

L2-Gradient Checking

Numerical Estimation of Gradients



$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \epsilon = 10^{-4}$$

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

Parameter Vector θ

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

$$\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$$

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

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

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

Code: For $i = 1:n$,

thetaPlus = theta;

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

thetaMinus = theta;

thetaMinus(i) = thetaMinus(i) - EPSILON;

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

end;

Check that gradApprox \approx DVec

↑ From backprop

Implementation Note:

- Implement backprop 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. Using backprop for learning
↳ DVec

Important

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

L3-Random Initialization

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

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

So $\Theta_{ij}^{(2)} = 0$ for all i, j, l .

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

NO \rightarrow Also $d_1^{(2)} = d_2^{(2)} \rightarrow \frac{d}{d\Theta_{01}^{(1)}} J(\Theta) = \frac{d}{d\Theta_{02}^{(1)}} J(\Theta) \rightarrow \Theta_{01}^{(1)} = \Theta_{02}^{(1)}$

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

\Rightarrow

Random Initialization: Symmetry breaking

Different ϵ than previous slides

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

eg $\Theta_{11} = \text{rand}(10, 11) \cdot (2 \cdot \text{INIT_EPSILON}) - \text{INIT_EPSILON}$

$\Theta_{12} = \text{rand}(1, 11) \cdot (2 \cdot \text{INIT_EPSILON}) - \text{INIT_EPSILON}$

L4 - Putting it together

Training a neural network

Pick a network architecture (connectivity pattern between neurons):

- Number of input units: Dimension of features $x^{(i)}$

- Number of output units: Number of classes

- Reasonable default: 1 hidden layer, or if ≥ 1 hidden layer, have same number of hidden units in every layer (usually the more the better)

Then:

- 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 $\frac{\partial J(\Theta)}{\partial \Theta_{jk}^{(l)}}$ for $j=1:m$

Perform forward propagation and backpropagation using example $(x^{(i)}, y^{(i)})$
(Get activations $a^{(l)}$ and delta terms $d^{(l)}$ for $l=2, \dots, L$.)

$$\Delta^{(l)} := \Delta^{(l)} + d^{(l+1)} (a^{(l)})^T$$

end;

compute $\frac{\partial J(\Theta)}{\partial \Theta_{jk}^{(l)}}$

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