

ML w5 cost function and back propagation

- binary classification
- multi class classification

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

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

Θ : $J(\theta)$ and (partial) derivative term $\frac{\partial}{\partial \theta_{j|i}}$, every i, j, l

L : Total number of layer in network

S_l : number of unit in layer l .

K : number of output unit / class

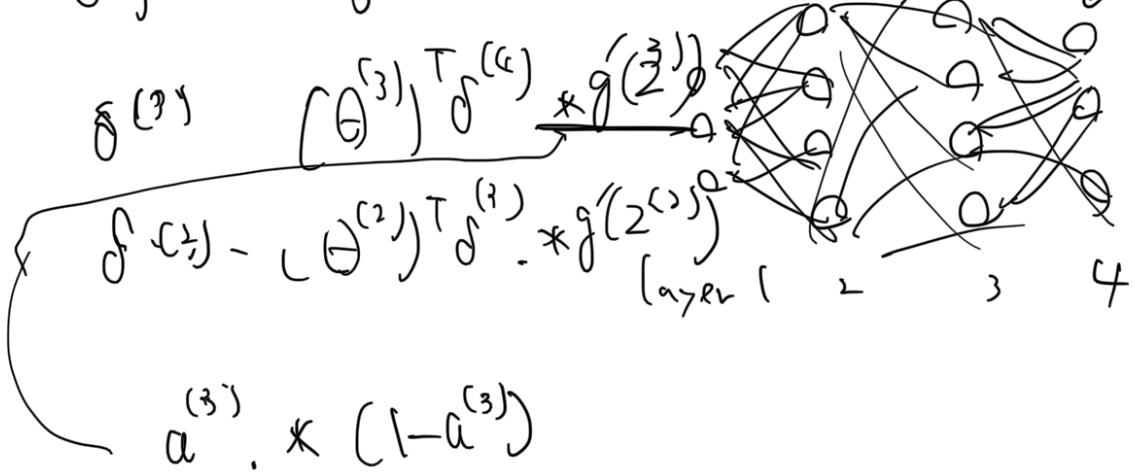
Back Propagation Algorithm

$$\text{minimize}_{\theta} J(\theta)$$

Gradient computation

Gradient computation (backward propagation algorithm)

$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad \text{or} \quad a_j^{(l)} \leftarrow \delta^{(3)} \cdot \partial^{(4)}$$



$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta^{(l+1)}$$

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

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

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

→ Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$

→ compute $f^{(L-1)}, f^{(L-2)}, \dots, f^{(2)}$

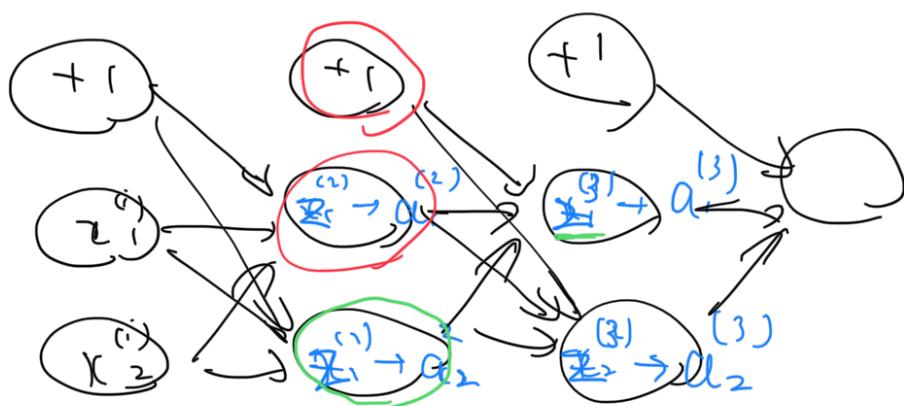
→ $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_r^{(l+1)}$

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

$\therefore m^{(l)} := \frac{1}{m} \Delta_{ii}^{(l)} \quad \text{if } j = 0$

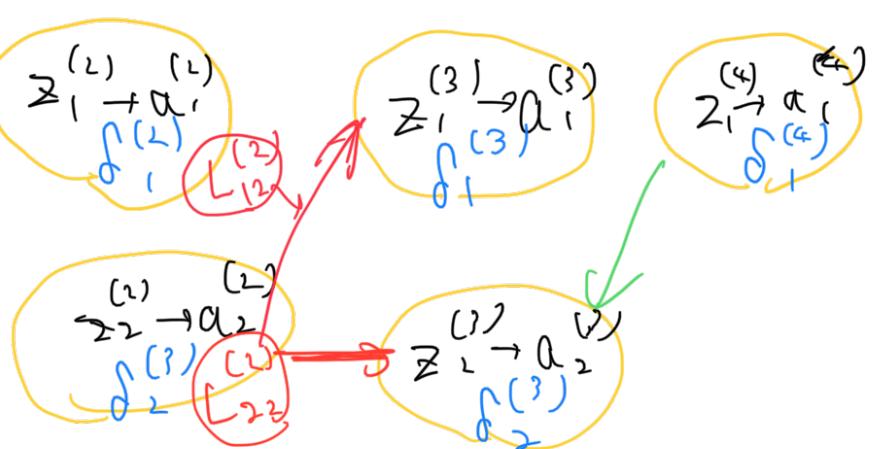
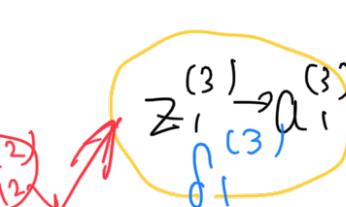
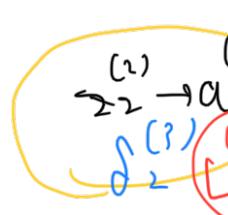
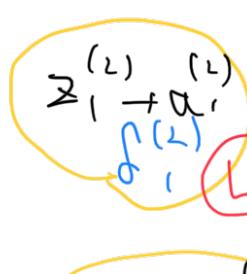
$\rightarrow V \cup m \cup$

Forward Propagation



$$(x^{(1)}, y^{(1)})$$

$$z_1^{(3)} = L_1^{(2)} \times 1 \quad L_1^{(2)} a_1 + L_{12}^{(2)} a_2$$



$$\delta_1^{(4)} = y^{(1)} - a_1^{(4)}$$

$$\delta_2^{(2)} = \textcircled{-}_2 \delta_1^{(1)} + \textcircled{L}_{22}^{(2)} \delta_2^{(3)}$$

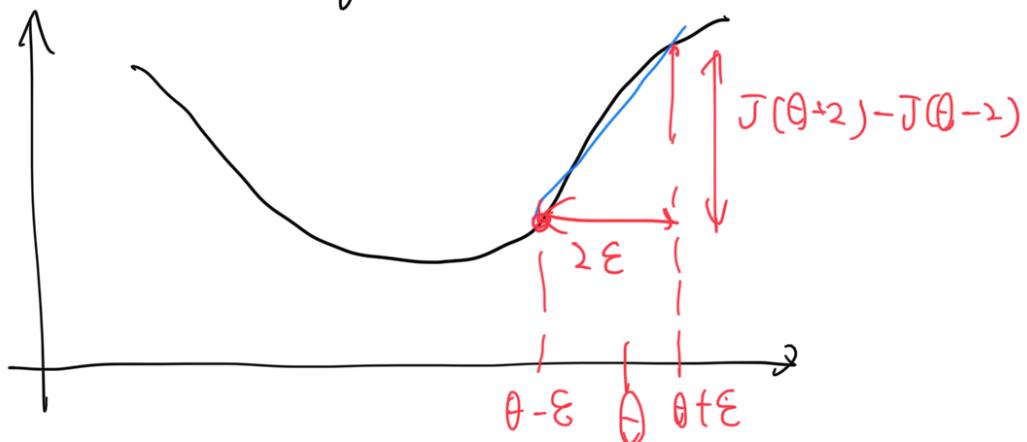


Back Propagation in Practice

(programming practice)

- reshape
- initial theta
- fminunc (

Gradient Checking .



$$\frac{\frac{d}{d\theta} J(\theta)}{\epsilon} = \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

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

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

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

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

$$\rightarrow \frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta, \theta_2 + \varepsilon, \theta_3, \dots, \theta_n) - J(\theta, \theta_2 - \varepsilon, \theta_3, \dots, \theta_n)}{2\varepsilon}$$



(activate code)

for $i = 1:n$

$\text{thetaPlus} = \text{theta}$

$\text{thetaPlus}(i) = \text{thetaPlus}(i) + \text{EPSILON}$

$\text{thetaMinus} = \text{theta}$

$\text{thetaMinus}(i) = \text{thetaMinus}(i) - \text{EPSILON}$

$\text{gradApprox}(i) = (J(\text{thetaPlus}) - J(\text{thetaMinus})) / (2 \times \text{EPSILON})$

end

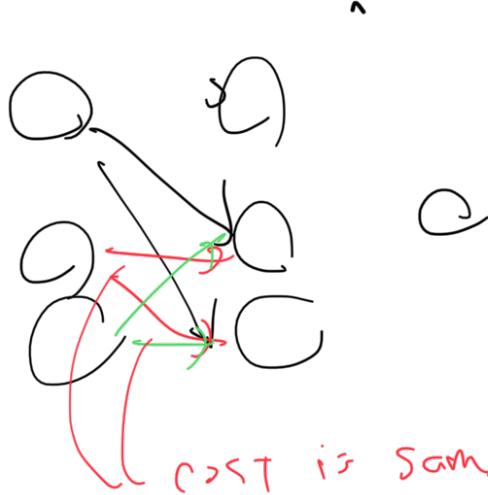
$((2 \times \text{EPSILON}))$

Check gradApprox DVec

Unrolled version of $(\theta^1, \theta^2, \theta^3)$

Backprop is computationally efficient

Random Initialization



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

$$\text{Theta1} = \text{rand}(10 \times 1) * (2 * \text{INIT_EPSILON})$$

$$\text{Theta2} = \underbrace{\text{rand}(1 \times 11) * (\text{INIT_EPSILON})}_{\text{Theta2}}$$

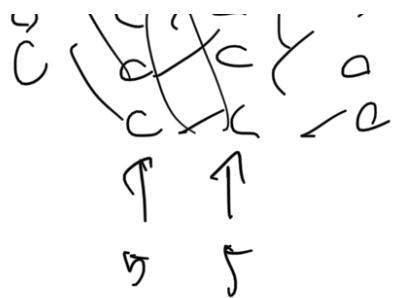
Putting it together

- Pick a nn architecture
(connectivity of input / hidden / output)

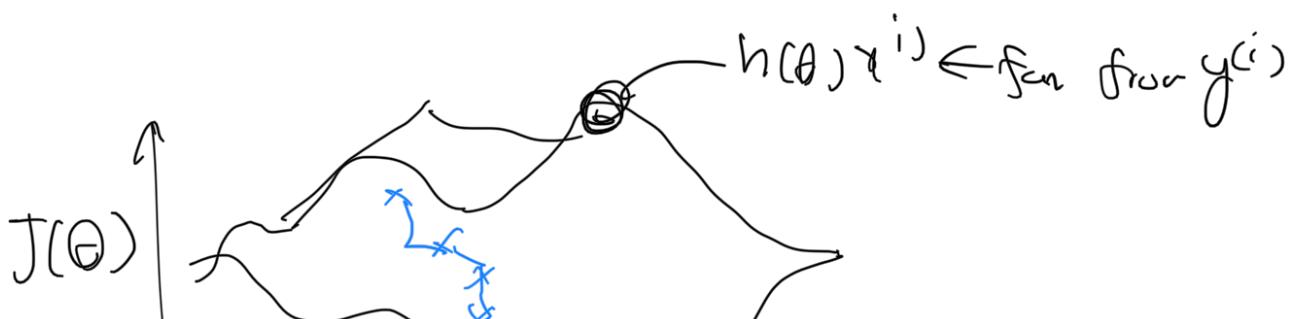
Reasonable default

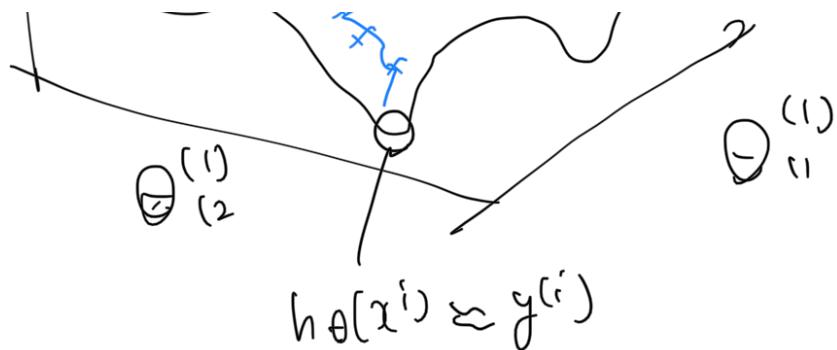
1: hidden layer, or if > 1 hidden layer, have same no. of hidden unit in every layer. (usually, the more the better)





- 1, Randomly initialize weights
- 2, Implement forward propagation to get $h_{\theta}(x)$
for any x)
- 3, Implementation code to compute cost function $J(\theta)$
- 4, Implement back prop to compute partial derivatives
 $\frac{\partial}{\partial \theta_j} J(\theta)$
- 5, Use gradient checking to compare $\frac{\partial}{\partial \theta_j} J(\theta)$
- 6, Use gradient descent or advanced optim. method
with back propagation to try minimize $J(\theta)$





...

...

1. You are training a three layer neural network and would like to use backpropagation to compute the gradient of the cost function. In the backpropagation algorithm, one of the steps is to update

$$\Delta_{ij}^{(2)} := \Delta_{ij}^{(2)} + \delta_i^{(3)} * (a^{(2)})_j$$

for every i, j . Which of the following is a correct vectorization of this step?

- $\Delta^{(2)} := \Delta^{(2)} + (a^{(3)})^T * \delta^{(3)}$
- $\Delta^{(2)} := \Delta^{(2)} + \delta^{(2)} * (a^{(2)})^T$
- $\Delta^{(2)} := \Delta^{(2)} + \delta^{(3)} * (a^{(3)})^T$
- $\Delta^{(2)} := \Delta^{(2)} + \delta^{(3)} * (a^{(2)})^T$

✓ 正解

This version is correct, as it takes the "outer product" of the two vectors $\delta^{(3)}$ and $a^{(2)}$ which is a matrix such that the (i, j) -th entry is $\delta_i^{(3)} * (a^{(2)})_j$ as desired.

2. Suppose Theta1 is a 5x3 matrix, and Theta2 is a 4x6 matrix. You set `thetaVec = [Theta1(:); Theta2(:)]`. Which of the following correctly recovers Theta2?

- `reshape(thetaVec(16 : 39), 4, 6)`
- `reshape(thetaVec(15 : 38), 4, 6)`
- `reshape(thetaVec(16 : 24), 4, 6)`
- `reshape(thetaVec(15 : 39), 4, 6)`
- `reshape(thetaVec(16 : 39), 6, 4)`

✓ 正解

This choice is correct, since Theta1 has 15 elements, so Theta2 begins at index 16 and ends at index $16 + 24 - 1 = 39$.

3. Let $J(\theta) = 2\theta^3 + 2$. Let $\theta = 1$, and $\epsilon = 0.01$. Use the formula $\frac{J(\theta+\epsilon)-J(\theta-\epsilon)}{2\epsilon}$ to numerically compute an approximation to the derivative at $\theta = 1$. What value do you get? (When $\theta = 1$, the true/exact derivative is $\frac{dJ(\theta)}{d\theta} = 6$.)

- 6.0002
 5.9998
 8
 6

 正解

We compute $\frac{(2(1.01)^3+2)-(2(0.99)^3+2)}{2(0.01)} = 6.0002$.

4. Which of the following statements are true? Check all that apply.

- Computing the gradient of the cost function in a neural network has the same efficiency when we use backpropagation or when we numerically compute it using the method of gradient checking.
 Using gradient checking can help verify if one's implementation of backpropagation is bug-free.

 正解

If the gradient computed by backpropagation is the same as one computed numerically with gradient checking, this is very strong evidence that you have a correct implementation of backpropagation.

- For computational efficiency, after we have performed gradient checking to
verify that our backpropagation code is correct, we usually disable gradient checking before using backpropagation to train the network.

 正解

Checking the gradient numerically is a debugging tool: it helps ensure a correct implementation, but it is too slow to use as a method for actually computing gradients.

- Gradient checking is useful if we are using one of the advanced optimization methods (such as in fminunc) as our optimization algorithm. However, it serves little purpose if we are using gradient descent.

5. Which of the following statements are true? Check all that apply.

- Suppose you are training a neural network using gradient descent. Depending on your random initialization, your algorithm may converge to different local optima (i.e., if you run the algorithm twice with different random initializations, gradient descent may converge to two different solutions).
- If we initialize all the parameters of a neural network to ones instead of zeros, this will suffice for the purpose of "symmetry breaking" because the parameters are no longer symmetrically equal to zero.

(✗) これを選択しないでください

The trouble with initializing the parameters to all zeros is not the specific value of zero but instead that every unit in the network will get the same update after backpropagation. Initializing the parameters to all ones has the same difficulty.

- If we are training a neural network using gradient descent, one reasonable "debugging" step to make sure it is working is to plot $J(\Theta)$ as a function of the number of iterations, and make sure it is decreasing (or at least non-increasing) after each iteration.
- Suppose you have a three layer network with parameters $\Theta^{(1)}$ (controlling the function mapping from the inputs to the hidden units) and $\Theta^{(2)}$ (controlling the mapping from the hidden units to the outputs). If we set all the elements of $\Theta^{(1)}$ to be 0, and all the elements of $\Theta^{(2)}$ to be 1, then this suffices for symmetry breaking, since the neurons are no longer all computing the same function of the input.

(✗) これを選択しないでください

Since the parameters are the same within layers, every unit in each layer will receive the same update during backpropagation. The result is that such an initialization does not break symmetry.