

# Neural Network: Learning

## CostFunction and Backpropagation

- denotation:

- $L$  = total number of layers in the network
- $s_l$  = number of units (not counting bias unit) in layer  $l$
- $K$  = number of output units/classes

- CostFunction:

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

- note

- at week 4 we realize that neural network works like linking many logistic regression. For a logistic regression we use  $J(\theta)$ , and here the double sum simply **adds up**  $J(\theta)$  calculated for each cell in **the output layer**. (that's why the loop  $k=1:K$ )

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

- the triple sum simply adds up the squares of all the individual  $\Theta$ s in the **entire** network.

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

- why do we add up all the theta in the entire network, but in the double sum there don't loop from  $i=1:L-1$ ?
  - remember, we are focusing on build up cost function which is used to judge the accuracy of prediction. Further, the cost function will also use to generate thetas.
  - so the double loop is to count all the predictions ( hypothesis ). But in regularization we would to generate all theta with better performance, so we should include computing all  $\Theta$ s.

## Backpropagation Algorithm

- goal:  $\min J(\theta)$  ⇨ base problem: Seeking partial derivative of all thetas.
- method: use Backpropagation
- algorithm here: (**unregularized**)

$$\frac{\partial}{\partial \Theta_{i,j}^{(l)}} J(\Theta)$$

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

## Backpropagation algorithm

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

Set  $\Delta_{ij}^{(l)} = 0$  (for all  $l, i, j$ ). (use to compute  $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$ )

For  $i = 1$  to  $m \leftarrow (\underline{x^{(i)}}, \underline{y^{(i)}})$ .

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

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

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

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

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

- o
- o crucial step:

4. Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$  using  $\delta^{(l)} = ((\Theta^{(l)})^T \delta^{(l+1)}) \cdot a^{(l)} \cdot (1 - a^{(l)})$

The delta values of layer  $l$  are calculated by multiplying the delta values in the next layer with the theta matrix of layer  $l$ . We then element-wise multiply that with a function called  $g'$ , or  $g$ -prime, which is the derivative of the activation function  $g$  evaluated with the input values given by  $z^{(l)}$ .

The  $g$ -prime derivative terms can also be written out as:

$$g'(z^{(l)}) = a^{(l)} \cdot (1 - a^{(l)})$$

5.  $\Delta_{i,j}^{(l)} := \Delta_{i,j}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$  or with vectorization,  $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

Hence we update our new  $\Delta$  matrix.

- $D_{i,j}^{(l)} := \frac{1}{m} (\Delta_{i,j}^{(l)} + \lambda \Theta_{i,j}^{(l)})$ , if  $j \neq 0$ .
- $D_{i,j}^{(l)} := \frac{1}{m} \Delta_{i,j}^{(l)}$  If  $j=0$

The capital-delta matrix  $D$  is used as an "accumulator" to add up our values as we go along and eventually compute our partial derivative. Thus we get  $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$

- Personal understanding:

- o compare to forward propagation, which compute each layer's units from left to right, backward propagation compute each layer's  $\delta$  from right to left.
- o **what  $\delta$  actually is?**

Intuitively,  $\delta_j^{(l)}$  is the "error" for  $a_j^{(l)}$  (unit  $j$  in layer  $l$ ).

1. in video, pro Andrew Ng explain it as

2. review our goal: find **partial derivative of theta** to minimize  $J(\Theta)$ . Therefore, the  $\delta$  formally define the derivative.

$$\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(t)$$

3. so our work still the same: find derivative, which equals to solve each  $\delta$ .

- o how to **solve  $\delta$** ?

- by calculus, we actually working on **chain rule**
- If we consider simple non-multiclass classification ( $k = 1$ ) and disregard regularization,
  - so here each  $\delta$ , we compute it by:

layer  $l$ .  $l+1$ . for each  $\delta_i^l$  in  $\delta^l$ .

we have.

$$\delta_i^l = \theta_{i1}^l \delta_1^{l+1} + \theta_{i2}^l \delta_2^{l+1} + \dots + \theta_{i, s_{l+1}}^l \delta_{s_{l+1}}^{l+1}$$

$$= [\theta_i^l]^T \delta^{l+1}$$

to vectorize:

compute  $\delta^l$

have.  $\delta^l = [\theta^l]^T \cdot \delta^{l+1}$

- more detail explanation: <https://www.zhihu.com/question/27239198/answer/89853077>

## Implement BP:

- skills:
  - unroll/reshape parameters
    - vector  $\leftrightarrow$  matrix
    - why?
      - matrix: better understanding when coding
      - vector: some advanced algorithms assume that you have all of your parameters unrolled into a big long vector
    - how? example here:

```
1 thetaVector = [ Theta1(:); Theta2(:); Theta3(:); ]
2 deltaVector = [ D1(:); D2(:); D3(:); ]
```

If the dimensions of Theta1 is 10x11, Theta2 is 10x11 and Theta3 is 1x11, then we can get back our original matrices from the "unrolled" versions as follows:

```
1 Theta1 = reshape(thetaVector(1:110),10,11)
2 Theta2 = reshape(thetaVector(111:220),10,11)
3 Theta3 = reshape(thetaVector(221:231),1,11)
4
```

### 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  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ .

Use forward prop/backprop to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\Theta)$ .

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

### Gradient Checking

- there may be some subtle bugs in code which hard to find, at the same time program still can give a theta which let  $j(\theta)$  decreased but actually wrong.
- to find the bug in advance, we should monitor our theta from backprop : **compare it with mathematically approximate theta**

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

$$\frac{\partial}{\partial \Theta_j} J(\Theta) \approx \frac{J(\Theta_1, \dots, \Theta_j + \epsilon, \dots, \Theta_n) - J(\Theta_1, \dots, \Theta_j - \epsilon, \dots, \Theta_n)}{2\epsilon}$$

- in vector or matrix :
- code example here ( let epsilon =  $10^{-4}$ )

```

1 epsilon = 1e-4;
2 for i = 1:n,
3     thetaPlus = theta;
4     thetaPlus(i) += epsilon;
5     thetaMinus = theta;
6     thetaMinus(i) -= epsilon;
7     gradApprox(i) = (J(thetaPlus) - J(thetaMinus))/(2*epsilon)
8 end;

```

- so before actually start to learning, use above **gradApprox** to verified backpropagation algorithm is correct.
- Once done the test, **turn off the testing loop**( which cost many computing resource)

#### • Random Initialize Theta:

- if assign initial theta with a same value, our neural network would degenerate to a logistic regression.
- so make it random in the range of epsilon

#### 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.

→ **Theta1** = `rand(10,11)*(2*INIT_EPSILON) - INIT_EPSILON;`

→ **Theta2** = `rand(1,11)*(2*INIT_EPSILON) - INIT_EPSILON;`

- how to choose epsilon?

<sup>2</sup>One effective strategy for choosing  $\epsilon_{init}$  is to base it on the number of units in the network. A good choice of  $\epsilon_{init}$  is  $\epsilon_{init} = \frac{\sqrt{6}}{\sqrt{L_{in} + L_{out}}}$ , where  $L_{in} = s_l$  and  $L_{out} = s_{l+1}$  are the number of units in the layers adjacent to  $\Theta^{(l)}$ .

- **but why ? \_wait to solve.**

## Assembling!

- before: pick a NN **Architecture**.
  - default: one hidden layer
  - if layer > 1, each hidden layer has the same amount of units.
  - input layer: dimension of features
  - output layer: number of classes
- Training:
  1. randomly initialize theta
  2. forward propagation
  3. cost function  $J(\theta)$
  4. **back propagation**
  5. gradient checking. Confirm BP works properly.
  6. Turn off gradient checking. Begin gradient descent or other function to minimize  $J(\theta)$

