

Neural Network: Learning

Termanteus

12 January 2020

Contents

I. Cost Function	3
II. Backpropagation	3
III. Octave: Unrolling Params	3
IV. Gradient Checking	4
V. Random initialization for Theta	4
VI. Training a Neural Network	4
1. Design a neural network	4
2. Training Neural Network	5

I. Cost Function

- Define:
 - L = total of layers in network
 - s_l = # units in layer l (exclude bias unit)
 - K = # output units/classes
- Cost function:

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

II. Backpropagation

- This is neural-network terminology for **minimizing cost function**.
- Goal: compute $\min_{\Theta} J(\Theta)$ \rightarrow compute partial derivative of $J(\Theta)$

Backpropagation algorithm

\rightarrow 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 $m \leftarrow (x^{(i)}, y^{(i)})$.

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

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

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

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

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ *$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$*

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

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

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

- With:

$$\delta^{(l)} = ((\Theta^{(l)})^T \delta^{(l+1)}) * g'(z^{(l)}); (g'(z^{(l)}) = a^{(l)} * (1 - a^{(l)}))$$

III. Octave: Unrolling Params

- Turn matrix to one vector: $A(:)$
- Combine multiple matrices to a matrix multiple column vectors: $[A(:), B(:), C(:)]$
- Turn one vector back to matrix: `reshape(A(from_element:to_element), rows, columns)`
- Apply to `fminunc(@costFunc, initTheta, options)`:
 - First put $\Theta_1, \Theta_2, \dots, \Theta_n$ to columns vector to put into `initTheta`
 - Inside `costFunc`, unroll to $\Theta_1, \Theta_2, \dots, \Theta_n$
 - Calculate $D^{(1)}, D^{(2)}, \dots$ and then unroll to `gradVec` to return.

IV. Gradient Checking

- Gradient checking will assure that the backpropagation implementation works as expected.

$$\frac{\delta}{\delta\Theta_j} \approx \frac{J(\Theta_1, \dots, \Theta_j + \varepsilon, \dots, \Theta_n)}{2\varepsilon}$$

Hence, we are only adding or subtracting epsilon to the Θ_j matrix. In octave we can do it as follows:

```
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;
9
```

Figure 1: Gradient Checking Algorithm

V. Random initialization for Theta

- Initializing all weights to zero doesn't work with neural network since all nodes will update to the same value repeatedly.

Hence, we initialize each $\Theta_{ij}^{(l)}$ to a random value between $[-\epsilon, \epsilon]$. Using the above formula guarantee bound. The same procedure applies to all the Θ 's. Below is some working code you could use to exp

```
1 If the dimensions of Theta1 is 10x11, Theta2 is 10x11 and Theta3 is 1x11.
2
3 Theta1 = rand(10,11) * (2 * INIT_EPSILON) - INIT_EPSILON;
4 Theta2 = rand(10,11) * (2 * INIT_EPSILON) - INIT_EPSILON;
5 Theta3 = rand(1,11) * (2 * INIT_EPSILON) - INIT_EPSILON;
6
```

Figure 2: Random Theta Algo

VI. Training a Neural Network

1. Design a neural network

- # input units = dimension of features
- # output units = # classes

- # hidden units per layer = the more the better (but the more the complex)
- **Default:** 1 hidden layer, if > 1 , hidden units/layer \geq the previous layer.

2. Training Neural Network

- Random initialize the weights
- Implement forward propagation to get $h_{\Theta}(x^{(i)})$ for any $x^{(i)}$
- Implement cost function
- Implement backpropagation to compute partial derivatives
- Use gradient checking to confirm the previous step works, then disable it..
- Use gradient descent or built-in optimization function to minimize the cost function with the weights in theta.
- **Notes:** Ideally, we want $h_{\Theta}(x^{(i)}) \approx y^{(i)}$. But $J(\Theta)$ is not convex so we might end up in a local minimum.