

Machine Learning PS4: Neural Net

Simon Seo

Some important functions in this implementation

- `g(x1, x2)` sigmoid function of dot product of x1 and x2
- `neuron(a_pre, w_intra, i)` calculate an activation from previous layer and weight matrix
- `forward_propagate(input_layer, * weight_matrix_args)` forward propagation algorithm from input and weight(s)
- `classifier_idx(input_layer, * weight_matrix_args)` classifies input data as {index of maximum value in output layer}
- `error(data, label, classifier, * weight_matrix_args)` calculates error (decimal) for given data, label, classifier, and weights
- `_cost_MLE(prediction, label)` maximum likelihood error: $-(y \cdot \log(h) + (1-y) \cdot \log(1-h))$
- `cost(output_matrix, label_list, cost_model)` calculates cost for given outputs, labels, and cost model

Error rate for the 5000 digits

The error calculated from the 5000 datasets was 2.48%

Loss function $J(\theta)$

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

Using the above MLE with natural log, the cost J calculated from the 5000 datasets was 0.28762951217843125.

Back Propagation algorithm for calculating partial derivatives

The pseudocode for computing the gradient for a given Θ is as follows:

```
computeGradient( $x, y, \theta$ ):  
 $\Delta_{st}^{(l)} = 0$  for all  $l, s, t$   
for each training data  $x^{(i)}$  and label  $y^{(i)}$  ( $i \in 1 \sim m$ ):  
     $a = \text{forwardPropagate}(x^{(i)}, \theta)$  //returns activations  $a^{(l)}$  for all  $l \in 1 \sim L$   
     $\delta = \text{backPropagate}(a, \theta, y^{(i)})$  //returns node error  $\delta^{(l)}$  for all  $l \in 2 \sim L$   
    for  $l \in 1 \sim L$ :  
        for  $t \in 1 \sim \text{dim}(a^{(l)}); s \in 1 \sim \text{dim}(a^{(l+1)})$ :  
             $\Delta_{st}^{(l)} += a_t^{(l)} \delta_s^{(l+1)}$   
            //each  $\Delta_{st}^{(l)}$  is the gradient of the cost function  $\left( \frac{\partial}{\partial \theta_{st}^{(l)}} J(\theta) \right)$   
return  $\Delta$ 
```

For each training data, the $\Delta_{st}^{(l)}$ term is the gradient for each $\theta_{st}^{(l)}$.

```
backPropagate( $a, \theta, y$ ):  
 $\delta_{st}^{(l)} = 0$  for all  $l, s, t$   
 $\delta^{(L)} = a^{(L)} - y$   
for  $l \in (L - 1) \sim 2$ :  
     $\delta^{(l)} = (\theta^{(l+1)T} \delta^{(l+1)}) .* (a^{(l)}) .* ([1] - a^{(l)})$   
    //  $\delta_j^{(l)}$  is the error contributed by node  $j$  of layer  $l$   
return  $\delta$ 
```