

Back Propagation in Netz

Nick Ewing

0.1 Descripton

The following definitions and algorithm characterize the implementation of `netz.core`.

0.2 Definitions

The following variables are used:

m : The number of training examples.

n : The number of input features for an example.

k : The number of outputs for an example.

$\{(x^i, y^i), \dots, (x^m, y^m)\}$: The training set, composed of m input-output example pairs where $x^{(i)}$ is an n dimensional vector containing the i^{th} input and $y^{(i)}$ a k dimensional vector containing the i^{th} expected output.

X : An $m \times n$ matrix where the i^{th} row contains $(x^{(i)})^T$.

Y : An $m \times k$ matrix where the i^{th} row contains $(y^{(i)})^T$.

L : The number of layers in the network.

s_l : The number of neurons in layer l . $s_1 = n$ and $s_L = k$.

λ : The regularization constant.

α : The learning rate constant.

γ : The learning momentum constant.

$\Theta_{ij}^{(l)}$: The synapse weight between neuron i of layer l and neuron j of layer $l + 1$. $\Theta^{(l)}$ is thus a $s_{(j+1)} \times s_j + 1$ matrix.

$A_{ij}^{(l)}$: The change in synapse weight between neuron i of layer l and neuron j of layer $l + 1$ of the last epoch. $A^{(l)}$ is thus a $s_{(j+1)} \times s_j + 1$ matrix.

$a^{(l)}$: A vector of length $s_l + 1$ when $l < L$ and s_l when $l = L$, containing the activation values for neurons in layer l where $a_0^{(l)}$ is the bias neuron when $l < L$.

$\delta^{(l)}$: A vector of length $s_l + 1$ when $l < L$ and s_l when $l = L$, containing the back propagated error values associated with neurons in layer l .

$\Delta_{ij}^{(l)}$: The sum change from all examples in synapse weight between neuron i of layer l and neuron j of layer $l + 1$. $\Delta^{(l)}$ is thus a $s_{(j+1)} \times s_j + 1$ matrix.

$D_{ij}^{(l)}$: The regularized mean change from all examples in synapse weight between neuron i of layer l and neuron j of layer $l + 1$. $D^{(l)}$ is thus a $s_{(j+1)} \times s_j + 1$ matrix.

$.*$: Element-wise matrix multiplication operator.

The sigmoid activation function:

$$g(z) = \frac{1}{1 - e^{-z}} \quad (1)$$

0.3 Algorithm

```

 $\Theta^{(j)} \leftarrow \text{RandomValuedMatrix}(s_{(j+1)}, s_j + 1)$  for  $j = 1 \dots L - 1$ 
 $A^{(j)} \leftarrow \text{ZeroValuedMatrix}(s_{(j+1)}, s_j + 1)$  for  $j = 1 \dots L - 1$ 
repeat
  for  $i \leftarrow 1 \dots m$  do
    {Forward propagate input activations through network}
     $a^{(1)} \leftarrow x^{(i)}$ 
     $a^{(l)} \leftarrow g(\Theta^{(l-1)} a^{(l-1)})$  for  $l = 2 \dots L$ 
    {Find error of output layer from training outputs}
     $\delta^{(L)} \leftarrow a^{(L)} - y^{(i)}$ 
    {Back propagate error through network}
     $\delta^{(l)} \leftarrow (\Theta^{(l)})^T \delta^{(l+1)} \cdot a^{(l)} \cdot (1 - a^{(l)})$  for  $l = L - 1 \dots 2$ 
    {Multiply errors by activations}
     $\Delta^{(l)} \leftarrow \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$  for  $l \leftarrow 1 \dots L - 1$ 
  end for
  {Divide errors by the number of training examples and add in regularization}
   $D_{ij}^{(l)} \leftarrow \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}$  if  $j \neq 0$ 
   $D_{ij}^{(l)} \leftarrow \frac{1}{m} \Delta_{ij}^{(l)}$  if  $j = 0$ 
  {Adjust weights}
  for  $j \leftarrow 1 \dots L - 1$  do
     $A^{(j)} \leftarrow \alpha D^{(j)} + \gamma A^{(j)}$ 
     $\Theta^{(j)} \leftarrow \Theta^{(j)} - A^{(j)}$ 
  end for
   $epoch \leftarrow epoch + 1$ 
   $MSE \leftarrow \text{CalculateMeanSquareError}(\Theta, X, Y)$ 
until  $MSE < \text{desired-error}$  or  $epoch > \text{max-epoch}$ 

```