

Explanatory Notes for 6.390

Shauntclair Ruiz (Current TA)

Fall 2022

Comments

A few important side notes on training. First, on derivatives:

Concept 1

Sometimes, depending on your **loss** and **activation** function, it may be easier to directly compute

$$\frac{\partial \mathcal{L}}{\partial \mathbf{Z}^L}$$

Than it is to find

$$\partial \mathcal{L} / \partial \mathbf{A}^L \text{ and } \partial \mathbf{A}^L / \partial \mathbf{Z}^L$$

So, our algorithm may change slightly.

Another thought: initialization.

Concept 2

We typically try to pick a **random initialization**. This does two things:

- Allows us to avoid weird **numerical** and **symmetry** issues that happen when we start with $\mathbf{W}_{ij} = 0$.
- We can hopefully find different **local minima** if we run our algorithm multiple times.
 - This is also helped by picking **random data points** in **SGD** (our typical algorithm).

Here, we choose our **initialization** from a **Gaussian** distribution, if you know what that is.

If you do not know a gaussian distribution, that shouldn't be a problem. It is also known as a "normal" distribution.

Pseudocode

Our training algorithm for backprop can follow smoothly from what we've laid out.

SGD-NEURAL-NET($\mathcal{D}_n, T, L, (m^1, \dots, m^L), (f^1, \dots, f^L), \text{Loss}$)

```

1  for every layer:
2      Randomly initialize
3          the weights in every layer
4          the biases in every layer
5
6  While termination condition not met:
7      Get random data point i
8      Kepp track of time t
9
10     Do forward pass
11         for every layer:
12             Use previous layer's output: get pre-activation
13             Use pre-activation: get new output, activation
14
15         Get loss: forward pass complete
16
17     Do back-propagation
18         for every layer in reversed order:
19             If final layer: #Loss function
20                 Get  $\partial \mathcal{L} / \partial A^L$ 
21
22             Else:
23                 Get  $\partial \mathcal{L} / \partial A^\ell$ : #Link two layers
24                      $(\partial Z^{\ell+1} / \partial A^\ell) * (\partial \mathcal{L} / \partial Z^{\ell+1})$ 
25
26                 Get  $\partial \mathcal{L} / \partial Z^\ell$ : #Within layer
27                      $(\partial A^\ell / \partial Z^\ell) * (\partial \mathcal{L} / \partial A^\ell)$ 
28
29             Compute weight gradients:
30                 Get  $\partial \mathcal{L} / \partial W^\ell$ : #Weights
31                      $\partial Z^\ell / \partial W^\ell = A^{\ell-1}$ 
32                      $(\partial Z^\ell / \partial W^\ell) * (\partial \mathcal{L} / \partial Z^\ell)$ 
33
34                 Get  $\partial \mathcal{L} / \partial W_0^\ell$ : #Biases
35                      $\partial \mathcal{L} / \partial W_0^\ell = (\partial \mathcal{L} / \partial Z^\ell)$ 
36
37             Follow Stochastic Gradient Descend (SGD): #Take step
38                 Update weights:
39                      $W^\ell = W^\ell - (\eta(t) * (\partial \mathcal{L} / \partial W^\ell))$ 
40
41                 Update biases:
42                      $W_0^\ell = W_0^\ell - (\eta(t) * (\partial \mathcal{L} / \partial W_0^\ell))$ 
43
44  Return final neural network with weights and biases

```

Last Updated: 11/09/22 05:37:01