

NLP assignment #1

Andrius Buinovskij - 18-940-270

Q1 b) i)

Well since all inputs are 1 and all weights are 1 then

$$\mathbf{s}_1^1 = \sum_{i=1}^3 \mathbf{x}_i \cdot \mathbf{w}_{i,1}^1 = 3 \quad (1)$$

$$\mathbf{s}_2^1 = \sum_{i=1}^3 \mathbf{x}_i \cdot \mathbf{w}_{i,2}^1 = 3 \quad (2)$$

Where \mathbf{s}_j^i is the sum input to the j 'th neuron in the i 'th layer, and \mathbf{s}^i is a vector whose length is equal to the number of neurons in the i 'th layer, with the input being the 0'th layer. So \mathbf{s}^1 is of length 2 since 1'st layer has 2 neurons.

Let \mathbf{n}_j^i be the output of the j 'th neuron in the i 'th layer, then of course \mathbf{n}^i is a vector of length equal to the number of neurons in the i 'th layer:

$$n_j^i = \text{ReLU}(\mathbf{s}_j^i) \quad (3)$$

So in our case We get

$$\mathbf{n}_1^1 = \text{ReLU}(\mathbf{s}_1^1) = \text{ReLU}(3) = 3 \quad (4)$$

$$\mathbf{n}_2^1 = \text{ReLU}(\mathbf{s}_2^1) = \text{ReLU}(3) = 3 \quad (5)$$

Same steps in the next layer:

$$\mathbf{s}_1^2 = 3 \cdot 1 + 3 \cdot 1 = 6 \quad (6)$$

And now We pass this through a sigmoid for our output instead of a ReLU so We get

$$\text{out} = \sigma(\mathbf{s}_1^2) = \frac{1}{1 + e^{-6}} = 0.99752737684 \quad (7)$$

Q1 b) ii)

$$\frac{d}{d\mathbf{w}_{j,k}^1} \mathbf{x}_i \cdot \mathbf{w}_{j,k}^1 = \mathbf{x}_i \quad (8)$$

Since the inputs are all 1 this simplifies to

$$\frac{d}{d\mathbf{w}_{j,k}^1} \mathbf{x}_i \cdot \mathbf{w}_{j,k}^1 = \frac{d}{d\mathbf{w}_{j,k}^1} 1 \cdot \mathbf{w}_{j,k}^1 = 1 \quad (9)$$

Now for the second layer:

$$\frac{d}{d\mathbf{w}_{1,1}^2} f = \frac{d}{d\mathbf{w}_{1,1}^2} \mathbf{n}_{1,1} \cdot \mathbf{w}_{1,1}^2 = \frac{d}{d\mathbf{w}_{1,1}^2} \text{ReLU}(\mathbf{s}_{1,1}) \cdot \mathbf{w}_{1,1}^2 = \text{ReLU}(\mathbf{s}_{1,1}) \quad (10)$$

Since $\text{ReLU}(\mathbf{s}_{1,1})$ is just a constant w.r.t. $\mathbf{w}_{1,1}^2$. In this case We simply get

$$\frac{d}{d\mathbf{w}_{1,1}^2} f = \text{ReLU}(\mathbf{s}_{1,1}) = 3 \quad (11)$$

And likewise for the other weight.

Q1 b) iii)

$$L_{BCE} = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})) \quad (12)$$

$$= -(0 \cdot \log(0.99752737684) + (1 - 0) \log(1 - 0.99752737684)) \quad (13)$$

$$= -(\log(0.00247262316)) \quad (14)$$

$$= 2.60684206722 \quad (15)$$

redone

Alright let's just roll all of these questions into one and do an iteration of backprop.

The forward pass is trivial.

Let \mathbf{x}_i^l be the output of the i 'th neuron in the l 'th layer.

Let $\mathbf{w}_{i,j}^l$ be the weight belonging to j 'th neuron multiplying the i 'th input in the l 'th layer. $\mathbf{w}_{:,j}^l$ is then simply the weight vector associated with the j 'th neuron in the l 'th layer.

Let \mathbf{s}_j^l be $(\mathbf{x}^{l-1})^\top \mathbf{w}_{:,j}^l$, the weighted sum of inputs to the j 'th neuron in the l 'th layer.

We can then say that $x_i^l = \sigma(\mathbf{s}_j^l)$, where σ is the nonlinearity of choice.

Then define

$$\delta_j^l = \frac{\partial L_{BCE}}{\partial \mathbf{s}_j^l} \quad (16)$$

And finally

$$\delta_j^{l-1} = \sum_{i=1}^{d^l} \frac{\partial L_{BCE}}{\partial \mathbf{s}_i^l} \frac{\partial \mathbf{s}_i^l}{\partial \mathbf{x}_j^{l-1}} \frac{\partial \mathbf{x}_j^{l-1}}{\partial \mathbf{s}_j^{l-1}} \quad (17)$$

$$= \sum_{i=1}^{d^l} \delta_i^l \frac{\partial \mathbf{s}_i^l}{\partial \mathbf{x}_j^{l-1}} \frac{\partial \mathbf{x}_j^{l-1}}{\partial \mathbf{s}_j^{l-1}} \quad (18)$$

Where d^l is the number of neurons in layer l . So now We have a recursive definition which uses dynamic programming. This could be further nuanced by expressing things in matrix notation, but it's good enough for present purposes.

Here is also a picture since just looking at symbols is a nightmare.

So now

$$\frac{\partial L_{BCE}}{\partial \mathbf{s}_1^2} = \frac{\partial L_{BCE}}{\partial \hat{y}} \frac{\partial \hat{y}}{\mathbf{s}_1^2} \quad (19)$$

$$\frac{\partial L_{BCE}}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} - (y \log(\hat{y}) + (1-y) \log(1-\hat{y})) \quad (20)$$

$$= -y \frac{1}{\hat{y}} - (1-y) \frac{\partial}{\partial \hat{y}} \log(1-\hat{y}) \quad (21)$$

$$= -y \frac{1}{\hat{y}} - (1-y) \frac{-1}{1-\hat{y}} \quad (22)$$

$$= -y \frac{1}{\hat{y}} + \frac{1-y}{1-\hat{y}} \quad (23)$$

$$\frac{\partial \hat{y}}{\mathbf{s}_1^2} = \frac{\partial}{\mathbf{s}_1^2} \sigma(\mathbf{s}_1^2) \quad (24)$$

$$= \sigma(\mathbf{s}_1^2) \cdot (1 - \sigma(\mathbf{s}_1^2)) \quad (25)$$

So then We have

$$\delta_1^2 = \frac{\partial L_{BCE}}{\partial \mathbf{s}_1^2} = \frac{\partial L_{BCE}}{\partial \hat{y}} \frac{\partial \hat{y}}{\mathbf{s}_1^2} \quad (26)$$

$$= \left(-y \frac{1}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \left(\sigma(\mathbf{s}_1^2) \cdot (1 - \sigma(\mathbf{s}_1^2)) \right) \quad (27)$$

Alright. We only have one more of these to figure out:

$$\delta_j^1 = \sum_{i=1}^{d^2} \delta_i^2 \frac{\partial \mathbf{s}_i^2}{\partial \mathbf{x}_j^1} \frac{\partial \mathbf{x}_j^1}{\partial \mathbf{s}_j^1} \quad (28)$$

But since $d^2 = 1$, i.e. there is only one neuron in the output layer, We have

$$\delta_j^1 = \sum_{i=1}^{d^2} \delta_i^2 \frac{\partial \mathbf{s}_i^2}{\partial \mathbf{x}_j^1} \frac{\partial \mathbf{x}_j^1}{\partial \mathbf{s}_j^1} \quad (29)$$

$$= \delta_1^2 \frac{\partial \mathbf{s}_1^2}{\partial \mathbf{x}_j^1} \frac{\partial \mathbf{x}_j^1}{\partial \mathbf{s}_j^1} \quad (30)$$

$$= \delta_1^2 \mathbf{w}_j^2 \sigma(\mathbf{s}_j^1) (1 - \sigma(\mathbf{s}_j^1)) \quad (31)$$

Now for the gradient update We will of course need

$$\frac{\partial L_{BCE}}{\partial \mathbf{w}_{i,j}^l} \quad (32)$$

But this can be easily derived since

$$\frac{\partial L_{BCE}}{\partial \mathbf{w}_{i,j}^l} = \frac{\partial L_{BCE}}{\partial \mathbf{s}_j^l} \frac{\partial \mathbf{s}_j^l}{\partial \mathbf{w}_{i,j}^l} \quad (33)$$

$$= \delta_j^l \mathbf{x}_i^{l-1} \quad (34)$$

Now We do a forward pass and We know that $\mathbf{s}_i^1 = 3$ and $\mathbf{s}_1^2 = 6$. Plugging in values We then get

$$\delta_1^2 = \left(-y \frac{1}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \left(\sigma(\mathbf{s}_1^2) \cdot (1 - \sigma(\mathbf{s}_1^2)) \right) \quad (35)$$

$$= \left(\frac{1}{1 - 0.99752737684} \right) \left(\sigma(6) \cdot (1 - \sigma(6)) \right) \quad (36)$$

$$= \left(\frac{1}{1 - 0.99752737684} \right) \left(0.99752737684 \cdot (1 - 0.99752737684) \right) \quad (37)$$

$$= 404.428792942 \cdot 0.00246650929 \quad (38)$$

$$= 0.99752737493 \quad (39)$$

Similarly We have

$$\delta_j^1 = \delta_1^2 \mathbf{w}_j^2 \sigma(\mathbf{s}_j^1) (1 - \sigma(\mathbf{s}_j^1)) \quad (40)$$

$$\delta_1^1 = \delta_1^2 \mathbf{w}_1^2 \sigma(\mathbf{s}_1^1) (1 - \sigma(\mathbf{s}_1^1)) \quad (41)$$

$$= 0.997527374931 \cdot \sigma(3) (1 - \sigma(3)) \quad (42)$$

$$= 0.997527374931 \cdot 0.95257412682 \cdot 0.04742587317 \quad (43)$$

$$= 0.04506495478 \quad (44)$$

δ_2^1 is equal to δ_1^1 .

Now that We have the deltas We can use

$$\frac{\partial L_{BCE}}{\partial \mathbf{w}_{i,j}^l} = \delta_j^l \mathbf{x}_i^{l-1} \quad (45)$$

For weights in the first layer this means

$$\frac{\partial L_{BCE}}{\partial \mathbf{w}_{i,j}^1} = \delta_j^1 \mathbf{x}_i^0 \quad (46)$$

$$= 0.04506495478 \cdot 1 \quad (47)$$

$$= 0.04506495478 \quad (48)$$

Where the 0'th layer is the input layer.

$$\frac{\partial L_{BCE}}{\partial \mathbf{w}_{i,j}^2} = \delta_j^2 \mathbf{x}_i^1 \quad (49)$$

$$= 0.99752737493 \cdot 3 \quad (50)$$

$$= 2.99258212479 \quad (51)$$

All that is left is to perform a step for each weight. All weights in the first layer simplify to

$$w_{i,j}^1 = 1 - 0.1 \cdot 0.04506495478 \quad (52)$$

$$= 0.99549350452 \quad (53)$$

Sidenote

So when actually implementing this, for the next assignemnt, We get

$$\delta_1^2 = \frac{\partial L_{BCE}}{\partial \mathbf{s}_1^2} = \frac{\partial L_{BCE}}{\partial \hat{y}} \frac{\partial \hat{y}}{\mathbf{s}_1^2} \quad (54)$$

$$= \left(-y \frac{1}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \left(\sigma(\mathbf{s}_1^2) \cdot (1 - \sigma(\mathbf{s}_1^2)) \right) \quad (55)$$

$$= \left(-y \frac{1}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \left(\hat{y} \cdot (1 - \hat{y}) \right) \quad (56)$$

$$= -\hat{y} \cdot (1 - \hat{y}) \cdot \frac{y}{\hat{y}} + \hat{y} \cdot (1 - \hat{y}) \cdot \frac{1-y}{1-\hat{y}} \quad (57)$$

$$= -(1 - \hat{y}) \cdot y + \hat{y} \cdot (1 - y) \quad (58)$$