

Practical assignment 1 feedback

You

July 17, 2018

Abstract

Your abstract.

1 Possible Changes

1. lecture 3 slide 36

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t (\text{grad} L + l \theta_l)$$

i think it should be

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t (\text{grad} L + l \theta^{(t)})$$

2. lecture 3 slide 59 typo save in the middle
3. paper assignment 1 pdf page 2

$$\delta_{N-1} = \frac{\partial \mathcal{L}}{\partial s_{N-1}} = \delta_N \cdot W_N \cdot \frac{\partial f(s_1)}{\partial s_1} \quad (1)$$

shouldn't be more precise to be:

$$\delta_{N-1} = \frac{\partial \mathcal{L}}{\partial s_{N-1}} = \delta_N \cdot W_N \cdot \frac{\partial f_{N-1}(s_1)}{\partial s_1} \quad (2)$$

and all the others similarly.

2 Feedback

In general i lost a lot of time trying to figure out what happens with the Loss functions. I think the slides do not cover the topic well enough. There is this suggested site <http://ufdl.stanford.edu/wiki/index.php/Softmax> in the practical-1.ipynb but I found these 2 slides very very helpful and way more clarifying than the recommended site. They can be found in <http://www.psi.toronto.edu/~jimmy/ece521/Lec9-nn2.pdf> (slides 4-5)

Binary cross-entropy loss

- Consider the Q distribution to be the discrete probability distribution of the observed binary labels $t \in \{0, 1\}$ in the training dataset
 - $Q(t = 1|x) = 1$ is either zero or one depending on the training example. (The dataset is observed, so there is no randomness)
- Choose the P distribution to be the model's prediction: $P(t = 1|x) = \hat{p}(x)$. The KL divergence between Q and P is in fact the cross-entropy loss:
$$KL(Q||P) = \sum_t Q(t|x) \log \frac{Q(t|x)}{P(t|x)} = - \sum_t Q(t|x) \log P(t|x) + \underbrace{\sum_t Q(t|x) \log Q(t|x)}_{\text{this is zero}} \\ = -Q(t = 1|x) \log P(t = 1|x) - Q(t = 0|x) \log P(t = 0|x) \\ = -(\hat{p} \log \hat{p} + (1 - \hat{p}) \log(1 - \hat{p}))$$
 - The cross-entropy loss function measures the distance between the empirical data distribution and the model predictive distribution

Multi-class cross-entropy loss and softmax

- It is easy to use the KL divergence interpretation to generalize the cross-entropy loss to a multi-class scenario:
 - Let there be K classes, with class labels $t \in \{1, \dots, K\}$
 - The multi-class cross entropy loss can be written using indicator function $I(\cdot)$:
$$KL(Q||P) = \sum_t Q(t|x) \log \frac{Q(t|x)}{P(t|x)} = - \sum_{k=1}^K I(k, t_{obs}) \log P(t = k|x)$$
 - Similarly, the multi-class generalization of the sigmoid function is the softmax function. The multi-class predictive distribution becomes:
$$P(t = k|x) = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}}$$
Here z_k are the outputs of a neural network or linear regression

Figure 1:

In practical 1 i got really confused and lost quite some time to figure out what SoftMaxLoss is. I could not find this term on the internet. I ended up implementing a Softmax layer and a Multi class cross entropy loss, which i think is the expected answer. I do not understand why it is called SoftMaxLoss instead of Softmax Layer and Cross entropy Loss function. There should be at least some clarifying comments, in my opinion. My guess is that in the practical session some explanation was given.

3 Conclusion

Overall, i think the slides are very good and i have learned a lot from the assignment 1. I self studied with the given material and whatever i could find on the internet. I had no collaboration with any other student or guy that could help me and of course i did not attend any practical or theoretical lecture. Under these conditions it took me 3 weeks (i guess avg of 3-4 hours a day) to solve the whole practical 1 (programming and pen and paper), hopefully the time would reduce dramatically if i was attending the practical sessions or if i had collaboration with other students.