

## Problem Set 12, Dec 5, 2019 (Neural Networks)

**Goals.** The goal of this exercise is to

- Better understand neural networks
- Implement the feed-forward function and backpropagation in a simple neural net.

**Setup, data and sample code.** Obtain the folder `labs/ex12` of the course github repository

[github.com/epfl/ML\\_course](https://github.com/epfl/ML_course)

In the following problems, we will use a very simple neural network. Let's assume we have a three-layer neural net with one input layer of size  $D = 4$ ,  $L = 1$  hidden layers of size  $K = 5$ , and one output layer of size 1, as shown in Figure 1.

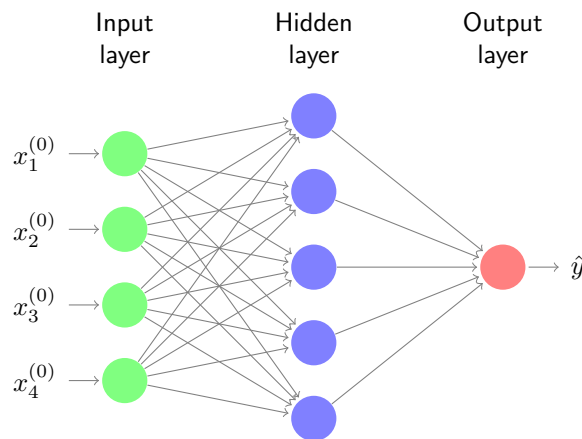


Figure 1: A simple neural network.

### Problem 1 (Feed-forward in neural networks):

In our simplified neural network, we have the feed-forward function shown below:

$$x_j^{(1)} = \phi \left( z_j^{(1)} \right) = \phi \left( \sum_{i=1}^D w_{i,j}^{(1)} x_i^{(0)} + b_j^{(1)} \right), \quad (1)$$

$$\hat{y} = \phi \left( z_1^{(2)} \right) = \phi \left( \sum_{i=1}^K w_{i,1}^{(2)} x_i^{(1)} + b_1^{(2)} \right). \quad (2)$$

Use Equation 1 and Equation 2 to fill in the corresponding template function in the notebook, and pass the test. For simplicity, in the following questions, let the bias term be 0 and use the Sigmoid as the activation function  $\phi(\cdot)$ .

### Problem 2 (Backpropagation in neural network):

Assume that we use the squared error as our loss function, as shown in Equation 3:

$$\mathcal{L} = \frac{1}{2}(\hat{y} - y)^2, \quad (3)$$

where we have only one sample in our case, and  $y$  is the true value while  $\hat{y}$  is the network prediction.

Evaluate the derivative of  $\mathcal{L}(\mathbf{w})$  with respect to weights  $w_{i,1}^{(2)}$  and  $w_{i,j}^{(1)}$ , and implement the corresponding function in the notebook.

### Problem 3 (Effect of regularization):

What is the effect of regularization on the weights? To get some insight, let  $\Theta$  be the vector of all weights in the neural network. Recall that we do not penalize the bias terms. Therefore, let us ignore them in the following. Let  $\Theta^*$  be a parameter that minimizes the cost function  $\mathcal{L}$  for the given test set (where the cost function does not include the regularization). We would like to study how the optimal weight changes if we include some regularization.

In order to make the problem tractable, assume that  $\mathcal{L}(\Theta)$  can be locally expanded around the optimal parameter  $\Theta^*$  in the form

$$\mathcal{L}(\Theta) = \mathcal{L}(\Theta^*) + \frac{1}{2}(\Theta - \Theta^*)^\top \mathbf{H}(\Theta - \Theta^*),$$

where  $\mathbf{H}$  is the Hessian whose components are the entries

$$\frac{\partial^2 \mathcal{L}}{\partial \Theta_i \partial \Theta_j}.$$

Now add a regularization term of the form  $\frac{1}{2}\mu \|\Theta\|_2^2$ .

1. Show that the optimum weight vector for the regularized problem is given by

$$\mathbf{Q}(\mathbf{\Lambda} + \mu \mathbf{I})^{-1} \mathbf{\Lambda} \mathbf{Q}^\top \Theta^*,$$

where  $\mathbf{H} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$  is the SVD of the symmetric matrix  $\mathbf{H}$ ,  $\mathbf{Q}$  is an orthonormal matrix, and  $\mathbf{\Lambda}$  is a diagonal matrix whose entries are non-negative and decreasing along the diagonal.

2. Show that  $(\mathbf{\Lambda} + \mu \mathbf{I})^{-1} \mathbf{\Lambda}$  is again a diagonal matrix whose  $i$ -th entry is now  $\lambda_i / (\lambda_i + \mu)$ .
3. Argue that along the dimensions of the eigenvectors of  $\mathbf{H}$  that correspond to large eigenvalues  $\lambda_i$  essentially no changes occur in the weight, but that along the dimensions of eigenvectors of very small eigenvalues the weight is drastically decreased.