

Homework 3

Ana Santos, 84364 Inês Branco, 81506 Vânia Nunes, 85235

April 2019

The link to the code: <https://drive.google.com/file/d/1vqvNrZs5v-PNK9hwzZ-2ShSvCO08gKSh/view?usp=sharing>

1 Problem 1

In this report we have considered the following notation:

- a) linear unit with quadratic error as loss function
- b) sigmoid unit with quadratic error as loss function
- c) sigmoid unit with cross entropy as loss function

As required in the problem, we used the gradient descent algorithm to train our data. The input corresponds to the 4 features of the iris flower: Sepal Length, Sepal Width, Petal Length and Petal Width. The purpose of the Irish-flower data set training is to be able to classify a flower in the category Iris-Virginica (output=1) or non Iris-Virginica (output=0). In total, there are 150 input vectors. There were 3 cases of combinations of activation and loss functions. First, the net for each case was computed:

$$net = \sum_{j=1}^N w_j x_{kj} \quad (1)$$

For both cases a) and b), it was computed a quadratic error as loss function:

$$E = \frac{1}{2} \sum_{k=1}^N (t_k - o_k)^2, \quad o_k = \sigma \left(\sum w_i x_i \right) \quad (2)$$

For c), a cross entropy error was used:

$$E = - \sum (t \ln o + (1 - t) \ln (1 - o)), \quad o = \sigma \left(\sum w_i x_i \right) \quad (3)$$

A loop was then performed in a number n of epochs. For each epoch, another loop over each point of the data set used for training and test was executed. The expected output for each point was computed, o_k , using the corresponding activation function for each case:

$$\sigma(x) = 1, \quad \text{for } a) \quad (4)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}, \quad \text{for } b) \text{ and } c) \quad (5)$$

To update the weights, we used the following equations:

$$w^{new} = w^{old} + \Delta w, \quad \Delta w = -\eta \frac{\partial E}{\partial w_i} \quad (6)$$

where $\frac{\partial E}{\partial w_i}$ takes the form:

$$\frac{\partial E}{\partial w_i} = - \sum (t - \sigma(net)) \sigma(net) (1 - \sigma(net)) x_i, \quad \text{for } a) \text{ and } b) \quad (7)$$

$$\frac{\partial E}{\partial w_i} = - \sum \left(\frac{t}{\sigma(net)} - \frac{1-t}{1-\sigma(net)} \right) \sigma(net) (1 - \sigma(net)) x_i \quad \text{for } c) \quad (8)$$

In each loop over the data set, it is found a value of Δw to update the weights. Over the various epochs, various values of Δw are computed in order to update the vector w_k for each point.

In order to study the convergence rate, we plotted the variation of the error in function of initialization parameters such as the learning parameter, η , the initialization vector of weights, w , and the number of epochs. The results are shown hereafter.

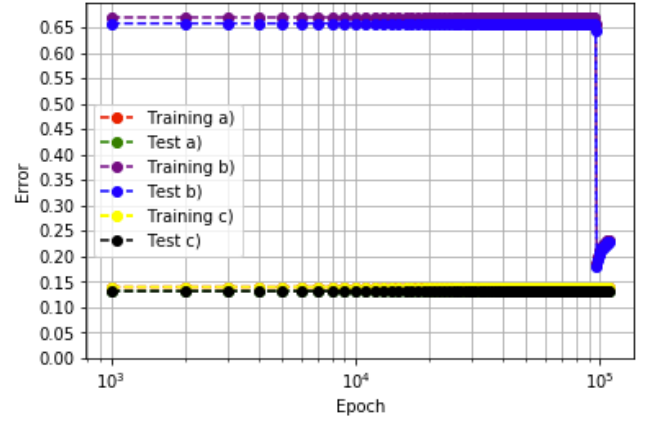


Figure 1: Plot of the variation of the average error (over all points of the data set) with the number of epochs.

Observing the plot, we realize that for both cases a) and c) the error has already converged in the epoch number 1000, whereas for case b) only between the epoch number 9×10^4 and 10^5 the error achieves a minimum. In order to find out which method has the best convergence rate, we have repeated the plot in figure 1 for the cases a) and c), with a lower range of the value of epoch (see the code). They both converge approximately at an epoch of 100. We conclude therefore that the convergence rate for a sigmoid unit with quadratic error is much slower than for the other two cases.

The learning rate, η , is a hyper-parameter that controls how much we are adjusting the weights of our network with respect to the loss gradient. The lower the value of η , the slower we travel along the downward slope. If the learning parameter is too high, the gradient descent can overshoot the minimum and fail to converge. Looking at figure 2, it is observable that for the cases a) and c) the error appears to remain almost constant for the values of η presented, whereas for b) a minimum is achieved for a value of η close but bigger than 0.01. From this point on, the error continues to increase until it levels to values close to 0.30.

Once again, to find a minimal value for the learning parameter, we re-plotted figure 2, for only the cases a) and c) (see code). With a better resolution, we observe that, for c), the error is minimum in 0.009, whereas for a) it is not observable.

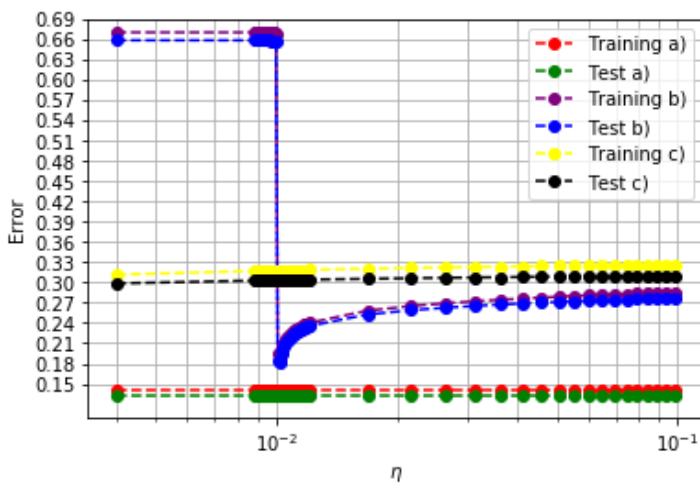


Figure 2: Plot of the variation of the average error (over all points of the data set) with the learning rate, η . The results were run for 10^5 epochs.

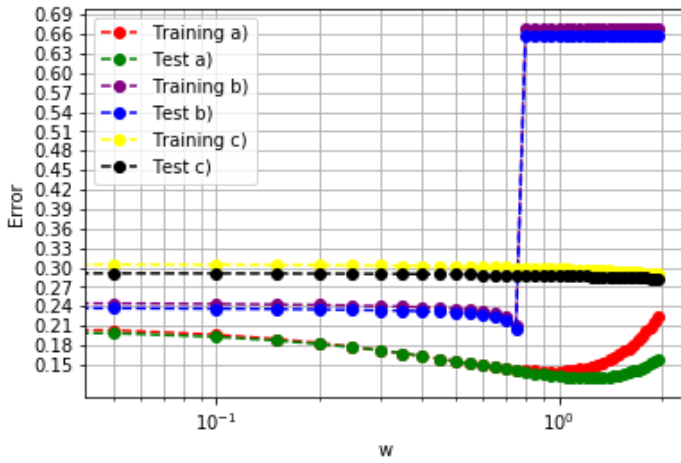


Figure 3: Plot of the variation of the average error (over all points of the data set) with epoch

For testing the impact of the initial weight vector, w , we studied the error variation with this parameter. The weight vector was initialized with all its elements with the same values. In figure 3, we see that in the case b), there is a convergence of the error for an initial weight of 7.5×10^{-1} ; in the case a), the minimum occurs in 1; and in case c), the error is minimum in 0 (see extra plots in the end of the code). For high values of epochs used, for the cases a) and c) the final values of the weight is independent of initial values of the weight and learning parameter. We've therefore computed, this final weights using less number of epochs (see plots in the end of the code).

Having in mind the optimal parameters already found, it was computed for each method the correspondent test and training errors:

- a) final error (training) = 0.140
- a) final error (testing) = 0.132
- b) final error (training) = 0.266
- b) final error (testing) = 0.259
- c) final error (training) = 0.262
- c) final error (testing) = 0.253

The best method revealed to be the one for case a).

In addition, it must be referred that the observed behaviour of training and test data was similar, which indicates there was not over-fitting of the results, and the trained algorithm was well applied to a generalized test data.

Also, about the plots, since one third of the output values are 1 and two thirds are 0, when the error function of b) is 0.66 it means that the output of the trained neuron is always 1, so the weights are all positive. The same argument can be made for 0.33, the output is always 0, the weights are all negative

2 Problem 2

To process the data, we decided to convert all the input parameters into numbers from 0 to 1. We also saw that some of the data was missing, in total 9 data points out of 277 had some information missing, so we decided to delete them. We divided the data set into testing and training, where about 75% of the data went into training and 25% into testing.

For all the graphs plotted we used one hidden layer, with a linear activation function and 2 inputs (similar to the class). We used 32 batch points and either 500 or 1000 epochs depending on how long it took to have smaller losses.

We tested this neural networks for the three activation functions (sigmoid, linear and relu), for the two regularizations (with the weights squared or with their absolute value), and for four different values of the regularization parameter λ (0.001, 0.01, 0.1, 1).

We made sure that all other initial parameters such as the initial weights were the same for all the plots.

By looking at the plots (in the code) we can see that using the mean square loss function, the neural networks that use the leaky relu converge the faster. For higher values of the regularization parameter λ the convergence seems to be faster for all functions (either activation, loss or regularization).

We can also see, by looking at the graphs (in the code), that when using the cross entropy loss function the convergence seems to be slower for all activation functions, and seems to be better for the sigmoid activation function (if the regularization used is $L2$, with the weights squared).

The values for the testing and training error do not diverge greatly so we assume that there was not over-fitting.

Since the mean-square error works better for linear and relu functions and the cross entropy for sigmoid function, the results are what was expected.