

# Optimization Methods for Machine Learning

## Homework 1

---

Ivan Bergonzani, Michele Cipriano

November 28, 2017

## 1 INTRODUCTION

The aim of the homework is to train and compare different neural networks on a regression problem. In particular, the regression task is performed on the Franke's function, building a dataset from it by sampling 100 random points  $(x^i, y^i)$  with noise, i.e. with  $y^i = F(x^i) + \varepsilon^i$  where  $\varepsilon^i$  is a random number in  $[-10^{-1}, 10^{-1}]$  and  $F$  is the Franke's function. The dataset has been split into a training set (70% of the dataset) and a test set (the remaining 30%).

Two architectures have been compared using different training methods and different hyperparameters. In particular, the multi-layer perceptron (MLP) and the radial basis function network (RBFN) have been trained with full minimization, two blocks method and decomposition method. In the full minimization, the training error is minimized w.r.t all the weights using the gradient descent algorithm. In the two blocks method, for the MLP the error is minimized via an extreme learning procedure and for the RBFN the error is minimized by first selecting the centers through a clustering algorithm and then by solving a linear least squares problem. In the decomposition method the error is minimized by alternating a convex minimization w.r.t. the output weights and a non-convex minimization w.r.t. all the other weights.

The best result has been obtained with **TODO**, which has an error of **TODO** on the test set. The project has been developed in Python with TensorFlow and Numpy for the learning algorithms and the computation of the tensors.

## 2 FULL MINIMIZATION

As mentioned before, the full minimization is done by using the gradient descent algorithm on the training error function, defined as:

$$E(\omega; \pi) = \frac{1}{2P} \sum_{p=1}^P \|f(x^p) - y^p\|^2 + \rho \|\omega\|^2$$

where  $\omega = (v, w, b)$  are the weights,  $\pi = (N, \rho, \sigma)$  are the hyperparameters,  $P$  is the dimension of the training set and  $f$  is the function computed by the neural network.

Let's consider a shallow MLP with linear output:

$$f(x) = \sum_{j=1}^N v_j g \left( \sum_{i=1}^n w_{ji} x_i - b_j \right)$$

with  $n$  dimension of the input,  $N$  dimension of hidden layer,  $v \in \mathbb{R}^N$  weights from the hidden to the output layer and  $g$  activation function, defined as:

$$g(t) = \tanh \left( \frac{t}{2} \right) = \frac{1 - e^{-\sigma t}}{1 + e^{-\sigma t}}$$

The training has been performed over 15000 iterations on the training set with a learning rate  $\eta = 0.001$ . The hyperparameters have been tuned with a grid search over all the possible combinations of the hyperparameters in:

$$(N, \rho, \sigma) \in \{25, 50, 75, 100\} \times \{10^{-3}, 10^{-4}, 10^{-5}\} \times \{1, 2, 3, 4\}$$

selecting the network with the lowest error on the test set. The best result obtained an error of 0.009 on the training set and an error of 0.006 on the test set as it is also possible to see in table 5.1. The hyperparameters found are  $N = 75$ ,  $\rho = 10^{-5}$  and  $\sigma = 3$ . In figure 5.1 the plot of the function computed by the MLP found.

TODO: explain overfitting/underfitting

An analogous experiment has been done with the radial basis function network, defined as:

$$f(x) = \sum_{j=1}^N v_j \phi(\|x^i - c_j\|)$$

where  $c_j \in \mathbb{R}^2$  is the center of the  $j$ -th hidden neuron,  $v \in \mathbb{R}^N$  weights from the hidden to the output layer and  $\phi$  is the activation function, defined as:

$$\phi(\|x - c_j\|) = e^{-(\|x - c_j\|/\sigma)^2}$$

with  $\sigma > 0$ . The training has been performed over 15000 iterations on the training set with a learning rate  $\eta = 0.001$ . The hyperparameters have been tuned with a grid search over all the possible combinations of the hyperparameters in:

$$(N, \rho, \sigma) \in \{25, 50, 70\} \times \{10^{-3}, 10^{-4}, 10^{-5}\} \times \{0.25, 0.5, 0.75, 1\}$$

selecting the network with the lowest error on the test set. The best result obtained an error of TODO on the training set and an error of TODO on the

test set (see table 5.1). The hyperparameters found are  $N = 0$ ,  $\rho = 0$  and  $\sigma = 0$ . In figure 5.1 the plot of the function computed by the RBFN found.

As already said before, both training procedures for MLP and RBFN have been performed using the gradient descent algorithm for 15000 iterations with a learning rate  $\eta = 0.001$ . This step has been implemented with TensorFlow's gradient descent optimizer, defined in `tf.train.GradientDescentOptimizer`. For each hyperparameter setting the training took 0.00s for the MLP and 0.00s for the RBFN on average.

TODO: explain overfitting/underfitting

TODO: MLP vs RBFN

### 3 TWO BLOCKS METHOD

The idea, here, is to train the network by first setting up a subset of weights and then by solving efficiently the least squares problem obtained. This speeds up the training of the network obtaining (similar?) results w.r.t the previous section.

For what regards the MLP all the weights  $w_{ji}$  and the biases  $b_j$ , with  $j \in \{1, \dots, N\}$  and  $i \in \{1, \dots, n\}$ , have been set up randomly. This reduces the minimization problem to the solution of:

$$\nabla_v E(\omega; \pi) = 0 \quad (3.1)$$

where the hyperparameters  $\pi$  are the values found in the previous section. Hence, the problem becomes a simple linear system:

$$\left( \frac{1}{2P} G^T G + \rho I \right) v = \frac{1}{2P} G^T y$$

where  $G \in \mathbb{R}^{P \times N}$  is defined as:

$$G_{rc} = g \left( \sum_{i=1}^n w_{ci} x_i^r - b_c \right)$$

This method has been performed 10000 times obtaining an error of TODO on the training set and an error of TODO on the test set (table 5.1). In figure 5.1 the plot of the function computed by the MLP found.

TODO: COMPARISON WITH Q1E1

An analogous procedure has been applied to RBFN. The centers  $c_j$ , with  $j \in \{1, \dots, N\}$  have been selected by using the K-means clustering algorithm. This allowed, as before, to reduce the problem to the solution of equation 3.1. The hyperparameters are the ones found in the previous section. Here  $G$  is defined as:

$$G_{ij} = \phi(\|x^i - c_j\|)$$

This method has been performed 10000 times obtaining an error of TODO on the training set and an error of TODO on the test set (table 5.1). In figure 5.1 the plot of the function computed by the RBFN found.

TODO: COMPARISON WITH Q1E2

TODO: implementation of both MLP and RBFN

Neural Network	N	$\sigma$	$\rho$	Training Error	Test Error	Time
Full MLP	75	3	$10^{-5}$	0.007	0.006	25.587s
Full RBFN	50	0.25	$10^{-5}$	0.005	0.005	34.770s
Extreme MLP	75	3	$10^{-5}$	0.002	0.002	0.001s
Unsupervised c. RBFN	50	0.25	$10^{-5}$	0.001	0.003	0.001s
TBD	TBD	TBD	TBD	TBD	<b>TBD</b>	TBD

Table 5.1: Results of the experiments described in the previous sections. The last column shows the amount of time needed to train each specific network.

## 4 DECOMPOSITION METHOD

TODO: E3

## 5 CONCLUSION

Conclusion. (Remember to say that all experiments have been performed on a kind of computer)

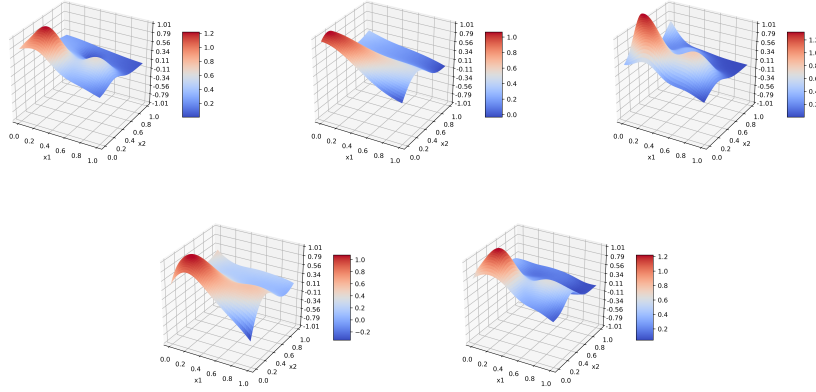


Figure 5.1: In order: the Franke's function and the functions obtained by Full MLP, Full RBFN, Extreme MLP, Unisupervised c. RBFN and TBD. The hyperparameters used for each network are specified in table 5.1.