

Optimization Methods for Machine Learning - Fall 2021

# Project # 1 - MLP and Generalized RBF Network

Laura Palagi & Giulia Di Teodoro

Posted on October 21, 2021 - due date November 18, 2021 (midnight 23:59)

## Instructions

The project will be done in groups formed by 2 to 4 people: each group must hand in their own answers. We will be assuming that, as participants in a Master's Degree course, every single student will be taking the responsibility to make sure his/her personal understanding of the solution to any work arising from such collaboration. Every student is expected to be able to explain to the teacher all the parts of the code submitted by him/her.

**Submission** Projects must be uploaded using a google form

<https://forms.gle/KLdDitg4jhkuyYv56>

Each team leader is in charge for uploading the files and he/she will receive an automatic acknowledgement by email after the upload.

The team must upload the python files following the instructions at the end of the project. A typed report in English and in pdf format must be uploaded too.

**The report must be of at most 4 pages excluded figures that must be put at the end. Do not include part of the python coding.**

**Evaluation criteria** Project is due at latest at midnight on the due date. For late submissions, the score will be decreased. It is worth 85% for the next 48 hours. It is worth 70% from 48 to 120 hours after the due date. It is worth 50% score after 120 hours delay. You have three questions Q1, Q2 - both with two points (Q1.1; Q1.2; Q2.1, Q2.2) - and Q3. You will be given the data set. Please note the teachers will check the overall quality of your results using a test set which is NOT provided to you (BLIND TEST SET).

The grade is "Italian style" namely in the range [0,30], being 18 the minimum degree to pass the exam. Answering only the first questions allows to get up to 24. Answering questions 1 and 2 allows to get up to 28. Answering Q1, Q2 and Q3 allows to get up to 30, while a bonus point can be taken by performing the additional bonus exercise. Please note that you need to write a last summary report of the project (mandatory) as explained at the end of this document.

The first project accounts for 35% of the total vote of the exam. For the evaluation of the first project the following criteria will be used:

1. 60% check of the implementation
2. 40% quality of the overall project as explained in the report.

In this assignment you will implement neural networks for regression. We want to reconstruct in the region  $[-2; 2] \times [-3; 3]$  a two dimensional function  $F : \mathbb{R}^2 \rightarrow \mathbb{R}$  which picture is sent in a separate file. You do not have the analytic expression of the function but only a data set obtained by randomly sampling 250 points  $x^p$  and adding to the value function  $F(x^p)$  a random noise in the range  $[-10^{-5}, 10^{-5}]$  so that the data set is

$$\{(x^p; y^p) : x^p \in \mathbb{R}^2; y^p \in \mathbb{R}; p = 1, \dots, 250\}$$

You must split randomly the data set into a training set and a test set. You will use the training set to evaluate the training error and the test set ONLY to evaluate the generalization error.

Choose the percentage of training data equal to 75% and use one of teammates' matricola numbers as the seed for the random split. Let  $P = 186$  be the number of instances in the training set.

Please note that in the last exercise of the project (bonus point) you can use all 250 points as training data and the overall performance is assessed by the teacher on a test set which is the same for all groups and is not made available to you beforehand.

In the exercise where you have to define a starting point for the optimisation procedure, we ask you to chose it randomly by fixing a seed. You can test different seeds, which can be described in the report, but you must fix the seed in the code to the value that gave you the best results.

### Question 1. (Full minimization)

You must construct a shallow Feedforward Neural Network (FNN) (one only hidden layer), both a MLP and a RBF network, that provides the model  $f(x)$  which approximates the true function  $F$ . We denote by  $\pi$  the hyper-parameters of the network to be settled by means of an heuristic procedure and  $\omega$  the parameters to be settled by minimizing the regularized training error

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

where the hyper parameter  $\rho$  stays in the range  $[10^{-5} \div 10^{-3}]$ .

1. (max score up to "22") Construct a shallow MLP with a linear output unit, namely

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

where  $\omega = (v, w, b)$ . The activation function  $g(\cdot)$  is the *hyperbolic tangent*

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

(with  $\sigma$  hyperparameter;  $g$  is available in Python with  $\sigma = 1$ : `Numpy.tanh`)

The hyper-parameters are

- the number of neurons  $N$  of the hidden layer
- the spread  $\sigma$  in the activation function  $g(t)$
- the regularization parameter  $\rho$

It must be possible to specify  $N$  as an input parameter of your training code.

Write a program which implements the regularized training error function  $E(v, w, b)$  (as obtained by (1) using  $f(x)$  of the MLP network) and uses a Python routine of the optimization toolbox (scipy.optimize) to determine the parameters  $v_j, w_{ji}, b_j$  which minimize it. The code must produce a plot of the approximating function found.

Analyse the occurrence of overfitting/underfitting varying the number of neurons  $N$  and the parameters  $\rho$  and  $\sigma$  (you can also decide to fix  $\sigma = 1$  if using the Python implementation of Tanh).

## 2. (max score up to "24") Construct a RBF network

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

where  $\omega = (v, c)$   $v \in \mathbb{R}^N$  and  $c_j \in \mathbb{R}^2$   $j = 1, \dots, N$ .

You must choose as RBF function  $\phi(\cdot)$  the Gaussian function

$$\phi(\|x - c_j\|) = e^{-(\|x - c_j\|/\sigma)^2} \quad \sigma > 0 \quad (3)$$

The hyper-parameters are

- the number of neurons  $N$  of the hidden layer
- the spread  $\sigma > 0$  in the RBF function  $\phi$
- the regularization parameter  $\rho$

It must be possible to specify  $N$  and  $\sigma$  as parameters that can be specified as input.

Write a program which implements the regularized training error function of the RBF network  $E(v, c)$  (as obtained by (1) using  $f(x)$  of the RBF network) and uses a Python routine of the optimization toolbox for its minimization with respect to both  $(v, c)$ . The code must produce a plot of the approximating function found.

Analyse the occurrence of overfitting/underfitting varying the number of units  $N$ , the spread parameters  $\sigma$  and  $\rho$ .

In the report, for both  $Q_{11}$  and  $Q_{12}$  you must state:

- the final setting for  $N$ ,  $\rho$  and  $\sigma$ ; how did you choose them and if you can put in evidence over/underfitting and if you can explain why;

- **1.** which optimization routine did you use for solving the minimization problem, **2.** the setting of its parameters (tolerance, max number of iterations etc) **3.** the returned message in output (successful optimization or others), **4.** starting/final value of the objective function, **5.** number of iterations, **6.** number of function/gradient evaluations, **7.** value of the norm of the gradient at the starting and final point if you evaluated it, **8.** computational time. Please note that it is not required to evaluate the gradient so derivative-free optimization routine may be applied or finite difference evaluation of the gradient can be used. Nevertheless using gradient-based method will improve performance and boost up your code. Explicit evaluation of the exact gradient will allow to obtain a greater score.

- the initial and final values of the error on the training set, defined as

$$E_{train} = \frac{1}{2P} \sum_{p=1}^{186} (f(x^p) - y^p)^2;$$

- the error on the test set. (NB: the training, test errors are computed using the formula (1) without the regularization term. This holds for all the exercises.)

- the plot of the function representing the approximating function obtained by the MLP and by the RBF networks;

- a comparison of performance between MLP and RBF networks both in terms of quality of approximation (training and test error) and in efficiency of the optimization (number of function/gradient evaluation and computational time needed to get the solution). Please put these values in a table as explained at the end.

## Question 2. (Two blocks methods)

1. (max score up to "26") Consider again the shallow MLP with linear output unit as defined at Ex. 1 of Question 1. Use the values of  $N, \rho, \sigma$  as you have fixed at Ex. 1.

Write a program which implements an *Extreme Learning* procedure, namely one that fix randomly the values of  $w_{ij}, b_j$  for all  $i, j$  and uses the "best" Python routine or any other optimization library to minimize the regularized quadratic convex training error  $E(v)$  of the only variables  $v \in \mathbb{R}^N$ .

As for the random generation of  $w, b$  you have no restrictions on the way of proceeding. We suggest to identify a range and repeat the random choice more than once. In the run file, the  $w, b$  still need to be defined in a random-like way, definitions "by hand" will not be accepted.

In the report you must state:

- which optimization routine you use for solving the quadratic minimization problem and the setting of its parameters. Compare the performance of the optimization process w.r.t Full optimization of Question 1 (Ex. 1)

- the values of the error on the training and test set; Compare these values with the ones obtained by the Full optimization of Question 1 (Ex. 1).
- the plot of the function representing the approximating function obtained by the Extreme Learning MLP in comparison with the true one.

2. (max score up to "28")

Consider again an RBF network as in Ex. 2 of Question 1.

The number of RBF units  $N$  of the hidden layer and the positive parameter  $\sigma$  in the RBF function are those chosen at Ex 2. Question 1.

Write a program which implements a method with unsupervised selection of the centers. You can select the centers by randomly picking  $N$  of the  $P$  points of the training set. We suggest to repeat the random choice more than once.

Find the optimal weights by minimizing the regularized error  $E(v)$  using a suitable Python routine or any other optimization library to minimize the quadratic convex function.

In the report you must state:

- which optimization routines of the Optimization toolbox you have used to solve the weights subproblems and the setting of its parameters. Compare the performance of the optimization process w.r.t Full optimization of Question 1 (Ex. 1)
- the value of the error on the training and test sets; Compare these values with the ones obtained by the Full optimization of Question 1 (Ex. 1).
- the plot of the function representing the approximating function obtained by the RBF with unsupervised selection of the centers in comparison with the true one.

### Question 3. (Decomposition method)

3. (max score up to "30")

Consider either the shallow MLP or the shallow RBF network. Set the hyper-parameters parameter at the values you selected in Question 1.

Write a program which implements a two block decomposition method, which alternates the convex minimization with respect to the output weights  $v$  and the non convex minimization with respect to the other parameters (either  $(w, b)$  for the MLP or the centers  $c$  for the RBF network respectively).

You must use appropriate Python routines of the optimization toolbox for solving the two block minimization problems. In this case, **it is required that you calculate the gradient of the nonconvex block**. Consider adjusting the tolerances of the optimizers as the number of iterations increases.

Compare the results with those obtained at the corresponding exercise of Question 1.

In the report you must state:

- which optimization routines of the Optimization toolbox you used to solve the two block subproblems and the setting of the parameters.
- the stopping criteria of the decomposition procedure; consider the use of early stopping rules.
- the number of outer iterations (number of subproblems solved), number of function/gradient evaluation and computational time needed to get it.
- the value of the error on the training and test sets;
- the plot of the function representing the approximating function obtained by the RBF with block decomposition in comparison with the true one.
- Comparison with the corresponding network obtained by random selection at the corresponding Exercise of Question 2. Comparison are both in the quality of the network (test error), and in the efficiency of the optimization procedure (number of function/gradient evaluation and computational time needed to get the solution). Please put these values in a table.

### Additional bonus exercise. Score up to "30 cum laude" (or 1 point bonus)

**Build your best model.** You can use all the 250 samples in the data set to construct the model and you may use different starting points than the one used in the preceding exercises (i.e. you can avoid fixing any randomization seed with your matricula). Submit your best model (either a MLP or RBF network). However, for this The teachers will use a new test set (randomly defined and not included within the data set) to check the quality of the results and the bonus point is assigned on the basis of an accuracy ranking.

### Final remarks for the report

In the final report (pdf) it is (**MANDATORY**) to gather into a final table (an example below) the comparison among all the implemented methods - in term of accuracy in learning and computational effort in training.

Ex		settings			Final train error	Final test error	Initial FOB	final FOB	optimization time
		$N$	$\sigma$	$\rho$					
Q1.1	Full MLP								
Q1.2	Full RBF								
Q2.1	Extreme MLP								
Q2.2	Unsupervised $c$ RBF								
Q2.3	The BEST MODEL								

### Instructions for python code

You are allowed to organize the code as you prefer **BUT** for each question  $ij = \{11, 12, 21, 22, 3\}$  you have to create a different folder where you must provide two files:

- A file called **run\_ij\_GroupName.py** (e.g. run\_11\_Example.py, run\_12\_Example.py). This file will be the only one executed in phase of verification of the work done. It can include all the classes, functions and libraries you used for solving the specific question *ij* but it has to print ONLY:

1. Number of neurons  $N$  chosen
2. Value of  $\sigma$  chosen
3. Value of  $\rho$  chosen
4. Values of other hyperparameters (if any)
5. Optimization solver chosen (e.g. L-BFGS, CG, NTC, Nelder-Mead....)
6. Number of function evaluations
7. Number of gradient evaluations
8. Time for optimizing the network (from when the solver is called, until it stops)
9. Training Error (defined as  $E_{train} = \frac{1}{2P_{train}} \sum_{i=1}^{P_{train}} (y_i - \tilde{y}_i)^2$ )
10. Test Error (defined as above but on the test set)

Please, in order to maintain the code as clean/clear as possible, do not define functions inside the run\_ij\_GroupName.py file, but make it call functions defined in other files.

- For each question you have to provide another file called **Test\_ij\_GroupName**. In this file you have to write a function called **ICanGeneralize** which takes in input a new matrix of points  $X_{test} \in R^{P \times n}$  and returns a vector  $Y_{pred}$ , which is the prediction returned by the best weights you found in point 1.1.
- DO NOT insert the data.csv file in the folders you will submit.