

Lab assignment 3: Radial basis functions neural networks

Academic year 2023/2024

Subject: Introduction to computational models
4th course Computer Science Degree (University of Córdoba)

15th November 2023

Resumen

This lab assignment serves as familiarisation for the student with radial basis functions (RBF) neural networks. In this way, a RBF neural network will be developed, using Python and the `scikit-learn` library¹. In this sense, the assignment will also serve as familiarisation with external libraries, widely used in machine learning (`numpy`, `pandas`...). In addition, we will introduce the problem of bias in machine learning models through `fairlearn`². The student must implement the algorithm and analyse the effect of different parameters over a given set of real-world datasets. Delivery will be made using the task in Moodle authorized for this purpose. All deliverables must be uploaded in a single compressed file indicated in this document. The deadline for the submission is **12th December 2023**. In case two students submit copied assignments, neither of them will be scored.

1. Introduction

The work to be done in this lab assignment consists in implementing a RBF neural network with a training stage divided into three steps:

1. Application of a clustering algorithm which will be used to establish the centres of the RBF (input-to-hidden-layer's weights).
2. The RBF radius adjustment is done by means of a simple heuristic (distance average to the rest of the centres).
3. Hidden-to-output's weights learning:
 - For regression problems, using the Moore-Penrose's pseudo-inverse.
 - For classification problems, using a logistic regression linear model.

The student should develop a Python's script able to train a RBF neural network with the aforementioned characteristics. This programme will be used to train models able to classify as accurate as possible a set of databases available in Moodle. Also, an analysis about the obtained results will be included. For the COMPAS dataset that we studied in lab assignment 2, we will also perform an algorithmic bias analysis to contrast the behaviour of the models in different demographic groups. **This analysis will greatly influence the qualification of this assignment.**

In the statement of the assignment, indicative values are provided for all parameters. However, it will be positively evaluated if the student finds other values for these parameters able to achieve better results.

¹<http://scikit-learn.org/>

²<https://fairlearn.org/>

Section 2 describes a series of general guidelines when implementing the training algorithm for RBF neural networks. Section 3 explains the experiments to be carried out once the algorithm is implemented. Finally, section 4 specifies the files to be delivered for this assignment.

2. Implementation of the RBF neural network training algorithm

2.1. Model's architecture to be considered

The RBF neural network models should have the following architecture:

- An input layer with as many neurons as input variables the dataset has.
- A hidden layer with a number of neurons specified by the user. It is important to highlight that, in the two previous lab assignment, the number of hidden layer was variable. However, for this lab assignment, we are going to consider just one hidden layer. The type of all the neurons in the hidden layer will be RBF (in contrast to the sigmoidal neurons used in the previous lab assignments).
- An output layer with as many neurons as output variables the dataset has.
 - When considering regression datasets, all the output neurons will be linear (i.e. similar to the sigmoidal neurons without the application of the $\frac{1}{1 + e^{-x}}$ transformation).
 - When considering classification datasets, all the output neurons will be softmax. The softmax function is already implemented by the logistic regression algorithm used for adjusting the weights of the output layer.

2.2. Weights adjustment

The instructions given in the class slides should be followed so that the training is carried out as follows:

1. Application of a clustering algorithm that will serve to establish the centres of the RBF (input-to-output layer weights). For classification problems, the centroid initialisation will be random and stratified, n_1 patterns³. For regression problems, n_1 will be randomly selected. After initialising the centroids, the `sklearn.cluster.KMeans` class will be used, with only one centroid initialisation (`n_init`) and a maximum of 500 iterations (`max_iter`).
2. To adjust the radius of the RBF, a simple heuristic will be applied (the half of the distance average to the rest of the centres). This is, the radius of the j -th neuron will be⁴:

$$\sigma_j = \frac{1}{2 \cdot (n_1 - 1)} \sum_{i \neq j} \|c_j - c_i\| = \frac{1}{2 \cdot (n_1 - 1)} \sum_{i \neq j} \sqrt{\sum_{d=1}^n (c_{jd} - c_{id})^2}. \quad (1)$$

3. Learning the weights from hidden-to-output layer.

- For regression problem, it is done using the Moore-Penrose pseudo-inverse. This is:

$$\beta_{((n_1+1) \times k)}^T = (\mathbf{R}^+)_{((n_1+1) \times N)} \mathbf{Y}_{(N \times k)} = \quad (2)$$

$$= \left(\mathbf{R}_{((n_1+1) \times N)}^T \times \mathbf{R}_{(N \times (n_1+1))} \right)^{-1} \mathbf{R}_{((n_1+1) \times N)}^T \mathbf{Y}_{(N \times k)} \quad (3)$$

³For this, the `sklearn.model_selection.train_test_split` method can be used. It performs one or more stratified dataset partitions, this is, keeping the ratio of patterns belonging to each class in the original dataset https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html

⁴Consider using the functions `pdist` and `squareform` of `scipy` to obtain the distances matrix

where \mathbf{R} is the matrix containing the outputs of the RBF neurons, β is a matrix containing a vector of parameters for each of the outputs to be predicted, and \mathbf{Y} is a matrix with the target outputs. To perform these operations, we will use the matrix functions of `numpy`, which is a dependence of `scikit-learn`.

- For classification problems, it is done using a logistic regression linear model. Using the `sklearn.linear_model.LogisticRegression` class, providing a value for the C parameter in order to apply regularisation. Note that in this library what we are specifying is the cost value C (importance of the approximation error versus the regularisation error), in such a way that $\eta = \frac{1}{C}$. We will use the `saga` solver and maximum of iterations (`max_iter`) of 10.

3. Experiments

We will test different configurations of the neural network and execute each configuration with five seeds (1, 2, 3, 4 and 5). Based on the results obtained, the average and standard deviation of the error will be obtained. For the regression problems, only the MSE will be shown. However, for classification problems, the CCR (the percentage of correct classified patterns) will be shown. To analyse algorithmic bias in the COMPAS dataset we will use the false positive rate, which is the most appropriate metric for this particular problem, but you can optionally include the false negative rate.

To assess how the implemented algorithm works, we will run it on three different regression datasets:

- *Sin-function dataset*: This dataset is composed of 120 training patterns and 41 testing patterns. It has been obtained by adding some random noise to the sin function (see Figure 1).

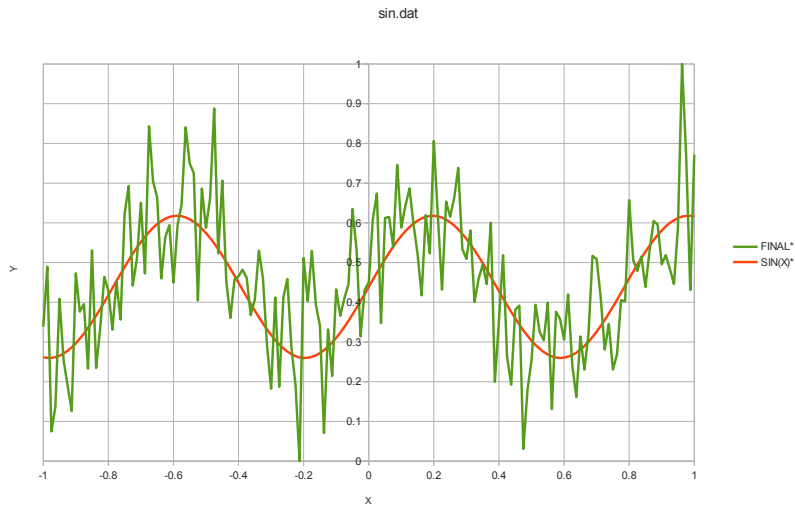


Figure 1: Data representation of the data included in the sin-function estimation problem.

- *Quake dataset*: this dataset is composed by 1633 training patterns and 546 testing patterns. It corresponds to a database in which the objective is to find out the strength of an earthquake (measured on the Richter scale). As input variables, we use the depth of focus, the latitude at which it occurs and the longitude⁵.

⁵see <https://sci2s.ugr.es/keel/dataset.php?cod=75> to seek more information.

- *Auto MPG dataset*: The Auto MPG dataset, available in the UCI Machine Learning Repository⁶, is a well-known dataset frequently used in machine learning and data analysis tasks. It contains information about various car models, including their attributes and fuel efficiency. This dataset is often used for regression and predictive modeling tasks, particularly to predict a car's miles per gallon (MPG) based on its features (cylinders, displacement, horsepower, etc.). The original dataset has 398 samples and 7 attributes. We binarized discrete attributes and removed model name so the final dataset attributes list is: cylinders, displacement, horsepower, weight, acceleration, modelYear, USA, Europe, Japan and MPG (target variable).

And two classification datasets:

- *ProPublica COMPAS*: This dataset is about the performance of COMPAS algorithm, a statistical method for assigning risk scores within the United States criminal justice system created by Northpointe. It was published by ProPublica in 2016⁷, claiming that this risk tool was biased against African-American individuals (we will deal with this in the following assignment). In this dataset, they analyzed the COMPAS scores for "risk of recidivism" so each individual has a binary "recidivism" outcome, that is the prediction task, indicating whether they were rearrested within two years after the first arrest (the charge described in the data). We reduced the original dataset from 52 to 9 attributes similarly to the original dataset: sex, age, age_cat, juv_fel_count, juv_misd_count, juv_other_count, priors_count, race, c_charge_degree. The prediction variable is whether the individual will be rearrested in two years or not.

1. sex: binary sex.
2. age: numerical age.
3. age_cat: categorical (age < 25, age ≥ 25 and age < 45, age ≥ 45).
4. juv_fel_count: a continuous variable containing the number of juvenile felonies.
5. juv_misd_count: a continuous variable containing the number of juvenile misdemeanors.
6. juv_other_count: a continuous variable containing the number of prior juvenile convictions that are not considered either felonies or misdemeanors.
7. priors_count: a continuous variable containing the number of prior crimes committed.
8. race⁸: binary attribute (0 means 'white' and 1 means 'black').
9. c_charge_degree: Degree of the crime. It is either M (Misdemeanor), F (Felony), or O (not causing jail time).

This database presents a class imbalance problem for each group, as there are 1361 negative vs 1313 positive labels for black people whereas there are 839 negative vs 544 positive labels for white people. Typically, this translates into models that will tend to label black people as reoffenders whilst it will tend to label white people as not reoffenders. This can be assessed through the false positive rate for each group. We will perform an exploratory analysis of this database within the practical sessions. In this lab assignment, the input variables of this dataset have been previously standardized in the `csv` version. The race variable is the last one in the `csv`.

- *noMNIST dataset*: originally, this dataset was composed by 200,000 training patterns and 10,000 test patterns, with a total of 10 classes. Nevertheless, for this lab assignment, the size of the dataset has been reduced in order to reduce the computational cost. In this sense, the dataset is composed by 900 training patterns and 300 test patterns. It includes a set of

⁶Raw dataset and attribute description available here: <https://archive.ics.uci.edu/dataset/9/auto+mpg>

⁷<https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>

⁸Recommended reading: There's No Scientific Basis for Race—It's a Made-Up Label

letters (from a to f) written with different typologies or symbols. They are adjusted to a squared grid of 28×28 pixels. The images are in grey scale in the interval $[-1,0; +1,0]^9$. Each of the pixels is an input variable (with a total of $28 \times 28 = 784$ input variables) and the class corresponds to a written letter (a, b, c, d, e y f , with a total of 6 classes). Figure 2 represents a subset of 180 training patterns, whereas figure 3 represents a subset of 180 letters from the test set. Moreover, all the letters are arranged and available in Moodle in the files `train_img_nomnist.tar.gz` and `test_img_nomnist.tar.gz`, respectively.



Figure 2: Subset of letters belonging to the training dataset.



Figure 3: Subset of letters belonging to the test dataset.

The average and standard deviation of two measures (regression) or four measures (classification) should be computed:

- Regression: average and standard deviation of training and testing MSE .
- Classification: average and standard deviation of training and testing CCR .

At least, the following configurations should be tried:

- *Network architecture*:
 - For all the datasets, consider a number hidden neurons (n_1) equal to the 5 %, 15 %, 25 % and 50 % of the total number of patterns of the dataset. In this stage, for classification problems use L1 regularisation and $\eta = 10^{-5}$.
- For the classification problems, once decided the best architecture, try the following values for η : $\eta = 10^{-3}, \eta = 10^{-2}, \eta = 10^{-1}, \eta = 1, \eta = 10^1, \eta = 10^2, \eta = 10^3$, along with the two types of regularisation (L2 y L1). What is happening? Compute the difference in number of coefficients for COMPAS and noMNIST dataset when the regularisation type is modified (L2 vs L1)¹⁰.

⁹Check <http://yaroslavvb.blogspot.com.es/2011/09/notmnist-dataset.html> for more information.

¹⁰The coefficients are in the `coef_` attribute of the logistic regression object. Consider that if the absolute value of a coefficient is lower than 10^{-5} , then the coefficient is null

- Implement the option `-v` so that, in classification problems, the classifier is trained using `sklearn.linear_model.LogisticRegressionCV` instead of using the class `sklearn.linear_model.LogisticRegression`. This will automatically calculate (without looking at the test values), the optimal value of η . Set the object so that the values to be tested (Cs) are $\eta = 10^{-3}$, $\eta = 10^{-2}$, $\eta = 10^{-1}$, $\eta = 1$, $\eta = 10^1$, $\eta = 10^2$ and $\eta = 10^3$ and that the optimization is done using a stratified k -fold with $k = 3$. Compare, on the two classification datasets, the results you get using this strategy (under the assumption of L1 and the best ratio of RBFs obtained above). Compare also the values you choose for *eta* with those obtained by the experimentation of the previous point. What can the differences be due to?
- For one of the classification problems, run the script considering the problem as if it was regression (i.e. the classification parameter is False and compute the *CCR* rounding the predictions to the closest integer). What is happening for this situation?

As a guideline, the training and generalisation errors achieved by a linear regression (using Weka) over the three regression datasets is shown:

- *sin dataset*: $MSE_{\text{train}} = 0,02968729$; $MSE_{\text{test}} = 0,03636649$.
- *Quake dataset*: $MSE_{\text{train}} = 0,03020644$; $MSE_{\text{test}} = 0,02732409$.
- *MPG dataset*: $MSE_{\text{train}} = 0,000000$; $MSE_{\text{test}} = 0,00750$.

Also, the training *CCR* and the test *CCR* achieved by a logistic regression (using Weka) over the two classification datasets is shown:

- *COMPAS dataset*: $CCR_{\text{train}} = 67,7348\%$; $CCR_{\text{test}} = 66,8514\%$.
- *noMNIST dataset*: $CCR_{\text{train}} = 80,4444\%$; $CCR_{\text{test}} = 82,6667\%$.

The student should be able to improve this error values with some of the configurations.

3.1. File format

The files containing the datasets will be CSV, in such a way that the values will be separated by commas. In this sense, there are no headers. In order to read the files properly, the function `read_csv` from `pandas` should be used. In the COMPAS database, the 'race' variable has been placed in the last column so that it can be easily processed and integrated with `fairlearn`.

4. Deliverables

The files to be submitted will be the following:

- Report in a pdf file describing the programme implemented, including results, tables and their analysis.
- Python script.

4.1. Report

The report for this lab assignment must include, at least, the following content:

- Cover with the lab assignment number, its title, subject, degree, faculty department, university, academic year, name, DNI and email of the student
- Index of the content with page numbers.
- Description of the steps for the RBF training stage (**1 page maximum**).

- Experiments and results discussion:
 - Brief description of the datasets used.
 - Brief description of the values of the parameters considered.
 - Results obtained, according to the format specified in the previous section.
 - Discussion/analysis of the results. The analysis must be aimed at justifying the results obtained instead of merely describing the tables. This analysis should include algorithmic bias assessment in COMPAS. Take into account that this part is extremely decisive in the lab assignment qualification. The inclusion of the following comparison items will be appreciated:
 - Test confusion matrix of the best neural network model achieved for the *noMNIST* database. Analysing the errors, **including the images of some letters for which the model mistakes**, to visually check if they are confusing. Comparison between the confusion matrix obtained for this assignment against the one obtained in the previous lab assignment.
 - Computational time needed for the training step for *noMNIST* dataset and comparison against the computational time spent in the previous lab assignment.
- Bibliographic references or any other material consulted in order to carry out the lab assignment different to the one provided by the lecturers (if any).

Although the content is important, the presentation, including the style and structure of the document will also be valued. The presence of too many spelling mistakes can decrease the grade obtained.

4.2. Executable and source code

Together with the report, the script file prepared to be run in the UCO's machines (concretely, test using `ssh` on `ts.uco.es`) must be included. In addition, all the source code must be included. The script developed should receive the following command-line arguments¹¹:

- Argument `-t, --train_file`: Indicates the name of the file that contains the training data to be used. This argument is compulsory, and without it, the program can not work.
- Argument `-T, --test_file`: Indicates the name of the file that contains the testing data to be used. If it is not specified, training data will be used as testing data.
- Argument `-s, --standardize`: Boolean indicating whether the data sets are to be standardized (inputs and outputs for regression, only inputs for classification). If not specified, we will assume that it is not standardized.
- Argument `-c, --classification`: Boolean that indicates whether it is a classification problem. If it is not specified, we will suppose that it is a regression problem.
- Argument `-r, --ratio_rbf`: Indicates the radius (by one) of RBF neurons with respect to the total number of patterns in training. If not specified, use 0,1.
- Argument `-l, --l2`: Boolean that indicated if L2 regularisation is used, instead of L1. If it is not specified, L1 will be used.
- Argument `-e, --eta`: Indicates the value for the *eta* (η) parameter. By default, use $\eta = 1e-2$.
- Argument `-f, --fairness`: Boolean that indicated if fairness metrics should be extracted from predictions. Assumes that the group is stored as the last variable of the input variables. By default, it is disabled.

¹¹To process the input sequence, the `click` library will be used.

- Argument `-o, --outputs`: Indicates the number of output columns of the dataset (always placed at the end). By default, use $o = 1$.
- Argument `-v, --logisticcv`: Boolean for activating the use of the class `LogisticRegressionCV` for classification problems. By default, it is assumed that it will not be applied (therefore, `LogisticRegression` is used).
- (Kaggle) Argument `-p, --pred`: Boolean that indicates if the prediction mode is used.
- (Kaggle) Argument `-m, --model_file`: Indicates the directory in which the trained models are saved (in the training mode, without the flag `p`) or the file containing the model that will be used (in the prediction mode, with the flag `p`).
- Argument `--help`: It shows the help of the program (use the one automatically generated by the `click` library)

An example of execution can be seen in the following output¹²:

```

1 i02gupep@NEWTS:~/imc/workspace/la3$ ./rbf.py --help
2 Usage: rbf.py [OPTIONS]
3
4 5 executions of RBFNN training
5
6 RBF neural network based on hybrid supervised/unsupervised training. We run
7 5 executions with different seeds.
8
9 Options:
10 -t, --train_file TEXT  Name of the file with training data.
11 -T, --test_file TEXT  Name of the file with test data. [required]
12 -s, --standarize      Standarize input variables (and outputs if it is a
13 regression problem). [default: False]
14 -c, --classification  The problem considered is a classification problem.
15 [default: False]
16 -r, --ratio_rbf FLOAT Ratio of RBF neurons (as a fraction of 1) with
17 respect to the total number of patterns. [default:
18 0.1]
19 -l, --l2              Use L2 regularization instead of L1 (logistic
20 regression). [default: False]
21 -e, --eta FLOAT      Value of the regularization parameter for logistic
22 regression. [default: 0.01]
23 -f, --fairness        Evaluates prediction using fairlearn metrics. It is
24 assumed that last input variable is the group
25 variable. [default: False]
26 -o, --outputs INTEGER Number of columns that will be used as target
27 variables (all at the end). [default: 1]
28 -p, --pred            Use the prediction mode. [default: False]
29 -v, --logisticcv      Consider LogisticRegressionCV. [default: False]
30 -m, --model TEXT      Directory name to save the models (or name of the
31 file to load the model, if the prediction mode is
32 active).
33 --help               Show this message and exit.
34
35 # Example with COMPAS dataset and fairness parameter
36 i02gupep@NEWTS:~/imc/workspace/la3$ ./rbf.py -t ./csv/train_compas.csv -T ./csv/
37 test_compas.csv -c --l2 -f
38 -----
39 Seed: 1
40 -----
41 Number of RBFs used: 405

```

¹²To make the developed code to work in the UCO machines, the packages `click` and the last version of the package `scikit-learn` should be installed, using the following commands:

```

pip install scikit-learn --user --upgrade
pip install click --user --upgrade

```



```

41 | /home/javi/anaconda3/lib/python3.11/site-packages/sklearn/linear_model/_sag.py:350:
    | ConvergenceWarning: The max_iter was reached which means the coef_ did not converge
42 | warnings.warn(
43 | Training MSE: 0.204360
44 | Test MSE: 0.212523
45 | Training CCR: 68.20%
46 | Test CCR: 66.69%
47 | -----
48 | Seed: 2
49 | -----
50 | Number of RBFs used: 405
51 | /home/javi/anaconda3/lib/python3.11/site-packages/sklearn/linear_model/_sag.py:350:
    | ConvergenceWarning: The max_iter was reached which means the coef_ did not converge
52 | warnings.warn(
53 | Training MSE: 0.204601
54 | Test MSE: 0.212833
55 | Training CCR: 68.23%
56 | Test CCR: 66.69%
57 | -----
58 | Seed: 3
59 | -----
60 | Number of RBFs used: 405
61 | /home/javi/anaconda3/lib/python3.11/site-packages/sklearn/linear_model/_sag.py:350:
    | ConvergenceWarning: The max_iter was reached which means the coef_ did not converge
62 | warnings.warn(
63 | Training MSE: 0.204681
64 | Test MSE: 0.212844
65 | Training CCR: 68.01%
66 | Test CCR: 66.85%
67 | -----
68 | Seed: 4
69 | -----
70 | Number of RBFs used: 405
71 | /home/javi/anaconda3/lib/python3.11/site-packages/sklearn/linear_model/_sag.py:350:
    | ConvergenceWarning: The max_iter was reached which means the coef_ did not converge
72 | warnings.warn(
73 | Training MSE: 0.204666
74 | Test MSE: 0.213130
75 | Training CCR: 67.93%
76 | Test CCR: 66.69%
77 | -----
78 | Seed: 5
79 | -----
80 | Number of RBFs used: 405
81 | /home/javi/anaconda3/lib/python3.11/site-packages/sklearn/linear_model/_sag.py:350:
    | ConvergenceWarning: The max_iter was reached which means the coef_ did not converge
82 | warnings.warn(
83 | Training MSE: 0.204510
84 | Test MSE: 0.212759
85 | Training CCR: 68.20%
86 | Test CCR: 66.69%
87 | *****
88 | Summary of results
89 | *****
90 | Training MSE: 0.204564 +- 0.000118
91 | Test MSE: 0.212818 +- 0.000194
92 | Training CCR: 68.11% +- 0.12%
93 | Test CCR: 66.72% +- 0.07%
94 | Training FN0: 59.85% +- 0.34%
95 | Training FN1: 32.55% +- 0.27%
96 | Test FN0: 51.55% +- 0.85%
97 | Test FN1: 37.66% +- 0.31%
98 | Training FP0: 14.11% +- 0.66%
99 | Training FP1: 31.02% +- 0.16%
100 | Test FP0: 21.59% +- 0.31%
101 | Test FP1: 28.89% +- 0.29%
102 |

```

```

103
104 # En los siguientes ejemplos, los CCRs salen 0 porque es un problema de regresión
105 i02gupep@NEWTS:~/imc/workspace/la3$ ./rbf.py -t ./train_mpg.csv -T ./test_mpg.csv -r 0.5
    -s
106 -----
107 Seed: 1
108 -----
109 Number of RBFs used: 156
110 Training MSE: 0.037012
111 Test MSE: 0.277343
112 Training CCR: 0.00%
113 Test CCR: 0.00%
114 -----
115 Seed: 2
116 -----
117 Number of RBFs used: 156
118 Training MSE: 0.040975
119 Test MSE: 0.140427
120 Training CCR: 0.00%
121 Test CCR: 0.00%
122 -----
123 Seed: 3
124 -----
125 Number of RBFs used: 156
126 Training MSE: 0.037219
127 Test MSE: 0.182746
128 Training CCR: 0.00%
129 Test CCR: 0.00%
130 -----
131 Seed: 4
132 -----
133 Number of RBFs used: 156
134 Training MSE: 0.044984
135 Test MSE: 0.145340
136 Training CCR: 0.00%
137 Test CCR: 0.00%
138 -----
139 Seed: 5
140 -----
141 Number of RBFs used: 156
142 Training MSE: 0.040233
143 Test MSE: 0.174549
144 Training CCR: 0.00%
145 Test CCR: 0.00%
146 *****
147 Summary of results
148 *****
149 Training MSE: 0.040084 +- 0.002914
150 Test MSE: 0.184081 +- 0.049390
151 Training CCR: 0.00% +- 0.00%
152 Test CCR: 0.00% +- 0.00%
153
154 i02gupep@NEWTS:~/imc/workspace/la3$ ./rbf.py -t ./train_mpg.csv -T ./test_mpg.csv -r
    0.15 -s
155 -----
156 Seed: 1
157 -----
158 Number of RBFs used: 46
159 Training MSE: 0.087675
160 Test MSE: 0.101936
161 Training CCR: 0.00%
162 Test CCR: 0.00%
163 -----
164 Seed: 2
165 -----
166 Number of RBFs used: 46
167 Training MSE: 0.105456

```

```

168 | Test MSE: 0.127247
169 | Training CCR: 0.00%
170 | Test CCR: 0.00%
171 | -----
172 | Seed: 3
173 | -----
174 | Number of RBFs used: 46
175 | Training MSE: 0.086466
176 | Test MSE: 0.107544
177 | Training CCR: 0.00%
178 | Test CCR: 0.00%
179 | -----
180 | Seed: 4
181 | -----
182 | Number of RBFs used: 46
183 | Training MSE: 0.086680
184 | Test MSE: 0.095453
185 | Training CCR: 0.00%
186 | Test CCR: 0.00%
187 | -----
188 | Seed: 5
189 | -----
190 | Number of RBFs used: 46
191 | Training MSE: 0.089423
192 | Test MSE: 0.116676
193 | Training CCR: 0.00%
194 | Test CCR: 0.00%
195 | *****
196 | Summary of results
197 | *****
198 | Training MSE: 0.091140 +- 0.007234
199 | Test MSE: 0.109771 +- 0.011175
200 | Training CCR: 0.00% +- 0.00%
201 | Test CCR: 0.00% +- 0.00%
202 |
203 | i02gupep@NEWTS:~/imc/workspace/la3$ ./rbf.py -t ./train_sin.csv -T ./test_sin.csv -r
    | 0.15 -o 1
204 | -----
205 | Seed: 1
206 | -----
207 | Number of RBFs used: 18
208 | Training MSE: 0.012100
209 | Test MSE: 0.104193
210 | Training CCR: 0.83%
211 | Test CCR: 2.44%
212 | -----
213 | Seed: 2
214 | -----
215 | Number of RBFs used: 18
216 | Training MSE: 0.011399
217 | Test MSE: 0.200716
218 | Training CCR: 0.83%
219 | Test CCR: 2.44%
220 | -----
221 | Seed: 3
222 | -----
223 | Number of RBFs used: 18
224 | Training MSE: 0.011953
225 | Test MSE: 0.102114
226 | Training CCR: 0.83%
227 | Test CCR: 2.44%
228 | -----
229 | Seed: 4
230 | -----
231 | Number of RBFs used: 18
232 | Training MSE: 0.012089
233 | Test MSE: 0.082532

```

```

234 Training CCR: 0.83%
235 Test CCR: 2.44%
236 -----
237 Seed: 5
238 -----
239 Number of RBFs used: 18
240 Training MSE: 0.011961
241 Test MSE: 0.092528
242 Training CCR: 0.83%
243 Test CCR: 2.44%
244 *****
245 Summary of results
246 *****
247 Training MSE: 0.011901 +- 0.000258
248 Test MSE: 0.116417 +- 0.042847
249 Training CCR: 0.83% +- 0.00%
250 Test CCR: 2.44% +- 0.00%

```

4.3. [OPTIONAL] Save the model to a file.

During the training stage, the script can save the model trained as a pickle¹³. This will allow to use the trained model to predict the outputs of the **Kaggle** dataset.

To save the model, it is necessary to use the `-m` parameter. An execution example is as follows:

```

1 i02gupep@NEWTS:~/imc/workspace/la3$ ./rbf.py -t train.csv -T val.csv -l -c -r 0.1 -m
  model
2 -----
3 Seed: 1
4 -----
5 Number of RBFs used: 157
6 Training MSE: 0.121783
7 Test MSE: 0.140203
8 Training CCR: 82.60%
9 Test CCR: 78.67%
10 -----
11 Seed: 2
12 -----
13 Number of RBFs used: 157
14 Training MSE: 0.122003
15 Test MSE: 0.139427
16 Training CCR: 82.35%
17 Test CCR: 79.05%
18 -----
19 Seed: 3
20 -----
21 Number of RBFs used: 157
22 Training MSE: 0.123603
23 Test MSE: 0.138931
24 Training CCR: 82.35%
25 Test CCR: 78.86%
26 -----
27 Seed: 4
28 -----
29 Number of RBFs used: 157
30 Training MSE: 0.122647
31 Test MSE: 0.138973
32 Training CCR: 82.16%
33 Test CCR: 78.67%
34 -----
35 Seed: 5
36 -----
37 Number of RBFs used: 157
38 Training MSE: 0.124017

```

¹³<https://docs.python.org/3/library/pickle.html>

```

39 | Test MSE: 0.142853
40 | Training CCR: 82.48%
41 | Test CCR: 78.67%
42 | *****
43 | Summary of results
44 | *****
45 | Training MSE: 0.122811 +- 0.000874
46 | Test MSE: 0.140077 +- 0.001461
47 | Training CCR: 82.39% +- 0.15%
48 | Test CCR: 78.78% +- 0.15%

```

Once the execution is finished, there will be a folder named “model” containing 5 pickles. Each one corresponds with the generated model for each seed. In order to obtain the predictions, one of these 5 pickles should be chosen.

```

1 | i02gupep@NEWTS:~/imc/workspace/la3$ ls model/
2 | 1.pickle 2.pickle 3.pickle 4.pickle 5.pickle

```

4.4. [OPTIONAL] Obtaining the predictions for Kaggle.

Once the model is saved to a pickle, it is possible to obtain the output predictions for the Kaggle dataset. For this, `-m` and `-p` parameters should be used. Below is an example:

```

1 | i02gupep@NEWTS:~/imc/workspace/la3$ ./rbf.py -T test.csv -p -m model/2.pickle
2 | Id,Category
3 | 0,1
4 | 1,0
5 | 2,0
6 | 3,0
7 | 4,0
8 | 5,0
9 | 6,1
10 | 7,0
11 | 8,0
12 | 9,0
13 | 10,0
14 | 11,0
15 | 12,1
16 | 13,0
17 | 14,1
18 | 15,0
19 | 16,0
20 | 17,1
21 | 18,0
22 | 19,0
23 | 20,0
24 | 21,0
25 |
26 |
27 | ...
28 |
29 | 1887,0
30 | 1888,0
31 | 1889,0
32 | 1890,1
33 | 1891,1
34 | 1892,0
35 | 1893,0
36 | 1894,0
37 | 1895,1
38 | 1896,0
39 | 1897,0
40 | 1898,1
41 | 1899,0

```

The output can be redirected to a `csv` file:

```
1 i02gupep@NEWS:~/imc/workspace/la3$ ./rbf.py -T kaggle.csv -p -m modelo/2.pickle > submission.csv
```

This file is ready to be uploaded to Kaggle.