

Exam Report

Data Mining and Neural Networks

Daniel Gerardo GIL SANCHEZ
daniel.gilsanchez@student.kuleuven.be
MSc Statistics

Supervisor: Prof. Johan Suykens

Academic year 2018-2019

Contents

Contents	i
1 Exercise Session 1	1
1.1 Function Approximation (noiseless case)	1
1.2 The role of the hidden layer and output layer	2
1.3 Function Approximation (noisy case)	5
1.4 Curse of dimensionality	13
2 Exercise Session 2	16
2.1 Santa Fe laser data - time-series prediction	16
2.2 Alphabet recognition	21
2.3 Breast Cancer Wisconsin - classification problem	23
3 Exercise Session 3	27
3.1 Dimensionality reduction by PCA analysis	27
3.2 Input selection by Automatic Relevance Determination (ARD)	28
References	30
4 Appendix	31

1 Exercise Session 1

1.1 Function Approximation (noiseless case)

To talk about overfitting and underfitting, it is important to define both situations first. Overfitting is when the training process is totally driven by the training error, so the network memorizes the training examples. When new data are presented, the error is going to be large because the network has not learned how to generalize to new situations. On the other hand, underfitting is when the network does not fit the training data well nor generalize to new situations.

Figure 1 shows some cases of underfitting and overfitting. In the upper left, there is a simple nonlinear function (displayed as red crosses) and a simple network with just one hidden neuron (black line). In this case, the fit is not good enough when the target function reaches either -1 or 1 in the vertical axis. The upper right plot shows a typical underfitting situation, when the fit of the function is only good in a small range in the domain of the input variable but behaves poorly outside this range. On the other hand, both plots in the bottom show overfitting in a simple nonlinear function as well as in a complex one (left and right respectively). In both situations the data are fitted (almost) perfectly, leaving no room to generalize new situations. It is important to mention that this case does not have any noise, so the fit is not wiggled as it usually is in overfitting situations.

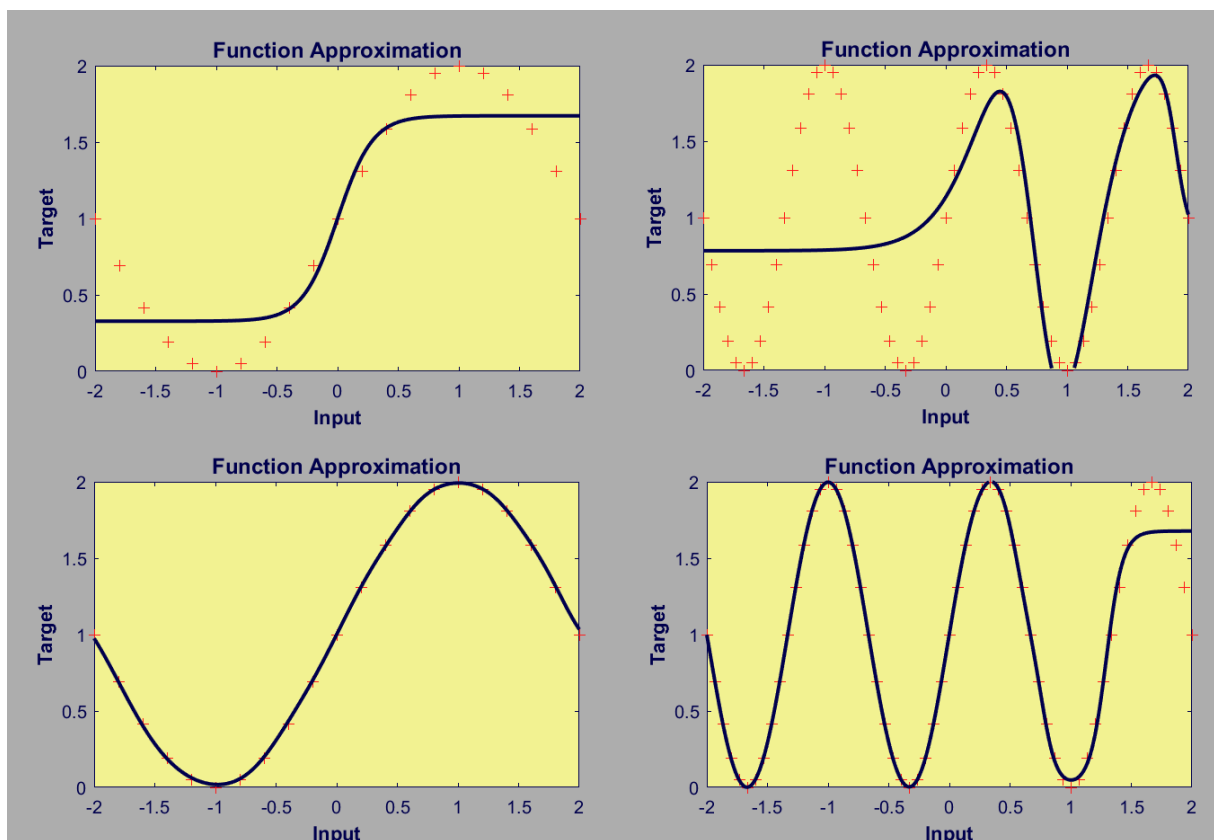


Figure 1: Underfitting (top) and Overfitting (bottom) scenarios

The reason why any of these situations appear in any model building process (not only neural networks) is because of the complexity of the model in relation with the data at hand. If the number of parameters in the model is close to the total number of observations, then it is very likely to have overfitting. Conversely, if the number of parameters in the model is much smaller than the total number of observations, then there is no chance of overfitting but

the likelihood of getting a underfitted model increases. Of course, it is difficult to know how complex a model should be to avoid both scenarios, so other methods are used to deal with these situations, such as early stopping and regularization.

1.2 The role of the hidden layer and output layer

1. Given a set of input and output patterns $\{(x_p, y_p)\}_1^P$, how would you solve the linear regression problem using a neural network?

This can be seen as the McCulloch-Pitts model with the only difference that instead of using a nonlinear activation function, such as $\tanh(\cdot)$, a linear activation is used (the **identity function**). This model is a weighted sum of the inputs, denoted as x in a n -dimensional space, plus a bias term that represents a threshold to activate a neuron. Following McCulloch-Pitts model for illustration, only **one hidden layer** is needed to learn a linear regression function with only **one neuron**. A representation of such neural network is:

$$y = f\left(\sum_{i=1}^n w_i x_i + b\right) = f(a) \quad (1)$$

Where w_i are the interconnection weights in the context of neural networks or the parameter of each variable in the context of regression. b is the bias term in neural network context or the intercept in regression jargon. x are the input patterns or covariates and y is the target or dependent variable. As it is mentioned above, f is the **identity function**, so equation 1 results in:

$$y = \sum_{i=1}^n w_i x_i + b \quad (2)$$

Which is similar to the equation used in regression. Figure 2 shows the architecture of a neural network that represents a linear regression¹.

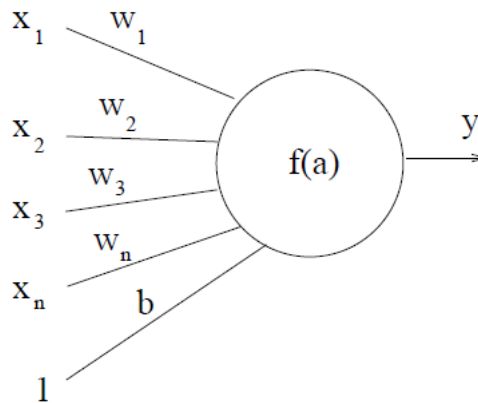


Figure 2: Neural network architecture, f is the identity function.

2. and 3. A vector of 21 observations equally spaced in the interval $[0,1]$ is created, and the response variable is calculated as $y_p = -\cos(0.8\pi x_p)$ with $p = 1, \dots, 21$. In Figure 3 this relationship is plotted. There, it can be seen that a linear regression may fit the data well in the range $[0.3,0.7]$ but not further. In fact, if the domain of the input variable was extended, it would be possible to see that the function reaches to 1 and decreases again until -1, so a linear model would not be a good idea. In this sense, if the data at hand is only in the range $[0,1]$, a model with the architecture chosen in the previous numeral could give a decent representation. However, if data outside the range $[0,1]$ were to be used, the model would not be useful anymore.

¹This image is taken from the lectures of this course

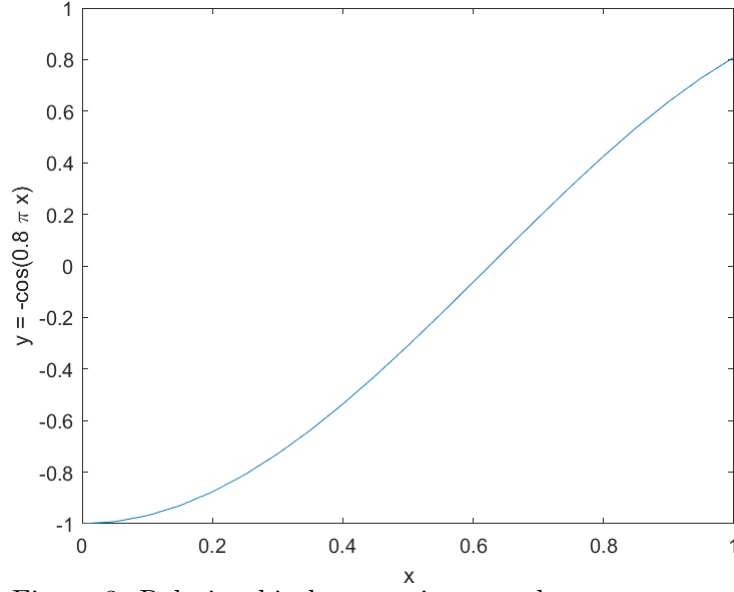


Figure 3: Relationship between input and output patterns

4. The next step is to train a neural net with one hidden layer and two neurons. The configuration of the network is changed to prevent it from rescaling the inputs and outputs. It is important to mention that the training process starts with different initial values for weights and biases and different division of the data into training, validation and test sets. For this reason, the network is trained several times to increase the likelihood of finding a good generalization and avoid local optima.

To do so, the neural net is trained 100 times using different starting values for weights and biases and different divisions of the dataset for training, validation and testing. Each time that the neural net reaches convergence, usually by early stopping, the Mean Squared Error (MSE) is saved and then compared to keep the solution with the lowest performance value. Notice, that in order to get reproducible results the seed in the random number generator is set to 0689432, which is my student number.

5. The activation of each neuron can be computed by using the estimated weights and bias in the calculation of the linear combination of inputs and plugging this result in the activation function, as follows:

$$a = \sigma(Vx + \beta) \quad (3)$$

Where a is a matrix with the activation value for each input pattern in each neuron. σ is the activation function. V and β are the weights and biases of the hidden neurons and x are the input patterns. In this case, σ is the *hyperbolic tangent sigmoid function* (*tansig*), V and β are:

$$V = \begin{bmatrix} 1.8896 \\ -5.2995 \end{bmatrix}, \beta = \begin{bmatrix} -1.1702 \\ -0.2583 \end{bmatrix} \quad (4)$$

Table 1 shows these activation values.

Input	Pattern (x)	Neuron 1	Neuron 2
1	0	-0.824	-0.253
2	0.05	-0.792	-0.48
3	0.1	-0.754	-0.657
4	0.15	-0.71	-0.783
5	0.2	-0.66	-0.866
6	0.25	-0.603	-0.919
7	0.3	-0.539	-0.952
8	0.35	-0.469	-0.971
9	0.4	-0.392	-0.983
10	0.45	-0.309	-0.99
11	0.5	-0.222	-0.994
12	0.55	-0.13	-0.996
13	0.6	-0.036	-0.998
14	0.65	0.058	-0.999
15	0.7	0.151	-0.999
16	0.75	0.242	-1
17	0.8	0.329	-1
18	0.85	0.41	-1
19	0.9	0.486	-1
20	0.95	0.555	-1
21	1	0.617	-1

Table 1: Activation value in each hidden neuron

6. Figure 4 plots the values of the output, and each activation value against the input patterns. It can be seen that the first neuron follows the pattern of the output, whereas the second neuron is inversely related with it.

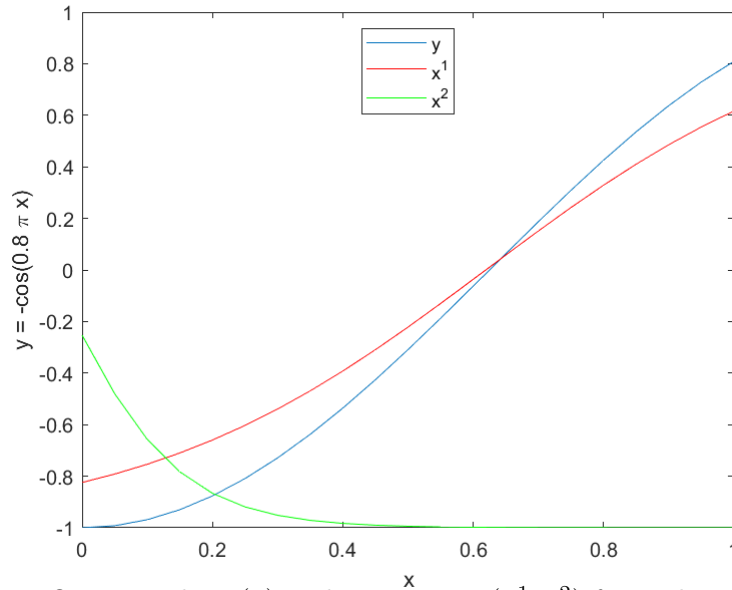


Figure 4: Output values (y) and activations (x^1, x^2) for each input pattern

7. Once the network is estimated, it is possible to calculate an estimation for the output

pattern. This is conducted by applying the following formula:

$$output = W\sigma(Vx + \beta) + \beta^* \quad (5)$$

Where σ , V , β and x are defined as in equation 3. W and β^* are the weights and bias of the output neuron, which in this case are:

$$W = \begin{bmatrix} 1.3383 \\ 0.1571 \end{bmatrix}, \beta^* = 0.1429 \quad (6)$$

Notice that the activation function of the output layer is the identity function, thus it is not expressed in equation 5. These results are shown in Figure 5 and compared with the true function. It is clear that the neural net perfectly fitted the true function.

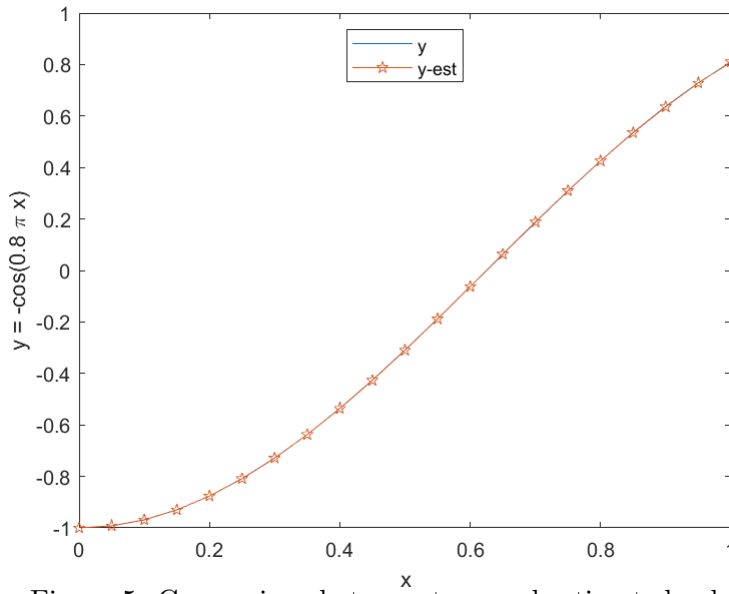


Figure 5: Comparison between true and estimated values

1.3 Function Approximation (noisy case)

A lot of different neural networks can be fitted by changing the options given in the exercise. To investigate how a neural net changes when different specifications are set, each of the options is analyzed by holding the remain options constant. For instance, to measure whether the fit of a neural net is better when more data points are used, the remain characteristics of the neural net such as the number of hidden neurons and training algorithm are constant.

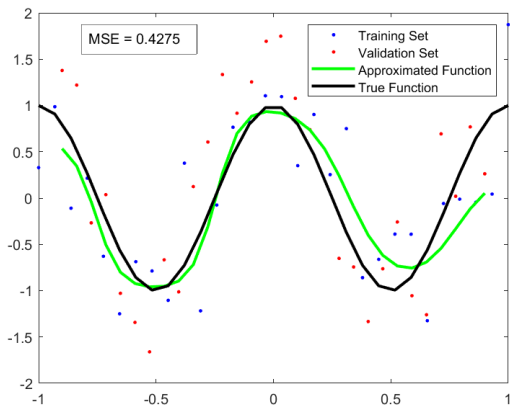
The comparison of the fit between different neural networks is conducted visually and through the Mean Squared Error (MSE) calculated on a test set. This test set is completely different and independent of the train set. In this way one can be sure that the test set in one neural net is not going to be used as a training set in other neural net.

In regard to local optima, it is well known that each time a neural net is trained, different initial weights and biases, and different divisions of the data into training, validation and test sets are considered. Thus, each neural network is trained 100 times, and the net with the lowest MSE is kept for comparison. As before, in order to get reproducible results the seed in the random number generator is set to 0689432.

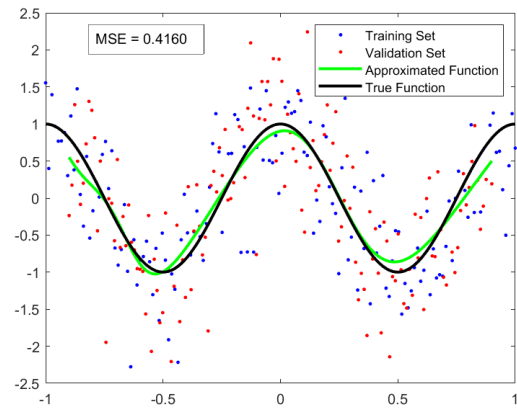
1. Size of the training data set

Three different sizes are tested in this comparison: 30, 150 and 1200 data points. In this case, the standard deviation is set to 0.6 to control the amount of noise, the number of hidden neurons is 5, the training algorithm is the Levenberg-Marquardt, which is the default in MATLAB, to improve generalization early stopping is used, thus no regularization is used and the selection of initial weights is at random (although a seed is set for reproducibility).

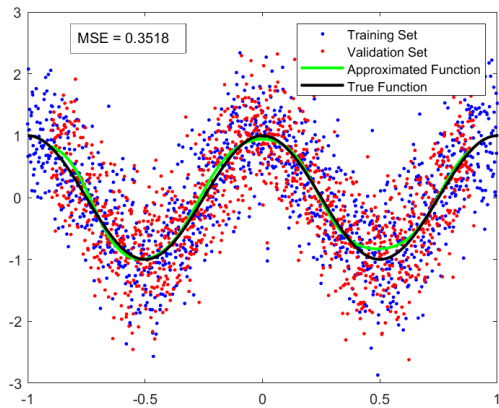
Figure 6 shows the fit and the MSE for each training dataset size. The lowest MSE is presented in the neural net with 1200 observations, where the approximated function is almost the same as the true function over the whole range of the input variable. On the other hand, in both neural nets with fewer observations, it is clear that the approximated function is very similar to the true function but, they are not as powerful in predicting new data (see MSE). This is a great example to show that more data does not necessarily lead to better estimations, since it all depends on several characteristics of the neural net and in this case the quality of the prediction was approximately the same, with way less data points (see plot 6a).



(a) $N = 30$



(b) $N = 150$



(c) $N = 1200$

Figure 6: Comparison of neural nets with different training sample sizes

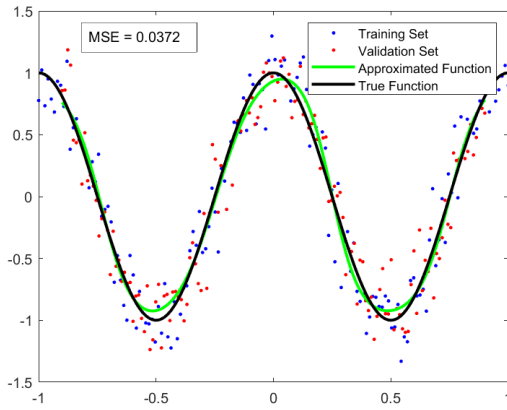
2. Amount of noise on the training data

In this case three different standard deviations are considered to add noise to the cosine wave: 0.2, 0.6 and 1.2. As before, the remaining characteristics are held constant and

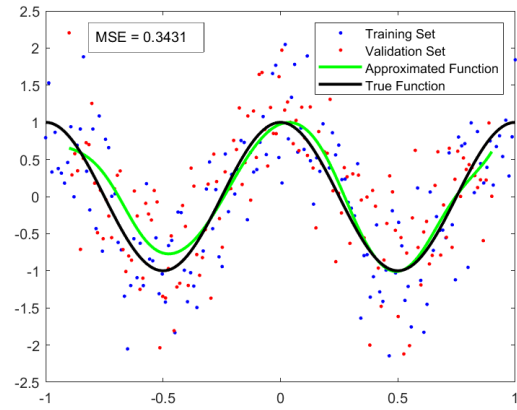
the size of the dataset is fixed at 150. Note that the choice of the sample size could be anything because the purpose of this section is to analyze the differences in the fit when different noises are added to the raw data.

Figure 7 shows the fit and the MSE for each training dataset with different noise. The lowest MSE is presented in the neural net with the smallest standard deviation. This result is expected due to the fact that there is not much variability, so it is easier to predict. On the other hand, the training data set that had a standard deviation of 1.2, presents the highest MSE and therefore is less reliable to make predictions. In regard to the visual representation, it is clear that in plot 7c the lack of fit is usually located from 0.25 onward.

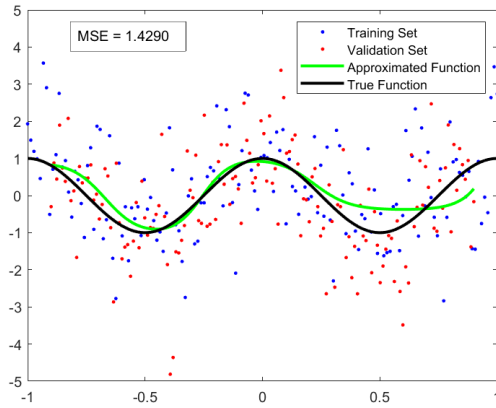
This example complements the results showed in the previous section where the fit was better when the number of input patterns were large. It actually tells that it does not really matter the number of observations but how much variance there is in them. The larger the variance, the harder it is to get a good fit.



(a) $sd = 0.2$



(b) $sd = 0.6$



(c) $sd = 1.2$

Figure 7: Comparison of neural nets with different noise in the training sample

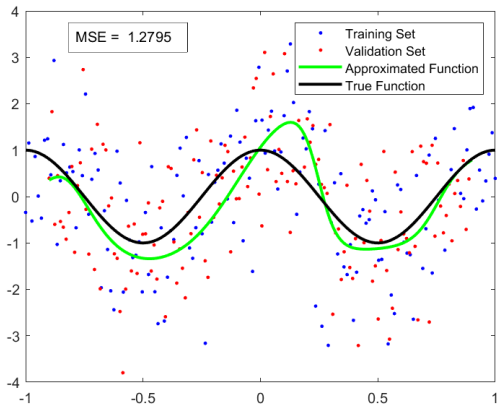
3. Number of hidden neurons

In this case three different scenarios are considered in regard to the number of hidden neurons: 6, 10 and 20. As before, the remaining characteristics are held constant, the size of the training dataset is 150 and the amount of noise considered is a standard

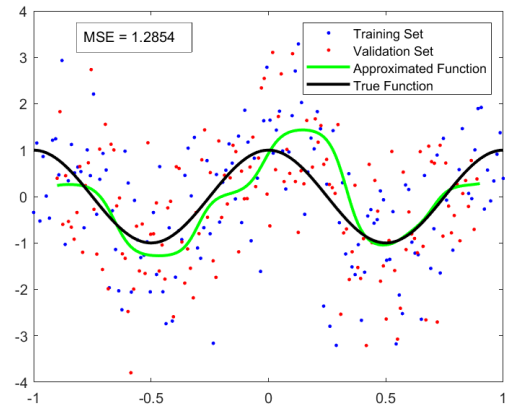
deviation of 1.2. There are two reasons to choose this amount of noise: first, because it showed the worst performance in comparison to the other two scenarios and second, because it is a more realistic example.

In regard to the number of hidden neurons, it was already shown in the previous comparison that 5 hidden neurons were not enough to get a decent representation of the true function. For this reason, the number of hidden neurons considered are larger than 5.

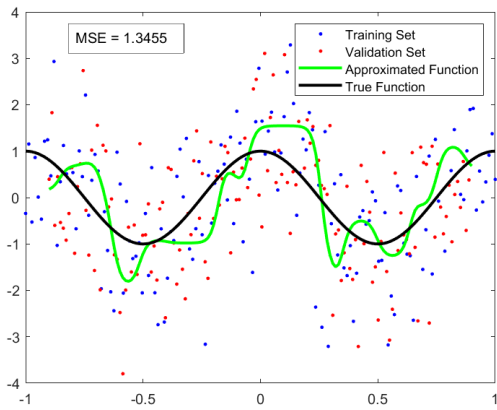
In Figure 8 the comparison is presented. The neural network that has the lowest MSE is the one that has 6 hidden neurons. Graphically, it can be seen that it is closer to the true function and it is very similar to the approximated function of the net with 10 hidden neurons. In contrast, the net with 20 hidden neurons is a clear example of overfitting for two reasons: first, the MSE is higher than in any other case and second, the approximation function is not as smooth as the true function, meaning that the net memorized the training patterns instead of having learned the general pattern.



(a) 6 hidden neurons



(b) 10 hidden neurons



(c) 20 hidden neurons

Figure 8: Comparison of neural nets with different number of hidden neurons

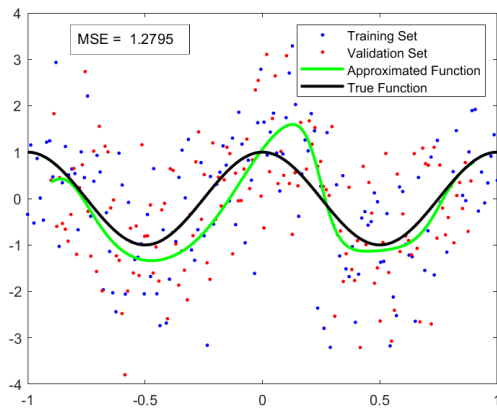
4. Training algorithms

In this case four different training algorithms are compared: backpropagation, conjugate gradient, quasi-Newton and Levenberg-Marquardt, which is the default in MATLAB. As before, the remaining characteristics are held constant, the size of the training dataset

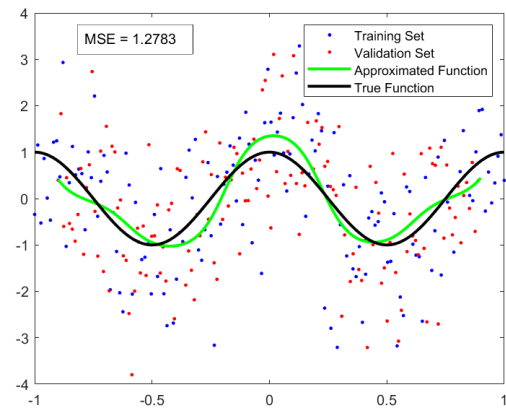
is 150, the amount of noise considered is a standard deviation of 1.2 and the number of hidden units is 6.

It is important to mention that backpropagation is a term that sometimes is used to refer to gradient descent algorithm, so in MATLAB the gradient descent algorithm was conducted.

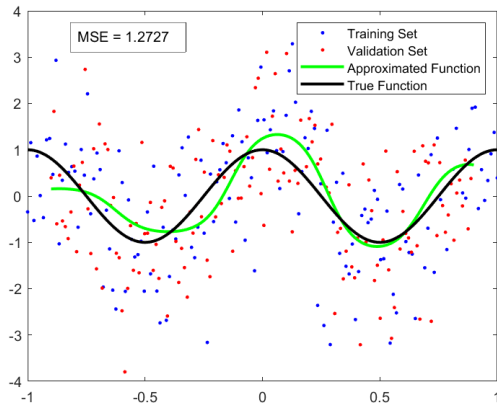
Figure 9 shows the difference in fit between each algorithm. Visually, it can be seen that all of the algorithms produce similar results. In regard to the MSE, the algorithm that presents the lowest value is the conjugate gradient and the highest value is the Quasi-Newton. It is important to mention that even though the difference in MSE is not large between algorithms, some of them were slower in comparison to Levenberg-Marquardt, which is well known to perform well in function fitting problems and it is considerably faster than the other ones in such situations (see Matlab, 2018). The message of this example is that not always the same algorithm performs better than other ones, it always depends on the problem and the data at hand. The best way to keep the best solution is to try different algorithms and use the one that has better performance.



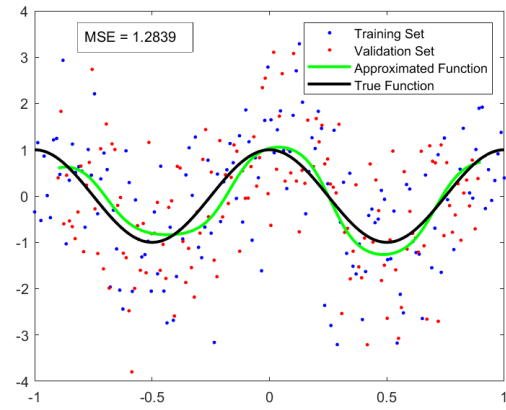
(a) Levenberg-Marquardt



(b) Gradient Descent (Backpropagation)



(c) Scaled Conjugate Gradient



(d) BFGS Quasi-Newton

Figure 9: Comparison of neural nets using different training algorithms

5. Local minimum vs Early stopping

In this case, two different stopping rules are considered: local minimum and Early stopping. Local minimum is performed by checking the gradient in each iteration and stopping

the algorithm when this value is smaller than a threshold, which by default is $1e - 5$. On the other hand, in early stopping rule, the error on the validation set is monitored during the training process and the algorithm stops when the error on the validation set start to rise.

To control which stopping rule to use in MATLAB, the number of validation checks is set to 0 in the local minimum case and 6 in early stopping (default value). This number of validations is the number of iterations that the error is increasing after reaching a local minimum. If no iteration is considered, the algorithm stops at the first local minimum and if six iterations are considered, the algorithm stops at a deeper point.

As before, the remaining characteristics are held constant, the size of the training dataset is 150, the amount of noise considered is a standard deviation of 1.2, the number of hidden units is 6 and the training algorithm is Levenberg-Marquardt, which is the faster one.

The difference between both stopping rules is presented in Figure 10. So far this criterion is the one that impacts the most the performance of a neural network. In the left plot, the approximated function is far from the true function and as a consequence the error in the test set increased a lot. This example is important in the sense that shows that the optimization process of a neural net is far from trivial and one should be aware of the default configurations that each software usually has, because it significantly impacts the final result.

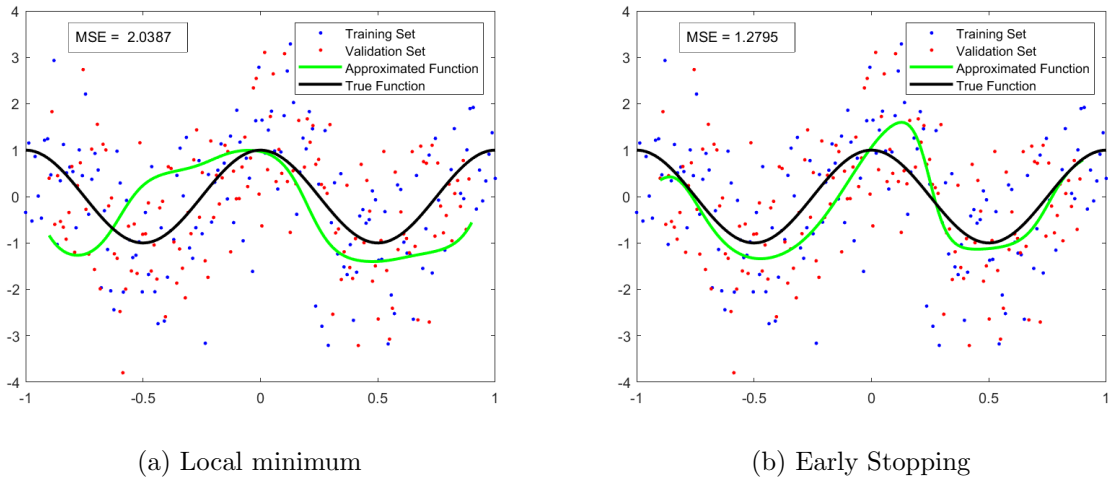


Figure 10: Comparison of neural nets using different stopping rules

6. Regularization constant

In this case different regularization constants are considered: from 0, where no regularization is conducted, to 0.5 where in the optimization process the original cost function and the weight decay term have the same importance. Four different constants are considered, namely 0, 0.1, 0.25 and 0.5.

As before, the remain characteristics are held constant, the size of the training dataset is 150, the amount of noise considered is a standard deviation of 1.2, the number of hidden neurons is 6, the training algorithm is Levenberg-Marquardt.

In Figure 11 the difference in fit between each regularization constant is presented. Graphically, there is a huge difference between no regularization and the constant equals to 0.5, because in the latter one (plot 11d) the approximated function is just a straight line

with a small slope. On the other hand, there is an improvement in the performance when the regularization is equal to 0.1 (plot 11b), the approximated function is almost the same as the true function. In this comparison is clear that the selection of the regularization constant is really important because it impacts the quality of the results. In fact, given that this choice is not straightforward, a Bayesian approach can be taken by considering all weights and biases as random variables.

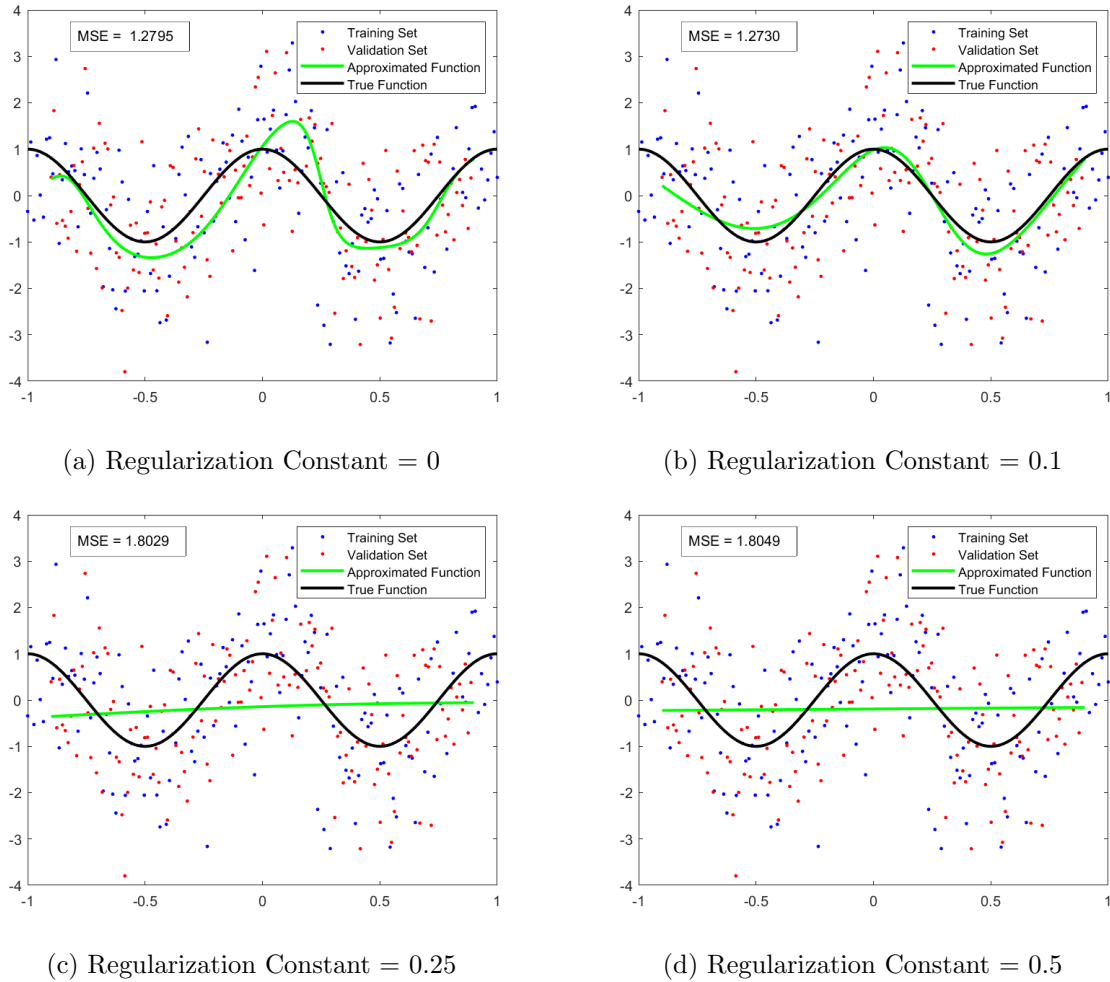


Figure 11: Comparison of neural nets using different regularization constant

7. Initial weights

As it was mentioned above, each time a neural net is trained, different initial values are given to weights and biases. Usually, these values are random, but they can be set if prior information is available. In this case different starting values are considered by controlling the seed in the random number generator. Four seeds are compared: 0689432, 12345, 67890 and 111111. Note that this comparison can be done by just training different times the same neural net, but to get reproducible results it is conducted in this way. As before, the remaining characteristics of the neural nets are held constant.

Figure 12 shows the fit of each neural net using different starting weights and biases. Visually, all approximated functions are quite similar, and in some cases they are closer to the true function (see Figures 12b and 12d). In regard to the MSE, the differences are

not large but still they give the message of the importance of considering different starting values, since it increases the likelihood of getting a better solution.

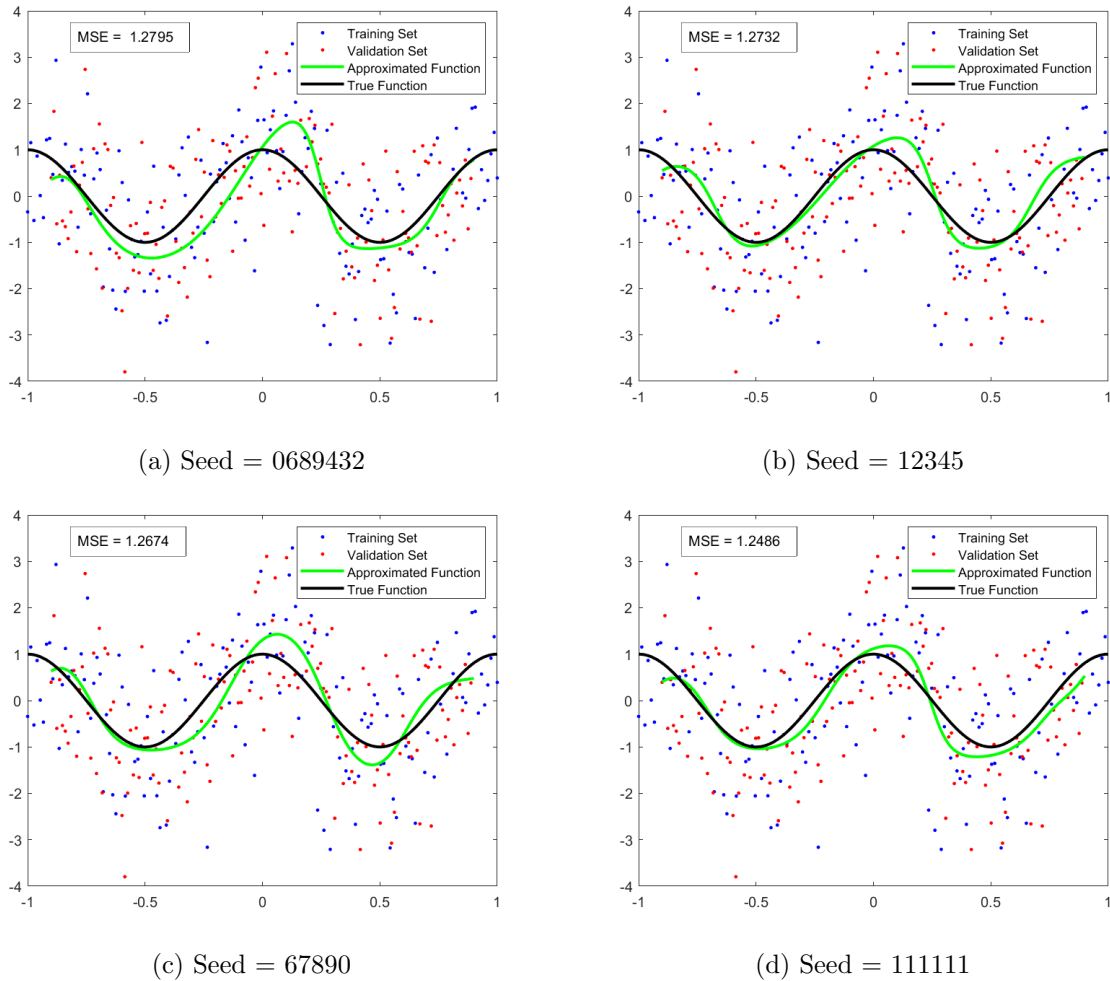


Figure 12: Comparison of neural nets using different initial weights and biases

All in all, different characteristics of a neural network were investigated and compared. Some interesting results were found in regard to the performance of a neural net in a test dataset. First, it does not really matter the number of observations available, but the variance there is them. The more variability there is, the harder it is to get a good solution. Second, the number of hidden units impacts directly the quality and performance of a neural net. Too complex neural nets tend to overfit the data which lead to bad predictions. Third, regularization is a useful tool to avoid overfitting and forget about the number of hidden units. Fourth, the choice of the training algorithm has to be influenced by the problem that is being analyzed, some of them are faster or lead to better results than others in determined situations. Finally, always train a neural network several times to avoid bad local optima solutions.

1.4 Curse of dimensionality

In this section, different neural networks are conducted to get a good approximation of the $\text{sinc}(t)$ function, which is defined as follows:

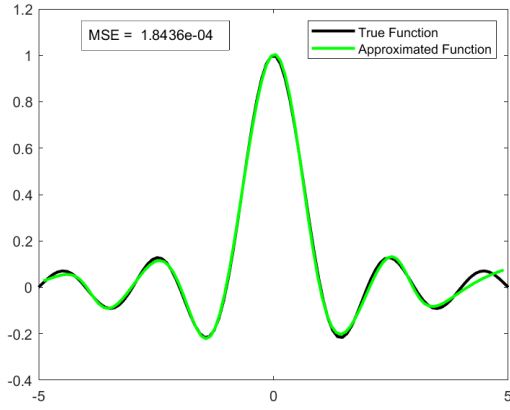
$$f(x) = \text{sinc} \left(\sqrt{\sum_{i=1}^m x_i^2} \right), \quad \text{sinc}(t) = \begin{cases} \frac{\sin(\pi t)}{\pi t} & t \neq 0 \\ 1 & t = 0 \end{cases} \quad (7)$$

Three different cases are considered: when m is equal to 1, 2 and 5. In each case, different training algorithms are used, and the Mean Squared Error is calculated to make a comparison between them and decide which leads to better results. This performance measure is computed on a test set, which is completely different and independent to the training set. In this way one can be sure that the test set in one neural net is not going to be used as a training set in other neural net. In addition, to avoid bad local optima solutions, each neural network is trained several times and the net with the lowest MSE is kept for comparison. As before, in order to get reproducible results the seed in the random number generator is set to 0689432.

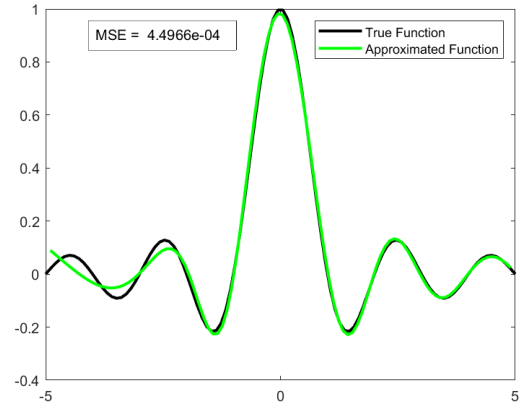
1. $m = 1$ (one-dimensional sinc function)

In this case, only one input is considered and it is created as linearly spaced vector between -5 and 5 of size 100. The test set is similarly created but the range considered is between -4.9 and 4.9 . The network used for this problem is a 1-5-1 network with *tansig* activation functions in the hidden layer and linear activation function in the output layer. Different number of hidden neurons are tried, from 2 to 10, and 5 is the final choice because it leads to good results in a rather simple network. Different training algorithms are considered, namely Levenberg-Marquardt, Bayesian Regularization, Quasi-Newton and Scaled conjugate gradient. Each neural net is trained 100 times, where different initial random weights and biases are used.

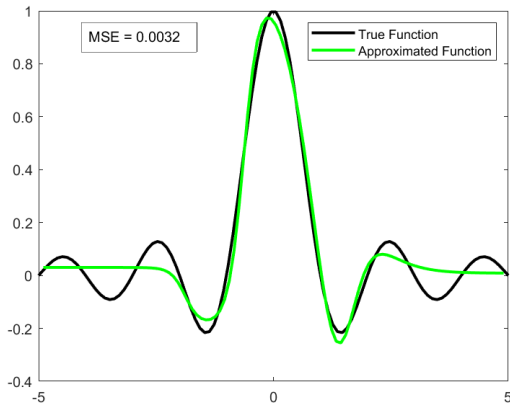
Figure 13 shows the difference in fit between each algorithm and the corresponding MSE. Visually, it can be seen that there are some differences in the fit, especially in the Quasi-Newton and the scaled conjugate gradient (plots 13c and 13d). In both, the approximated function is close to the true function in the range $[-2, 2]$, but outside this range the quality of the prediction is not as good. In regard to the MSE, the algorithm that presents the lowest value is the Levenberg-Marquardt and the highest value is the scaled conjugate gradient. It is important to mention that even though the difference in MSE is not large between the LM algorithm and the Bayesian regularization, the latter took almost 10 times more time to get the solution. Therefore, the LM algorithm is preferred in this case.



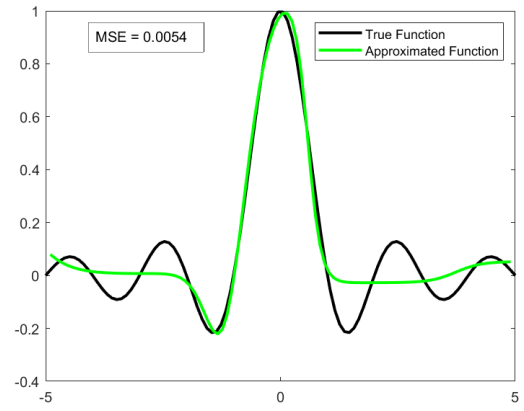
(a) Levenberg-Marquardt



(b) Bayesian regularization



(c) BFGS Quasi-Newton



(d) Scaled Conjugate Gradient

Figure 13: Comparison of neural nets using different training algorithms, when $m = 1$

2. $m = 2$ (mexican hat function)

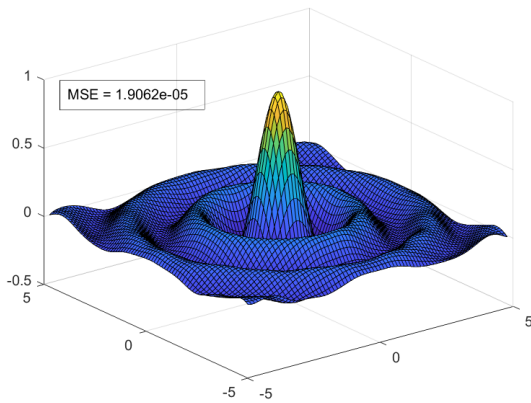
In this case, the input dataset is harder to get because it is necessary to get all possible pair combinations between the first and the second variable to build a grid. To do so, for x_1 a sequence of 101 numbers between -5 and 5 is created in a vector and then each element is repeated 101 times, resulting in a vector of 10201 elements. Then for x_2 another sequence of 101 numbers between -5 and 5 is created and the whole vector is repeated 101 times, resulting also in a vector of 10201 elements. This is a small example of how both input patterns look like:

$$x = (x_1, x_2) = \begin{bmatrix} -5 & -5 & -5 & \dots & 5 & 5 & 5 \\ -5 & -4.9 & -4.8 & \dots & 4.8 & 4.9 & 5 \end{bmatrix}$$

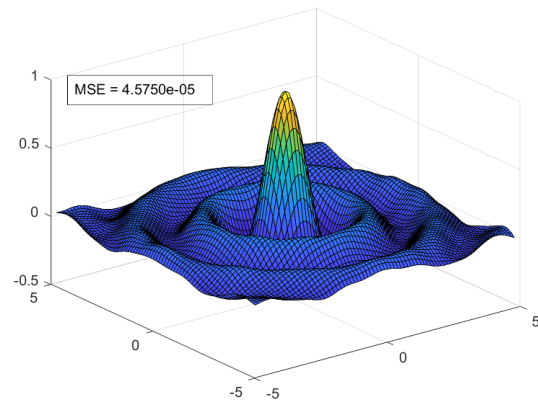
In a similar fashion, a grid of 79×79 observations is created for the test set with a sequence of numbers between -4.9 and 4.9 , resulting in two vectors with 6241 elements each. The network used for this problem is way more complicated than in the previous case because of the size of each input. A 2-50-1 network with tanh activation functions in the hidden layer and linear activation function in the output layer is trained. Different number of hidden neurons are tried, from 5 to 100, and 50 is the final choice because it leads to good results. As before, different training algorithms are considered and each

neural net is trained only 20 times because the execution time for each training process is larger.

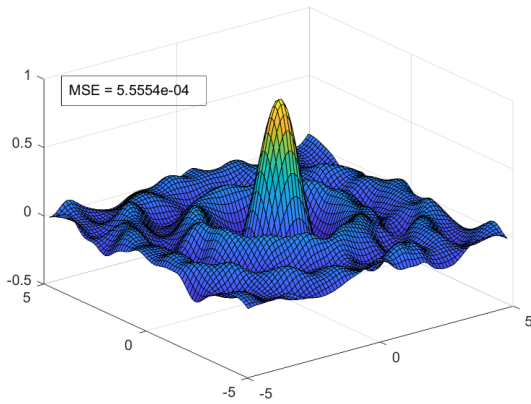
Figure 14 shows the difference in fit between each algorithm and the corresponding MSE. Visually, it can be seen that there are a lot of differences in the fit, especially in the Quasi-Newton and the scaled conjugate gradient (plots 14c and 14d). In both, the approximated function presents a lot of noise around the highest peak. In regard to the MSE, the algorithm that presents the lowest value is the Levenberg-Marquardt and the highest value is the scaled conjugate gradient. As in the case with just 1 input, the LM algorithm is preferred.



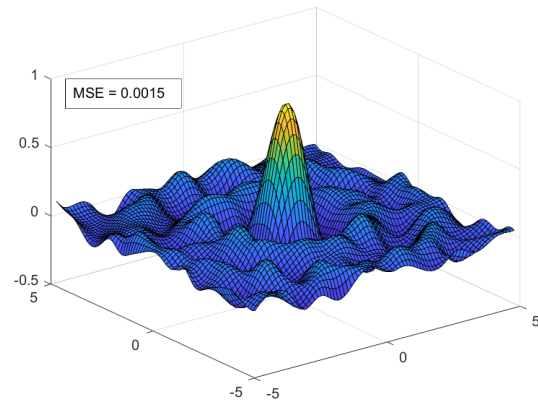
(a) Levenberg-Marquardt



(b) Bayesian regularization



(c) BFGS Quasi-Newton



(d) Scaled Conjugate Gradient

Figure 14: Comparison of neural nets using different training algorithms, when $m = 2$

3. $m = 5$

In this case, the input dataset is more difficult to get because it is necessary to get all possible combinations between the five input variables to build a grid in a sixth dimensional space. This process is conducted in a similar way as in the previous case, for each variable a sequence of just 5 values between -5 and 5 is created and then each vector is repeated the corresponding number of times to get all combinations, resulting in vectors of 15625 elements.

In regard to the test set, the process is repeated by creating a sequence of four numbers between -4.75 and 4.75 , resulting in vectors with 4096 elements each. The network used

for this problem is way more complicated than in the previous case because of the size of each input. A 5-100-1 network with tansig activation functions in the hidden layer and linear activation function in the output layer is trained. Different number of hidden neurons are tried, from 50 to 150, and 100 is the final choice because it leads to a MSE smaller than 0.01 without being "too" complicated. As before, different training algorithms are considered and each neural net is trained only 10 times because the execution time for each training process is larger.

Since it is not possible to show graphically the true function nor the approximated function, Table 2 shows the difference in MSE for each neural net. Surprisingly, the training of all networks started with the same initial weights and biases and the optimal solution was the same, or at least they got the same minimum, for the last three networks. However, the Quasi-Newton algorithm was the faster one and therefore it is chosen.

	MSE
Levenberg-Marquardt	0.0027
Bayesian Regularization	0.0018
BFGS Quasi-Newton	0.0018
Scaled Conjugate Gradient	0.0018

Table 2: Comparison of MSE between neural networks trained with different algorithms

To conclude, several neural networks were fitted to obtain a good approximation of the *sinc* function in 2, 3 and 6 dimensions. In the first case, 5 hidden neurons and the Levenberg-Marquardt algorithm were needed to obtain a good approximation of the true function in the interval $[-4.9, 4.9]$ of the input variable. In the second case, 50 hidden neurons and the LM algorithm were needed to get a decent representation of the true function. And in the third case, 100 hidden neurons and the Quasi-Newton were needed to get a good MSE. It is important to mention that in the last scenario, three different algorithms resulted in the same MSE but the faster one was preferred. In regard to the curse of dimensionality, it is clear that the number of hidden neuron increases as the dimension of the input increases, however the final number of parameters needed to obtain a good approximation is not as large as it would be needed if a polynomial was considered.

2 Exercise Session 2

2.1 Santa Fe laser data - time-series prediction

In time series problems, the use of the neural networks is different. The training process is the same as in every other kind of problem, feedforward mode, but the prediction process requires some changes. Once the training is finished, the neural network is used in an iterative way as a recurrent network, this is also called close loop, and it consists in using the prediction at some point t , to predict the value in $t + 1$.

The goal of this section is to predict the next 100 points from a chaotic laser dataset. To do this, different Non-linear Autoregressive (NAR) neural networks with different specifications are trained and then compared to get the best possible prediction. This comparison is conducted visually and by measuring the performance on a test set via the Mean Square Error (MSE).

In regard to local optima, it is well known that each time a neural net is trained, different initial weights and biases, and different divisions of the data into training, validation and test

sets are considered. Thus, each neural network is trained 100 times, and the net with the lowest MSE is kept for comparison. As before, in order to get reproducible results the seed in the random number generator is set to 0689432.

1. Lag order

Four different lags are tried in this comparison: 2, 7, 14 and 21. The lag 2 is chosen because it is usual in time series problems, is like analyzing the series in differences. Lag 7, 14 and 21 are chosen because every 7 time periods there is a peak in the series, this is similar to analyze the data in seasonal differences. To compare the performance of these neural networks, all characteristics are constant: 10 hidden units, the training algorithm is the Levenberg-Marquardt, the number of validation checks is 6 (default), no regularization is considered, and the selection of initial weights, which is controlled by the seed in the random number generator.

In Figure 15 the quality of the prediction of each neural net is presented. Visually, the predictions of all nets are close to the true value up to time 60 approximately. The main difference is noticeable after time 60, where some of them are far from the true value (Figure 15a) and some are close (Figure 15c). In regard to the MSE, the best prediction is given by the neural network that has a lag of 14, which means that the time series is persistent because the present value depends on past values up to 14 time periods.

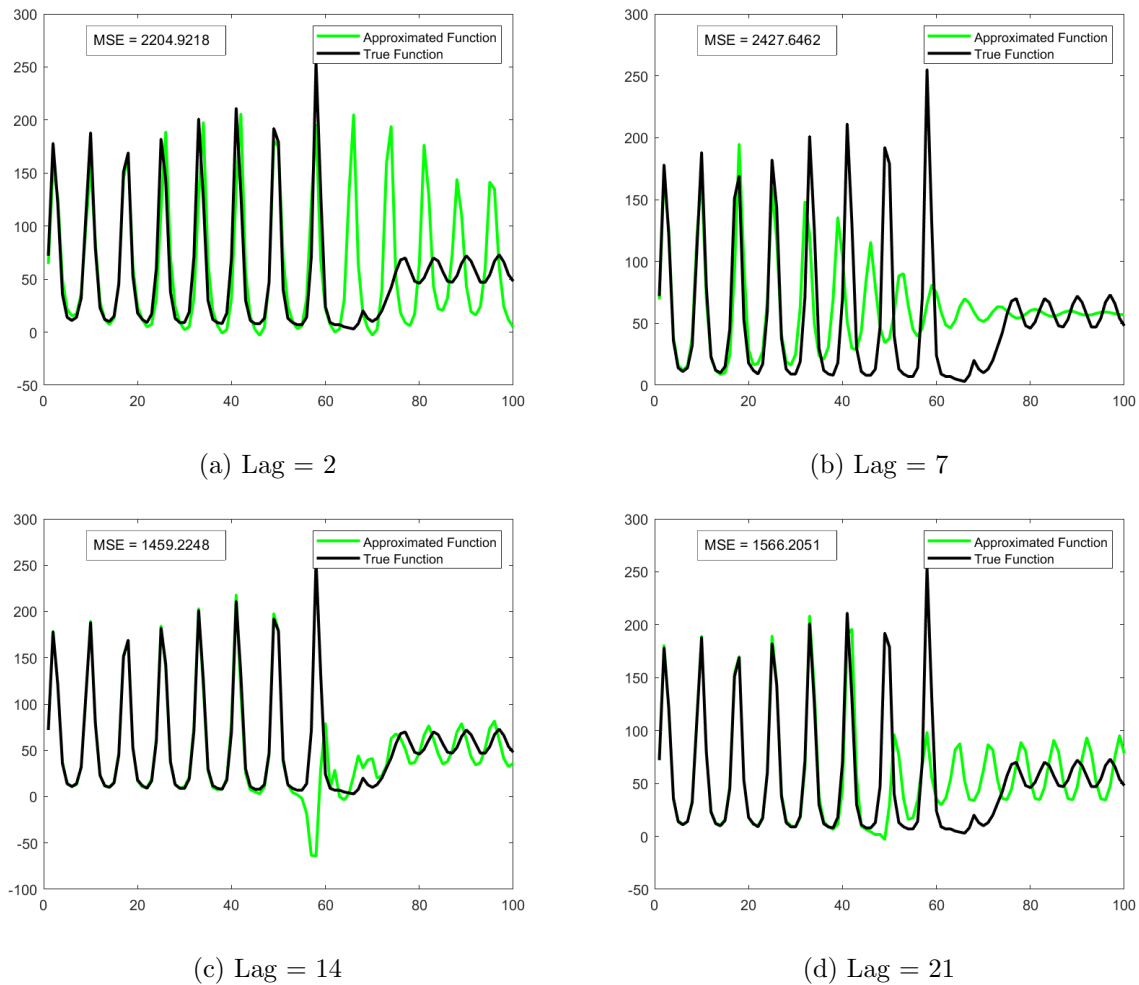
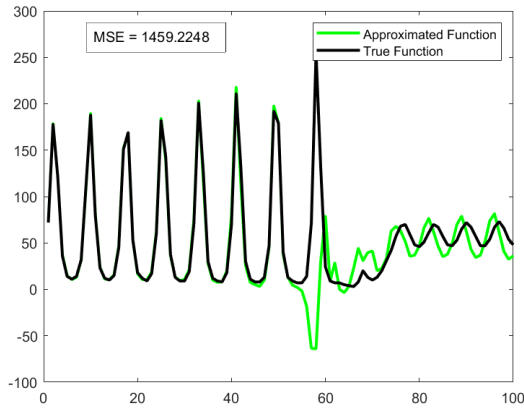


Figure 15: Comparison of neural nets using different lags

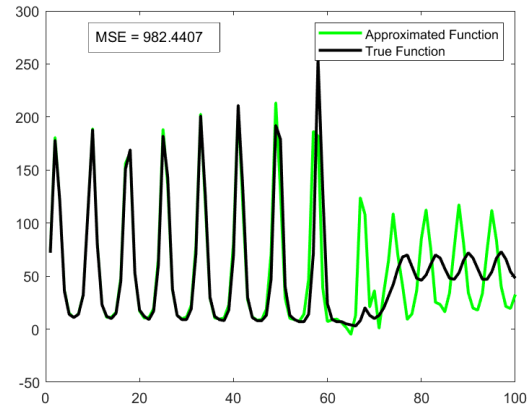
2. Number of hidden neurons

In this case four different scenarios are considered: 10, 20, 30 and 50 neurons. As before, the remaining characteristics are held constant, and the lag considered is 14. It is important to mention that more than 4 scenarios were considered in the training process (8, 9, 11, 12, 15 neurons), however only these 4 cases are shown because they produce an overview of the overall prediction trend.

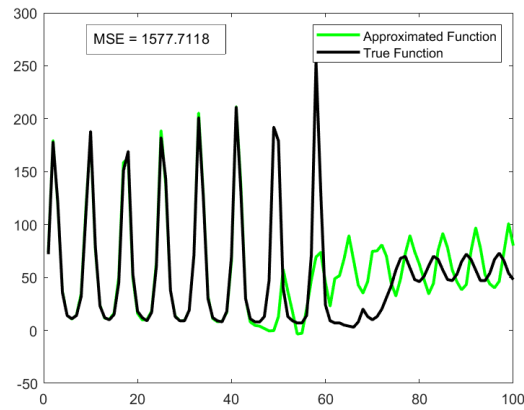
Figure 16 shows the prediction of each neural net and its corresponding performance measure. In terms of MSE, the best prediction is given by the net that has 50 hidden units followed by the one that has 20 units. Visually, the net that has only 10 hidden neurons is closer to the true value from the 70th period onward, but its prediction in the last peak, in the interval 50 to 62, is not very accurate. Something similar happens to the net with 30 neurons, the prediction in the last peak is poor. So, even though the prediction of the net with 50 units is not as close as, say, the one that has 10 units in the last part of the time window, it is preferred because it is capable of predicting correctly the last peak.



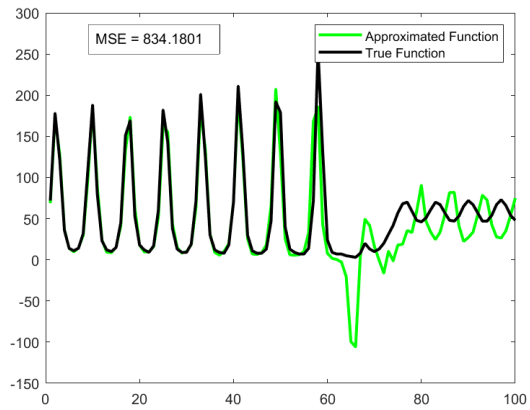
(a) 10 hidden neurons



(b) 20 hidden neurons



(c) 30 hidden neurons



(d) 50 hidden neurons

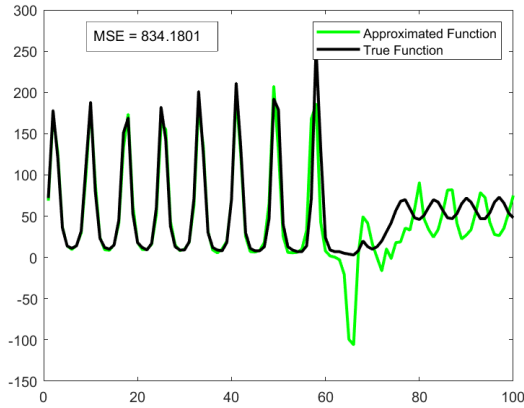
Figure 16: Comparison of neural nets using different number of hidden units

3. Training Algorithms

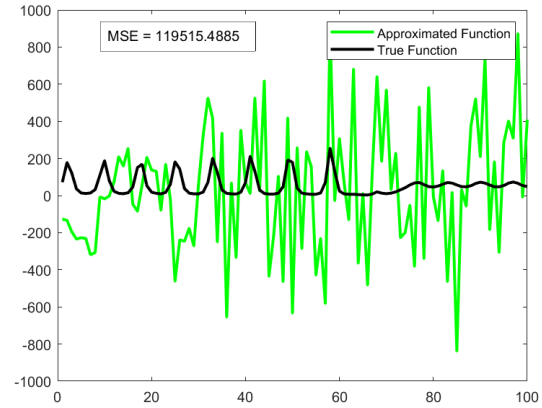
In this case four different algorithms are compared: Levenberg-Marquardt (default), Back-propagation, Conjugate Gradient and Quasi-Newton. As before, the remaining characteristics are held constant, the lag is 14 and the number of hidden units is 50. It is important

to mention that backpropagation is a term that sometimes is used to refer to gradient descent algorithm, so in MATLAB the gradient descent algorithm was conducted.

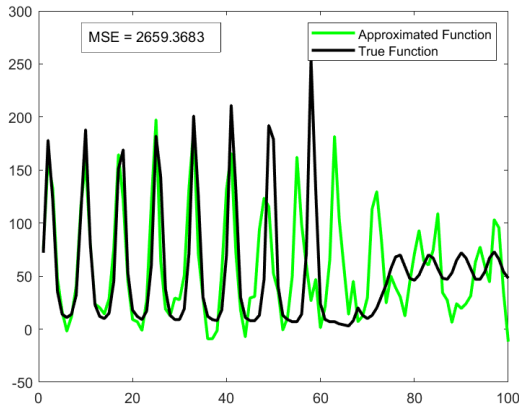
Figure 17 shows the difference in prediction between each algorithm. Visually, the results are surprising, at least for the gradient descent algorithm (plot 17b). In this case, the training process was different than in the other algorithms, because the MSE increased at each iteration instead of decrease, as it is supposed to. As a result, the prediction could not be far from the true value. In regard to the scaled conjugated gradient and the Quasi-Newton algorithms, the quality of the prediction is not better than in LM, and both took more time to get the solution.



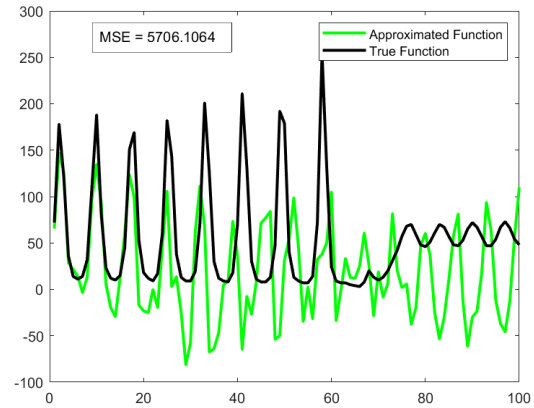
(a) Levenberg-Marquardt



(b) Gradient Descent (Backpropagation)



(c) Scaled Conjugate Gradient



(d) BFGS Quasi-Newton

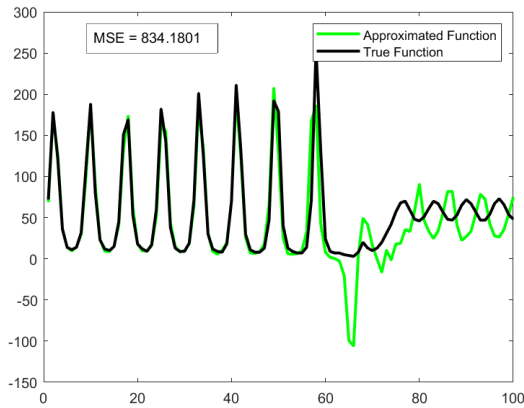
Figure 17: Comparison of neural nets using different training algorithms

4. Regularization Constant

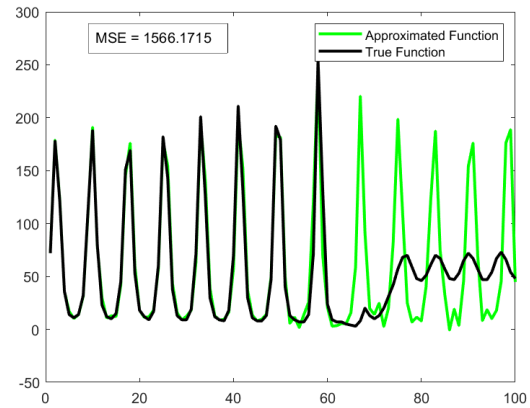
In this case different regularization constants are considered: from 0, where no regularization is conducted, to 0.5 where in the optimization process the original cost function and the weight decay term have the same importance. Four different constants are considered, namely 0, 0.1, 0.25 and 0.5. As before, the remaining characteristics are held constant, the lag is 14, the number of hidden units is 50 and the training algorithm is the Levenberg-Marquardt.

In Figure 18 the difference in prediction between each regularization constant is shown. Visually, the prediction is practically the same for all nets up to the 40th time period. After

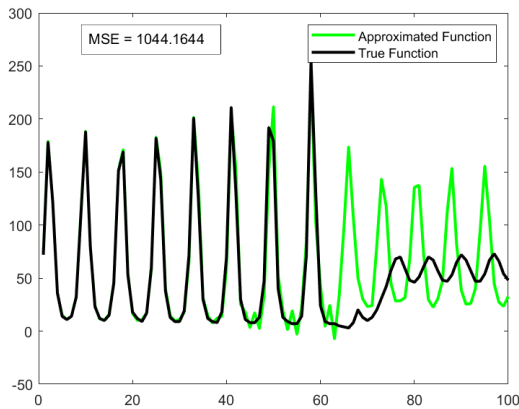
this point the difference in quality is mostly explained in the last 30 time periods. The network that presents the lowest MSE is the one that has the regularization constant equal to 0.5 (Plot 18d) and the largest is the one that has a regularization constant equal to 0.1 (Plot 18b). Notice, that even though the last net present the lowest MSE, the prediction from time 50 to 70 is not very good. In regard to the first neural net, with no regularization, it is possible to see that the prediction stick to the true value in more time points that in the last network. The reason why its corresponding MSE is larger is because of the prediction in the interval 60 to 70. For this reason, the best neural net to make predictions the first one with no regularization.



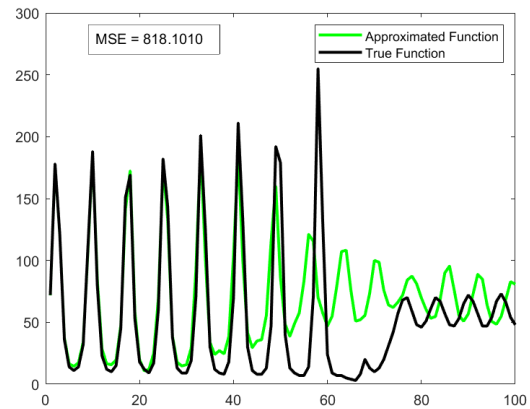
(a) Regularization Constant = 0



(b) Regularization Constant = 0.1



(c) Regularization Constant = 0.25



(d) Regularization Constant = 0.5

Figure 18: Comparison of neural nets using different regularization constant

In conclusion, different characteristics were compared and some interesting results were found in regard to the quality of the prediction. First, the lag needs to be chosen carefully because the quality of the prediction depends on this value. Second, it looks like the gradient descent algorithm is not applicable to this kind of problem because in the training process, the MSE increased at each iteration instead of decrease, as it is expected. Third, the best prediction is given by a neural net with 50 hidden neurons, using the LM algorithm and no regularization on a dataset with lag 14. It is important to mention that in all of these cases the correlogram of the residuals were checked but, none of them were white noise, so the results can be improved.

2.2 Alphabet recognition

This is an example of how to train a neural network to recognize alphabet characters. The input is a matrix of 26 columns (one for each letter) and 35 rows that correspond to a 5x7 bitmap. The values of this matrix are either 0 or 1, indicating what cell should be colored, so to speak. The output is an identity matrix of 26x26 that indicates which is the final letter. Figure 19 shows how the input data look like.

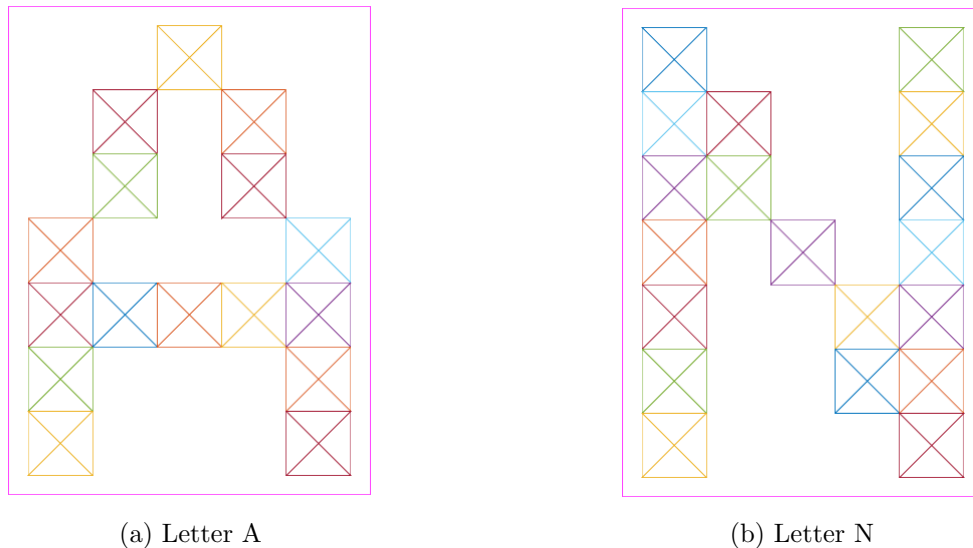


Figure 19: Example of how the input looks like

A neural network with one hidden layer and 25 hidden neurons is used to tackle this problem. The characteristics of the net are the default in the software, meaning that it is trained with the Levenberg-Marquardt algorithm, the dataset is randomly divided in training, validation and test sets, and the algorithm stops when the network is no longer likely to improve on the training or validation sets.

Since the neural net is trained only once in the example and given that a seed is set in the beginning of the code, it is certain to say that the training of the neural net stopped because the gradient of the cost function (MSE in this case) was smaller than $1e - 7$, which is the default value.

In addition, another neural net with the same characteristics is trained with noisy data. To do so, 30 copies of each letter are created and random values are added to each of them. Figure 20 shows similar examples with random noise.

The training of this second neural network stopped because the validation error started to increase. By default, when the validation error increases for 6 continuous iterations the training stops to prevent overfitting.

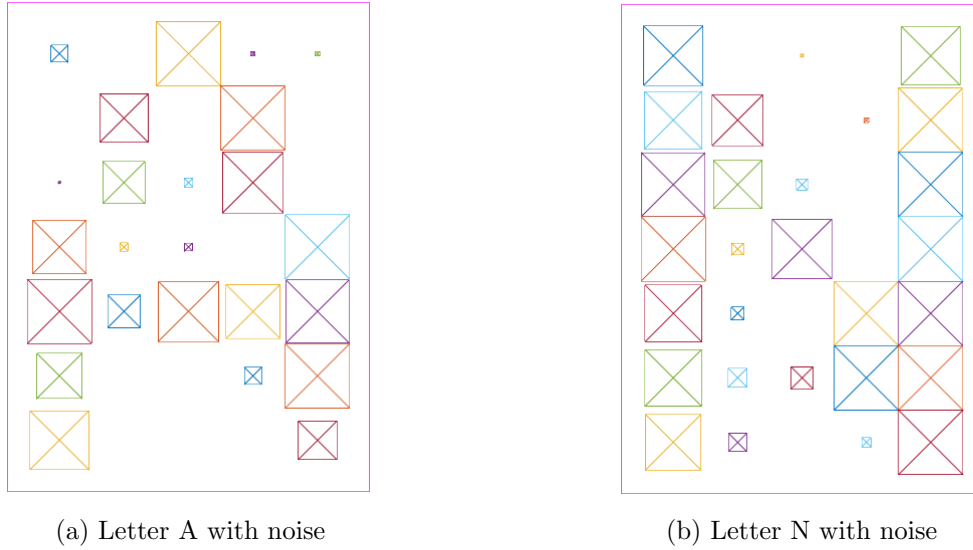


Figure 20: Example of how the input with noise looks like

Once both neural nets are trained, their performance is tested by predicting different test sets. To do so, 21 different noisy test datasets are used and in each of them the amount of noise varies by changing the standard deviation of the noise (from 0 to 1). This performance is measured as a percentage of recognition errors, and it is calculated by first transforming the maximum predicted value in each column of the output pattern to 1 and the remaining values to zero, so each column has only a single 1 which corresponds to the classified letter². Then the absolute difference between the predicted and the true values is computed, resulting in a matrix of zeros and ones. In this matrix when the classification is correct, the whole column is filled with zeros, whereas each time there is a misclassification, there are two 1s in the respective column. Thus, to compute the final percentage of misclassifications, all of the values in this matrix are summed and then divided by two times the number of columns to classify (2×780).

As a result, the second neural net presents the lowest percentage of recognition errors, meaning that it is the best to predict new data (see Figure 21). This is an expected result because the second neural net was trained with noisy data whereas the first one was not.

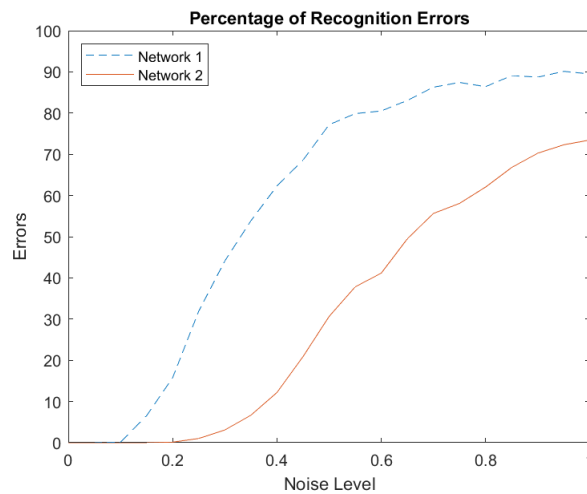


Figure 21: Comparison of both neural nets in terms of percentage of recognition errors.

²In this example, there are 780 (30×26) columns in the output pattern, that corresponds to 30 noisy copies of each letter and 26 possible letters

2.3 Breast Cancer Wisconsin - classification problem

The main difference between a function approximation and a classification problem is the architecture. Usually, the activation function used in the output layer is the softmax transfer function, which is a function that returns values in the range $[0,1]$ inclusive. This value can be interpreted as the probability of an observation to belong to a specific class.

The purpose of this section is to create a nonlinear classification rule to decide whether a patient has breast cancer or not. To do so, several neural networks with different specifications are trained and then compared to get the best solution possible. This comparison is conducted measuring the performance on a test set via cross-entropy (CE), the percentage of misclassifications and the ROC curve.

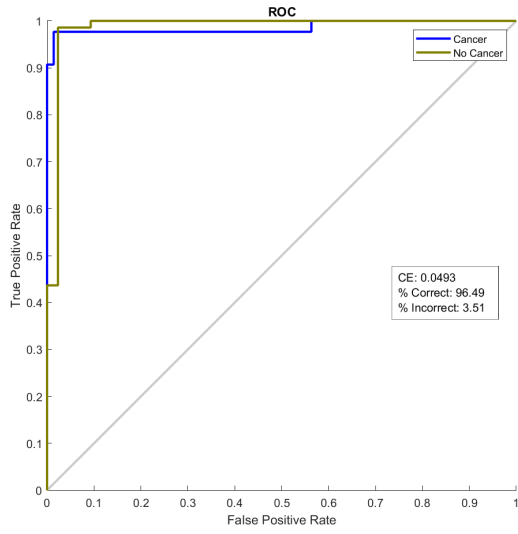
The original dataset, which has information of 569 patients, is divided in two subsets: a training set with 80% of the patients (455 observations) and a test set with the remaining 20%, which corresponds to 114 patients. The purpose of this division is to be sure that a test set in one neural network is not going to be used as a training set in other neural net.

In regard to local optima, it is well known that each time a neural net is trained, different initial weights and biases, and different divisions of the data into training, validation and test sets are considered. Thus, each neural network is trained 100 times and the net with the lowest cross-entropy is kept for comparison. As before, in order to get reproducible results the seed in the random number generator is set to 0689432.

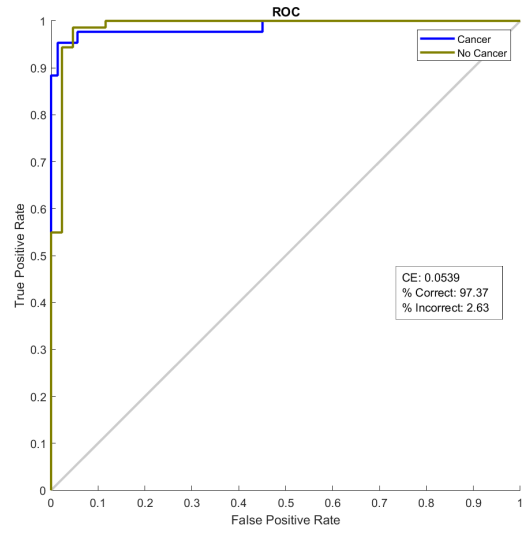
1. *Number of hidden neurons*

Four different architectures are considered: 2, 5, 10 and 15 hidden neurons. To compare the performance of these nets, the remaining characteristics are held constant, such as the training algorithm, which is the conjugate gradient, the number of validation checks, the regularization constant, which is equal to zero and the selection of initial weights (controlled by the seed in the random number generator).

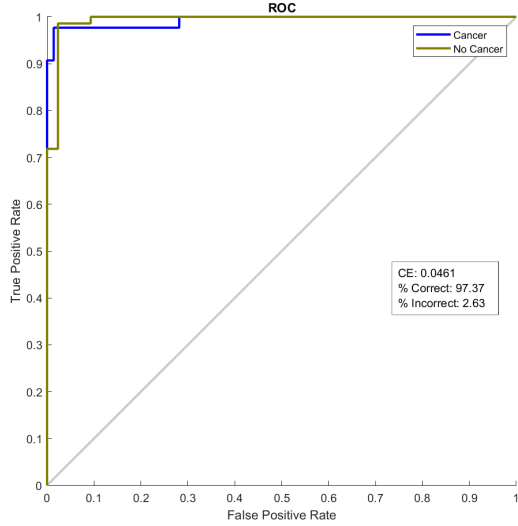
Figure 22 shows the ROC curve with the percentage of correct classifications and misclassification and the respective cross-entropy, for all neural networks on the test set. The first conclusion that can be drawn is that all nets give good results, because the area under the ROC curve is almost 1, which is the perfect scenario. The percentage of correct and incorrect misclassifications is the same for the nets with 5, 10 and 15 hidden neurons. The net with the smallest cross-entropy is the one with 10 units. These results lead to conclude that the neural network with 5 hidden neurons performs best than the others because it produces solid results with less number of parameters (is the most parsimonious).



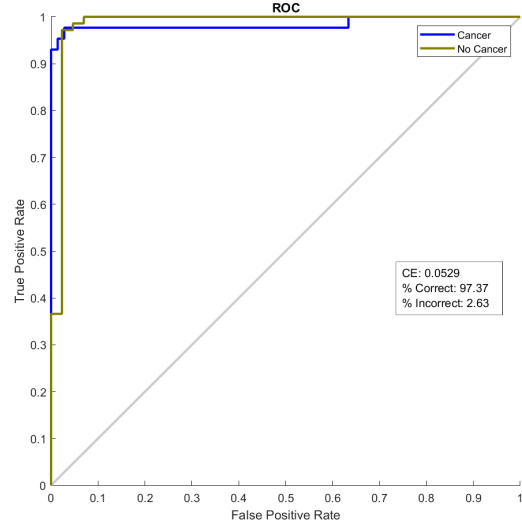
(a) 2 hidden neurons



(b) 5 hidden neurons



(c) 10 hidden neurons



(d) 15 hidden neurons

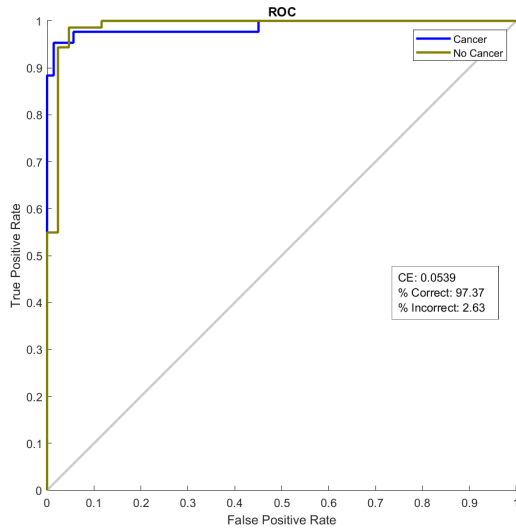
Figure 22: Comparison of neural nets using different number of hidden neurons

2. Training algorithms

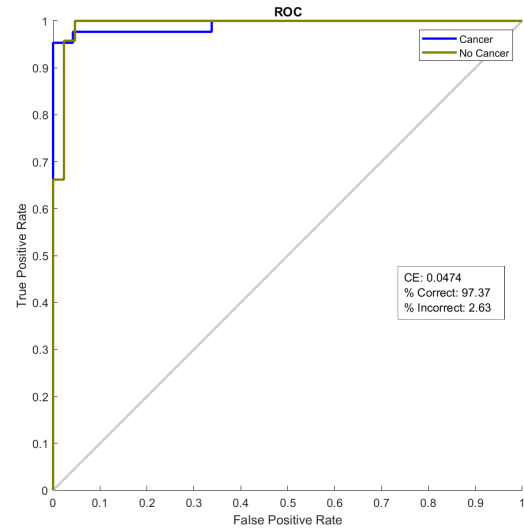
In this case four different training algorithms are compared: backpropagation, Levenberg-Marquardt, quasi-Newton and conjugate gradient, which is the default in MATLAB for pattern recognition. As before, the remaining characteristics are held constant and the number of hidden units considered is 5. It is important to mention that backpropagation is a term that sometimes is used to refer to gradient descent algorithm, so in MATLAB the gradient descent algorithm was conducted.

Figure 23 shows the difference in classification between each algorithm. Visually, all nets seem to provide similar results because the area under the ROC curve is almost 1, which is the desirable output. In regard to the percentage of classifications, the conjugate gradient and the Levenberg-Marquardt present the highest value (97.37%). On the other

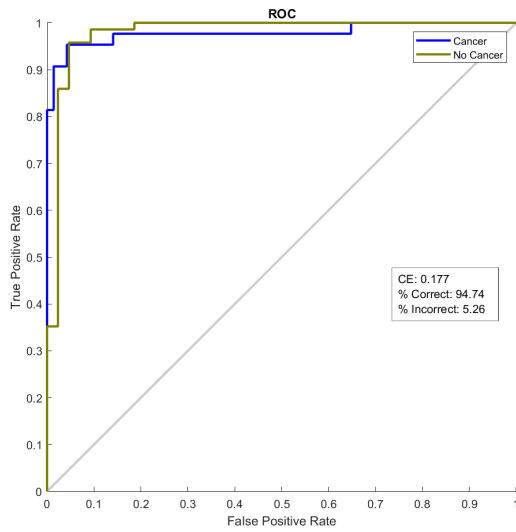
hand, the gradient descent algorithm has the highest percentage of misclassifications (5.26%), although this value is acceptable in a problem like this one. These results lead to conclude that the algorithm with the best performance and speed is the conjugate gradient, an expected result given the simulations presented in Matlab (2018).



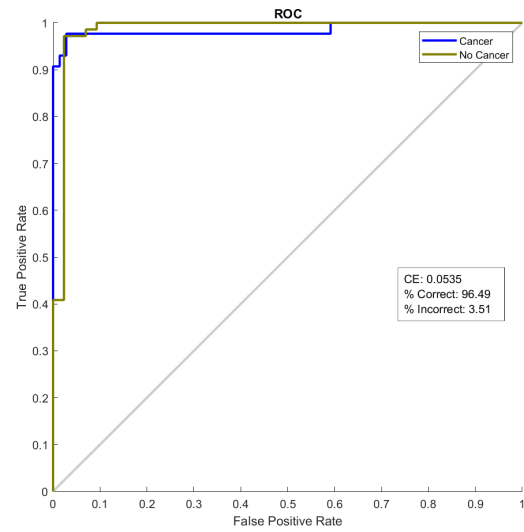
(a) Scaled Conjugate Gradient



(b) Levenberg-Marquardt



(c) Gradient Descent (Backpropagation)



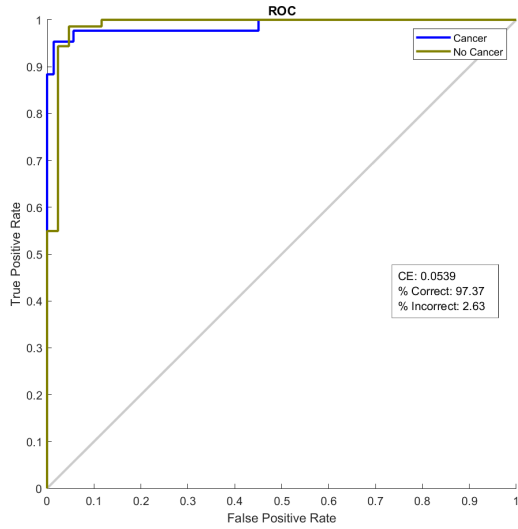
(d) BFGS Quasi-Newton

Figure 23: Comparison of neural nets using different training algorithms

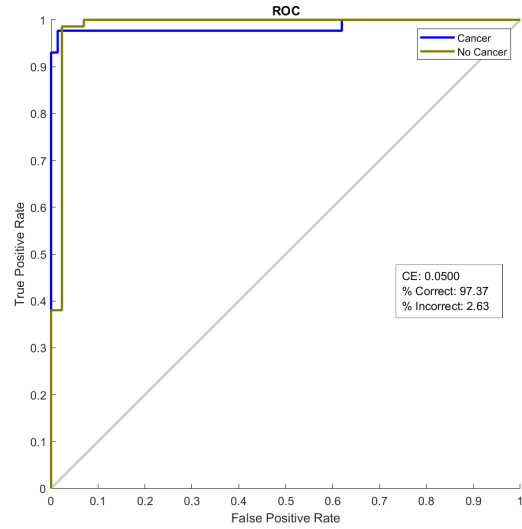
3. Regularization constant

The percentage of good classifications in the previous sections are really good for a classification problem like this one. This may suggest that those neural networks can be overfitted. To avoid this problem, different regularization constants are considered: 0, where no regularization is conducted, 0.01, 0.1 and 0.2. As before, the remain characteristics are held constant, the number of hidden neurons is 5 and the training algorithm is the conjugated gradient.

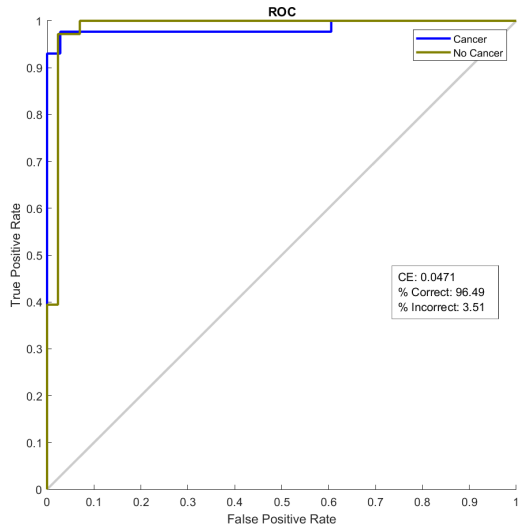
In Figure 24 the difference in performance for each regularization constant is presented. The first conclusion that can be drawn is that all nets give good results, because the area under the ROC curve is almost 1, which is the perfect scenario. The percentage of correct and incorrect misclassifications is the same for the nets with regularization constant equal to 0, 0.01 and 0.2, suggesting that the neural nets in the previous sections were not overfitted. These results lead to conclude that any regularization constant will provide solid results in terms of classification of new patients.



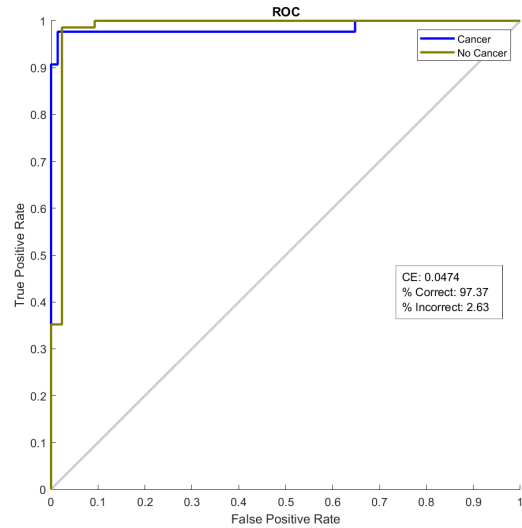
(a) Regularization constant = 0



(b) Regularization constant = 0.01



(c) Regularization constant = 0.1



(d) Regularization constant = 0.2

Figure 24: Comparison of neural nets using different regularization constants

To conclude, different characteristics of a neural network were investigated and compared. In this case of pattern recognition, in contrast to the previous case of function approximation, different parametrizations of the neural net give the same number of misclassification, which by the way is really small (just 3 cases over 114). The number of hidden neurons considered

was the smallest possible because a model with fewer parameters is easier to estimate and maintain in the long run. With respect to the training algorithm, the performance was equally good but the conjugate gradient was chosen because is faster. No regularization constant should be considered in this case because the performance of the neural net was not affected by it. Therefore, the best model to classify new patients is a net with 5 hidden neurons, using the conjugate gradient with no regularization.

3 Exercise Session 3

3.1 Dimensionality reduction by PCA analysis

In this section a bio-medical dataset is used to compare results in performance between a neural net trained with original observations and a neural net trained with the solution of a prior Principal Component Analysis (PCA). The input dataset consists in 21 spectral measurements of 264 blood samples and the idea is to predict three kinds of cholesterol for each blood sample: LDL, VLDL and HDL.

The first step is to conduct a PCA with the goal of having less (uncorrelated) input variables that explain a reasonable amount of the variance of the original dataset. To do so, each spectral measurement is standardized by subtracting its mean and dividing it by its standard deviation. In this way, all measurements can be compared between each other, and the solution of the PCA is easier to get.

Then the PCA analysis is performed on the standardized variables. The criteria to select the number of principal components is to retain all of those whose contribution to the total variation is more than 0.001. As a result, 4 principal components are retained, implying that this solution at least explains 98.30% of the total variability obtained in the original input dataset³.

Having done this, the next step is to train different neural networks. The first neural network considered is trained with the original dataset, which has 264 blood samples with 21 spectral measurements. This dataset is then divided between training, validation and test set, with 132, 66 and 66 blood samples respectively. The architecture of this net is 21-5-3, which means that there are 21 input measurements, just only hidden layer with 5 hidden neurons, and the output layer with 3 measurements (LDL, VLDL and HDL). The training algorithm is Levenberg-Marquardt and no regularization is conducted. To avoid bad local optima solutions, the network is trained 100 times to ensure that different initial weights and biases are considered. However, in order to get reproducible results the seed in the random number generator is set to 0689432.

In a similar fashion, a second neural network is also trained with the original dataset. All of the considerations mentioned above remain the same in this case, but the training algorithm. In this neural network a Bayesian regularization algorithm is used with the purpose of getting a solution that generalizes well. This algorithm updates the weights and biases according to the Levenberg-Marquardt algorithm, but the cost function (also called error or energy function) is modified by having now a linear combination of squared errors and weights. This algorithm modifies this linear combination so that at the end a network with good generalization results.

Additionally, a third neural network is trained with the solution of the PCA analysis. In this case, the input dataset does not have 21 spectral measurements, but 4 principal components that explain approximately 98% of the total variation. This dataset is also divided in training, validation and test sets as before and the remain characteristics and the architecture is the same as in the second neural network, namely it has 5 hidden neurons and it uses a Bayesian regularization algorithm.

³Assuming that all of the remain principal components contribute 0.001 to the total variation, the solution with just 4 components explains at least $1 - (17 * 0.001)$, with 17 principal components left out of the solution.

Table 3 shows the performance of each neural network in terms of Mean Squared Error (MSE) calculated on the training and the test set. Surprisingly, the first neural network (LM algorithm) is the one that has the best performance in both training and test datasets. It seems that the Bayesian regularization algorithm does not improve the quality of prediction in comparison with the LM algorithm, so no regularization is needed in this case. It is important to mention though, the difference in test MSE between the first and the second neural networks is not large, and therefore it could happen that if other dataset were used to train both networks, the second net could perform better than the first one. As it was mentioned before, bayesian regularization uses the LM algorithm to update weights and biases, so the performance of the second net can be influenced by the data at hand.

Network	Training	Test
LM algorithm	0.2684	0.2446
Bayesian Regularization	0.3055	0.2487
PCA	0.3442	0.2659

Table 3: Comparison of MSE between neural networks on training and test datasets

In regard to the third neural network, trained with 4 principal components instead of 21 original measurements, the performance in both training and test MSE is worse than Bayesian regularization on the original data. In this sense and given that there are not many inputs, the **best solution is the second neural network**. Although there are other motivations to choose the neural net with original data: first, it is clear that PCA is trained way more fast that the original dataset (because of the number of inputs), however in this case is not enough to be faster because the dataset is not large and the computing time difference is not huge. Second, when a PCA solution is considered to train a network, it is important to have in mind that each time a new blood sample comes into the dataset, the solution may change in the number of components retained or in different projections resulting of an update of the analysis. So, if the neural network is trained in online mode (weights and biases are updated each time new samples appear), the PCA solution will require an extra step.

3.2 Input selection by Automatic Relevance Determination (ARD)

In this section, the UCI ionosphere dataset is considered to conduct input selection by Automatic Relevance Determination. The dataset consists of 33 input variables and 351 observations in each of them. The output is a 2-class variable with values -1 or 1.

In this methodology, the weights and biases are considered as random variables and therefore they have a corresponding distribution. A different Gaussian prior distribution with a separate hyperparameter is considered, corresponding to each group of weights for each input variable. These hyperparameters are all set to 0.01, but will be re-estimated after the net is trained. The architecture of the neural network is 33-2-1 with tansig activation function in the hidden layer and logistic activation function in the output layer. The net is trained 5 times using the scaled conjugate gradient algorithm and at the end of each cycle the hyperparameters are re-estimated. As before, in order to get reproducible results the seed in the random number generator is set to 0689432.

The first column of Table 4 shows the hyperparameters for the corresponding input to hidden units weights. Each of these values represents an inverse variance, so the larger the value, the smaller the posterior variance will be. Variables 8, 15, 17, 19 and 28, present the largest values and therefore they are considered as less relevant. This is confirmed by looking how

the corresponding weights are close to zero (columns 2 and 3). On the other hand, variables 1, 9, 22 and 32 are considered as the most relevant because they present the largest inverse variance values. In Table 4, the less relevant are in red and the most relevant are in blue.

Inputs	First Training			Second training		
	Hyperpar.	Neuron 1	Neuron 2	Hyperpar.	Neuron 1	Neuron 2
1	0.008	1.721	−13.17	0.008	2.298	15.892
2	0.047	1.345	−6.379	0.038	−1.489	7.056
3	0.122	−3.005	−2.693	7.505	−0.213	−0.267
4	0.139	3.713	0.714	0.075	4.268	2.873
5	0.261	2.578	−0.933	3.241	0.551	0.506
6	0.09	4.332	−1.753	0.262	2.169	1.664
7	0.081	4.565	1.941	0.102	4.353	−0.726
8	167.914	0.005	0.018			
9	0.027	−1.937	−8.442	0.079	1.494	4.731
10	0.066	−5.3	1.383	0.02	0.062	−9.864
11	7.491	0.37	0.104	365.834	−0.005	−0.001
12	1.916	0.125	−0.925	0.34	−1.634	1.755
13	0.081	−3.154	−3.794	467.502	0.008	−0.005
14	0.537	0.807	−1.641	0.042	0.618	6.835
15	1769.108	0.000	0.003			
16	0.262	−0.068	−2.715	552.955	−0.003	0.009
17	811.9	0.005	−0.001			
18	8.795	−0.328	0.019	0.102	−1.202	4.207
19	292.539	0.014	−0.009			
20	8.552	−0.341	0.052	0.015	5.453	−10.176
21	0.301	−0.723	2.427	0.349	−0.835	−2.17
22	0.033	4.62	6.224	0.019	3.14	−9.606
23	3.559	−0.242	−0.677	0.171	−2.189	2.536
24	0.247	1.437	−2.427	0.053	−0.722	6.065
25	3.016	−0.208	−0.734	0.355	−0.047	2.317
26	0.203	−2.974	−0.101	0.01	−8.555	11.59
27	8.866	0.144	−0.224	0.028	1.07	−8.21
28	624.25	−0.006	0.005			
29	0.066	3.116	−4.527	0.017	−0.273	10.73
30	0.102	−2.381	−3.708	0.162	−1.132	3.252
31	4.301	0.464	0.375	17.741	0.035	−0.103
32	0.036	3.748	6.474	0.032	2.071	−7.549
33	0.061	−5.616	0.958	0.268	−0.989	−2.475

Table 4: Hyperparamaters and corresponding weights to neuron 1 and 2 in each neural network

The neural net is trained a second time under the same considerations mentioned before with the only difference that variables 8, 15, 17, 19 and 28 are removed from the analysis. Columns 4, 5, 6 in Table 4 shows the resulting hyperparameters and their corresponding weights. In this case, variables 11, 13, 16 and 31 are now considered as less relevant and variables 1, 20, 26 and 29 are the most relevant. The fact that there are other variables that are considered as less (more) relevant, can be explained by a possible correlation that exists between these variables and the ones removed. Having removed them, reveal new correlations between these variables and the ones that are now present in the dataset.

References

Matlab (2018). Choose a multilayer neural network training function. <https://nl.mathworks.com/help/deeplearning/ug/choose-a-multilayer-neural-network-training-function.html>. Last accessed on 4-1-2019.

4 Appendix

• Script Exercise Session 1

```
1 % Function Approximation (noiseless case)
2 nnd11gn
3
4 %% The role of the hidden layer and output layer
5 clear
6 % 2
7 x = linspace(0, 1, 21);
8 % 3
9 y = -cos(0.8 * pi * x);
10 figure; plot(x, y, '-'); xlabel('x'); ylabel('y = -cos(0.8 \pi x)');
11
12 % 4
13 % To avoid local optima initialize the NN several times
14 % Set seed value to obtain reproducible results
15 setdemorandstream(0689432) % This number is my student number
16 net = fitnet(2);
17 net = configure(net, x, y);
18 net.inputs{1}.processFcns = {};
19 net.outputs{2}.processFcns = {};
20 n_start = 100;
21 NN = cell(1, n_start);
22 y_est = zeros(n_start, length(x));
23 perfs = zeros(1, n_start);
24 for i = 1:n_start
25     fprintf('Training %d/%d\n', i, n_start);
26     NN{i} = train(net, x, y);
27     y_est(i, :) = NN{i}(x);
28     perfs(i) = mse(net, y_est(i,:), y);
29 end
30 [m,i] = min(perfs)
31 % 5.
32 [biases, weights] = hidden_layer_weights(NN{i});
33 net_act_fun = hidden_layer_transfer_function(NN{i});
34 input_act = net_act_fun(weights * x + biases);
35 % 6.
36 figure; plot(x, y, '-'); xlabel('x'); ylabel('y = -cos(0.8 \pi x)');
37 hold on;
38 plot(x, input_act(1,:), 'r-'); plot(x, input_act(2,:), 'g-');
39 hold off;
40 legend({'y', 'x^1', 'x^2'}, 'Location', 'north')
41 % 7.
42 [biases_out, weights_out] = output_layer_weights(NN{i});
43 act_fun_out = output_layer_transfer_function(NN{i});
44
45 % Calculation of output
46 output = act_fun_out(weights_out * net_act_fun(weights * x + biases)) + biases_out;
47 figure; plot(x, y, '-'); xlabel('x'); ylabel('y = -cos(0.8 \pi x)');
48 hold on;
49 plot(x, output, 'p-');
50 hold off;
51 legend({'y', 'y-est'}, 'Location', 'north')
52
53 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
54 %% Function Approximation (noisy case)
55 clear
56 %% Comparison of dataset size
57 % Size of the training data set (e.g. 30, 150, 1200 data points).
58 clear
59 setdemorandstream(0689432);
60
61 tr_x1 = linspace(-1, 1, 30);
62 tr_x2 = linspace(-1, 1, 150);
63 tr_x3 = linspace(-1, 1, 1200);
64 val_x1 = linspace(-0.9, 0.9, 30);
65 val_x2 = linspace(-0.9, 0.9, 150);
66 val_x3 = linspace(-0.9, 0.9, 1200);
67
68 % The remain characteristics are held constant.
69 % Amount of noise = standard deviation of 0.6
70 tr_y1 = cos(2 * pi * tr_x1) + 0.6 * randn(size(tr_x1));
71 tr_y2 = cos(2 * pi * tr_x2) + 0.6 * randn(size(tr_x2));
72 tr_y3 = cos(2 * pi * tr_x3) + 0.6 * randn(size(tr_x3));
73 val_y1 = cos(2 * pi * val_x1) + 0.6 * randn(size(val_x1));
74 val_y2 = cos(2 * pi * val_x2) + 0.6 * randn(size(val_x2));
75 val_y3 = cos(2 * pi * val_x3) + 0.6 * randn(size(val_x3));
76 [sum(tr_y1);sum(tr_y2);sum(tr_y3);sum(val_y1);sum(val_y2);sum(val_y3)]
77
78 % Training algorithm= Levenberg-Marquardt
79 % Hidden neurons = 5
80 % Early stopping, no regularization
81 net1 = fitnet(5, 'trainlm');
82 net2 = fitnet(5, 'trainlm');
83 net3 = fitnet(5, 'trainlm');
84
85 % Training net
86 setdemorandstream(0689432);
87 n_start = 100;
88 NN1 = cell(1, n_start);
89 perfs1 = zeros(1, n_start);
90 val_t1 = zeros(n_start, length(val_x1));
91 for i = 1:n_start
```

```

92     fprintf('Training %d/%d\n', i, n_start);
93     NN1{i} = train(net1, tr.x1, tr.y1);
94     val.t1(i, :) = NN1{i}(val.x1);
95     perfs1(i) = mse(net1, val.t1(i, :), val.y1);
96 end
97 [m1,i1] = min(perfs1)
98
99 % plot results (1.3.1)
100 figure;
101 plot(tr.x1, tr.y1, 'o','MarkerSize',2,'MarkerFaceColor','blue',...
102      'MarkerEdgeColor','blue');
103 hold on;
104 plot(val.x1, val.y1, 'o','MarkerSize',2,'MarkerFaceColor','red',...
105      'MarkerEdgeColor','red');
106 plot(val.x1, val.t1(i1,:), '-', 'LineWidth',2,'Color','green');
107 plot(tr.x1, cos(2 * pi * tr.x1), 'g-', 'LineWidth',2,'Color','black');
108 hold off;
109 legend('Training Set','Validation Set','Approximated Function','True Function');
110 dim = [0.2 0.6 0.3 0.3];
111 annotation('textbox',dim,'String','MSE = 0.4275','FitBoxToText','on');
112
113 % Training net2
114 setdemorandstream(0689432);
115 n_start = 100;
116 NN2 = cell(1, n_start);
117 perfs2 = zeros(1, n_start);
118 val.t2 = zeros(n_start, length(val.x2));
119 for i = 1:n_start
120     fprintf('Training %d/%d\n', i, n_start);
121     NN2{i} = train(net2, tr.x2, tr.y2);
122     val.t2(i, :) = NN2{i}(val.x2);
123     perfs2(i) = mse(net2, val.t2(i, :), val.y2);
124 end
125 [m2,i2] = min(perfs2)
126
127 % plot results (1.3.2)
128 figure;
129 plot(tr.x2, tr.y2, 'o','MarkerSize',2,'MarkerFaceColor','blue',...
130      'MarkerEdgeColor','blue');
131 hold on;
132 plot(val.x2, val.y2, 'o','MarkerSize',2,'MarkerFaceColor','red',...
133      'MarkerEdgeColor','red');
134 plot(val.x2, val.t2(i2,:), '-', 'LineWidth',2,'Color','green');
135 plot(tr.x2, cos(2 * pi * tr.x2), 'g-', 'LineWidth',2,'Color','black');
136 hold off;
137 legend('Training Set','Validation Set','Approximated Function','True Function');
138 dim = [0.2 0.6 0.3 0.3];
139 annotation('textbox',dim,'String','MSE = 0.4160','FitBoxToText','on');
140
141 % Training net3
142 setdemorandstream(0689432);
143 n_start = 100;
144 NN3 = cell(1, n_start);
145 perfs3 = zeros(1, n_start);
146 val.t3 = zeros(n_start, length(val.x3));
147 for i = 1:n_start
148     fprintf('Training %d/%d\n', i, n_start);
149     NN3{i} = train(net3, tr.x3, tr.y3);
150     val.t3(i, :) = NN3{i}(val.x3);
151     perfs3(i) = mse(net3, val.t3(i, :), val.y3);
152 end
153 [m3,i3] = min(perfs3)
154
155 % plot results (1.3.3)
156 figure;
157 plot(tr.x3, tr.y3, 'o','MarkerSize',2,'MarkerFaceColor','blue',...
158      'MarkerEdgeColor','blue');
159 hold on;
160 plot(val.x3, val.y3, 'o','MarkerSize',2,'MarkerFaceColor','red',...
161      'MarkerEdgeColor','red');
162 plot(val.x3, val.t3(i3,:), '-', 'LineWidth',2,'Color','green');
163 plot(tr.x3, cos(2 * pi * tr.x3), 'g-', 'LineWidth',2,'Color','black');
164 hold off;
165 legend('Training Set','Validation Set','Approximated Function','True Function');
166 dim = [0.2 0.6 0.3 0.3];
167 annotation('textbox',dim,'String','MSE = 0.3518','FitBoxToText','on');
168
169 %% Comparison of noise
170 clear
171 setdemorandstream(0689432);
172
173 % Size of the training data set is set to 150 observations.
174 tr.x = linspace(-1, 1, 150);
175 val.x = linspace(-0.9, 0.9, 150);
176
177 % Amount of noise (0.2, 0.6, 1.2)
178 tr.y1 = cos(2 * pi * tr.x) + 0.2 * randn(size(tr.x));
179 tr.y2 = cos(2 * pi * tr.x) + 0.6 * randn(size(tr.x));
180 tr.y3 = cos(2 * pi * tr.x) + 1.2 * randn(size(tr.x));
181 val.y1 = cos(2 * pi * val.x) + 0.2 * randn(size(val.x));
182 val.y2 = cos(2 * pi * val.x) + 0.6 * randn(size(val.x));
183 val.y3 = cos(2 * pi * val.x) + 1.2 * randn(size(val.x));
184 [sum(tr.y1);sum(tr.y2);sum(tr.y3);sum(val.y1);sum(val.y2);sum(val.y3)]
185 % The remain characteristics are held constant.
186 % Training algorithm= Levenberg-Marquardt
187 % Hidden neurons = 5

```

```

188 % Early stopping, no regularization
189 net1 = fitnet(5, 'trainlm');
190 net2 = fitnet(5, 'trainlm');
191 net3 = fitnet(5, 'trainlm');
192
193 % Training net1
194 setdemorandstream(0689432);
195 n_start = 100;
196 NN1 = cell(1, n_start);
197 perfs1 = zeros(1, n_start);
198 val_t1 = zeros(n_start, length(val.x));
199 for i = 1:n_start
200     fprintf('Training %d/%d\n', i, n_start);
201     NN1{i} = train(net1, tr.x, tr.y1);
202     val_t1(i, :) = NN1{i}(val.x);
203     perfs1(i) = mse(net1, val_t1(i, :), val.y1);
204 end
205 [m1,i1] = min(perfs1)
206
207 % plot results (1.3.4)
208 figure;
209 plot(tr.x, tr.y1, 'o', 'MarkerSize',2, 'MarkerFaceColor', 'blue', ...
210      'MarkerEdgeColor', 'blue');
211 hold on;
212 plot(val.x, val.y1, 'o', 'MarkerSize',2, 'MarkerFaceColor', 'red', ...
213      'MarkerEdgeColor', 'red');
214 plot(val.x, val_t1(i1,:), '-', 'LineWidth',2, 'Color', 'green');
215 plot(tr.x, cos(2 * pi * tr.x), 'g-', 'LineWidth',2, 'Color', 'black');
216 hold off;
217 legend('Training Set', 'Validation Set', 'Approximated Function', 'True Function');
218 dim = [0.2 0.6 0.3 0.3];
219 annotation('textbox',dim, 'String', 'MSE = 0.0372', 'FitBoxToText', 'on');
220
221 % Training net2
222 setdemorandstream(0689432);
223 n_start = 100;
224 NN2 = cell(1, n_start);
225 perfs2 = zeros(1, n_start);
226 val_t2 = zeros(n_start, length(val.x));
227 for i = 1:n_start
228     fprintf('Training %d/%d\n', i, n_start);
229     NN2{i} = train(net2, tr.x, tr.y2);
230     val_t2(i, :) = NN2{i}(val.x);
231     perfs2(i) = mse(net2, val_t2(i, :), val.y2);
232 end
233 [m2,i2] = min(perfs2)
234
235 % plot results (1.3.5)
236 figure;
237 plot(tr.x, tr.y2, 'o', 'MarkerSize',2, 'MarkerFaceColor', 'blue', ...
238      'MarkerEdgeColor', 'blue');
239 hold on;
240 plot(val.x, val.y2, 'o', 'MarkerSize',2, 'MarkerFaceColor', 'red', ...
241      'MarkerEdgeColor', 'red');
242 plot(val.x, val_t2(i2,:), '-', 'LineWidth',2, 'Color', 'green');
243 plot(tr.x, cos(2 * pi * tr.x), 'g-', 'LineWidth',2, 'Color', 'black');
244 hold off;
245 legend('Training Set', 'Validation Set', 'Approximated Function', 'True Function');
246 dim = [0.2 0.6 0.3 0.3];
247 annotation('textbox',dim, 'String', 'MSE = 0.3431', 'FitBoxToText', 'on');
248
249 % Training net3
250 setdemorandstream(0689432);
251 n_start = 100;
252 NN3 = cell(1, n_start);
253 perfs3 = zeros(1, n_start);
254 val_t3 = zeros(n_start, length(val.x));
255 for i = 1:n_start
256     fprintf('Training %d/%d\n', i, n_start);
257     NN3{i} = train(net3, tr.x, tr.y3);
258     val_t3(i, :) = NN3{i}(val.x);
259     perfs3(i) = mse(net3, val_t3(i, :), val.y3);
260 end
261 [m3,i3] = min(perfs3)
262
263 % plot results (1.3.6)
264 figure;
265 plot(tr.x, tr.y3, 'o', 'MarkerSize',2, 'MarkerFaceColor', 'blue', ...
266      'MarkerEdgeColor', 'blue');
267 hold on;
268 plot(val.x, val.y3, 'o', 'MarkerSize',2, 'MarkerFaceColor', 'red', ...
269      'MarkerEdgeColor', 'red');
270 plot(val.x, val_t3(i3,:), '-', 'LineWidth',2, 'Color', 'green');
271 plot(tr.x, cos(2 * pi * tr.x), 'g-', 'LineWidth',2, 'Color', 'black');
272 hold off;
273 legend('Training Set', 'Validation Set', 'Approximated Function', 'True Function');
274 dim = [0.2 0.6 0.3 0.3];
275 annotation('textbox',dim, 'String', 'MSE = 1.4290', 'FitBoxToText', 'on');
276
277
278 %%% Comparison of number of hidden neurons
279 clear
280 setdemorandstream(0689432);
281
282 % Size of the training data set is set to 150 observations.
283 tr.x = linspace(-1, 1, 150);

```

```

284 val_x = linspace(-0.9, 0.9, 150);
285
286 % Amount of noise (sd=1.2)
287 tr_y = cos(2 * pi * tr_x) + 1.2 * randn(size(tr_x));
288 val_y = cos(2 * pi * val_x) + 1.2 * randn(size(val_x));
289 [sum(tr_y);sum(val_y)]
290 % The remain characteristics are held constant.
291 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
292 net = fitnet(6, 'trainlm');
293 %net = fitnet(10, 'trainlm');
294 %net = fitnet(20, 'trainlm');
295
296 % Training net
297 setdemorandstream(0689432);
298 n_start = 100;
299 NN1 = cell(1, n_start);
300 perfs1 = zeros(1, n_start);
301 val_t1 = zeros(n_start, length(val_x));
302 for i = 1:n_start
303     fprintf('Training %d/%d\n', i, n_start);
304     NN1{i} = train(net, tr_x, tr_y);
305     val_t1(i, :) = NN1{i}(val_x);
306     perfs1(i) = mse(net, val_t1(i, :), val_y);
307 end
308 [m1,i1] = min(perfs1)
309
310 % plot results (1.3-7)
311 figure;
312 plot(tr_x, tr_y, 'o','MarkerSize',2,'MarkerFaceColor','blue',...
313      'MarkerEdgeColor','blue');
314 hold on;
315 plot(val_x, val_y, 'o','MarkerSize',2,'MarkerFaceColor','red',...
316      'MarkerEdgeColor','red');
317 plot(val_x, val_t1(i1,:), '-', 'LineWidth',2, 'Color', 'green');
318 plot(tr_x, cos(2 * pi * tr_x), 'g-', 'LineWidth',2, 'Color', 'black');
319 hold off;
320 legend('Training Set','Validation Set','Approximated Function','True Function');
321 dim = [0.2 0.6 0.3 0.3];
322 annotation('textbox',dim,'String',sprintf('MSE = %.4f',m1),'FitBoxToText','on');
323
324 %%% Comparison of training algorithms
325 clear
326 setdemorandstream(0689432);
327 % Size of the training data set is set to 150 observations.
328 tr_x = linspace(-1, 1, 150);
329 val_x = linspace(-0.9, 0.9, 150);
330 % Amount of noise (sd=1.2)
331 tr_y = cos(2 * pi * tr_x) + 1.2 * randn(size(tr_x));
332 val_y = cos(2 * pi * val_x) + 1.2 * randn(size(val_x));
333 [sum(tr_y);sum(val_y)]
334 % The remain characteristics are held constant.
335 % Early stopping, no regularization
336 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
337 net = fitnet(6, 'trainlm');
338 %net = fitnet(6, 'traingd'); % Backpropagation = Gradient descent
339 %net = fitnet(6, 'trainscg'); % Conjugate gradient
340 %net = fitnet(6, 'trainbfg'); % Quasi-Newton
341 % Training net
342 setdemorandstream(0689432);
343 n_start = 100;
344 NN1 = cell(1, n_start);
345 perfs1 = zeros(1, n_start);
346 val_t1 = zeros(n_start, length(val_x));
347 for i = 1:n_start
348     fprintf('Training %d/%d\n', i, n_start);
349     NN1{i} = train(net, tr_x, tr_y);
350     val_t1(i, :) = NN1{i}(val_x);
351     perfs1(i) = mse(net, val_t1(i, :), val_y);
352 end
353 [m1,i1] = min(perfs1)
354
355 % plot results (1.3-7)
356 figure;
357 plot(tr_x, tr_y, 'o','MarkerSize',2,'MarkerFaceColor','blue',...
358      'MarkerEdgeColor','blue');
359 hold on;
360 plot(val_x, val_y, 'o','MarkerSize',2,'MarkerFaceColor','red',...
361      'MarkerEdgeColor','red');
362 plot(val_x, val_t1(i1,:), '-', 'LineWidth',2, 'Color', 'green');
363 plot(tr_x, cos(2 * pi * tr_x), 'g-', 'LineWidth',2, 'Color', 'black');
364 hold off;
365 legend('Training Set','Validation Set','Approximated Function','True Function');
366 dim = [0.2 0.6 0.3 0.3];
367 annotation('textbox',dim,'String',sprintf('MSE = %.4f',m1),'FitBoxToText','on');
368
369 %%% Comparison of stopping rule
370 clear
371 setdemorandstream(0689432);
372 % Size of the training data set is set to 150 observations.
373 tr_x = linspace(-1, 1, 150);
374 val_x = linspace(-0.9, 0.9, 150);
375 % Amount of noise (sd=1.2)
376 tr_y = cos(2 * pi * tr_x) + 1.2 * randn(size(tr_x));
377 val_y = cos(2 * pi * val_x) + 1.2 * randn(size(val_x));
378 [sum(tr_y);sum(val_y)]
379 % The remain characteristics are held constant.

```

```

380 % Early stopping, no regularization
381 net = fitnet(6, 'trainlm');
382 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
383 net.trainParam.max_fail = 0
384 %net.trainParam.max_fail = 6
385 % Training net
386 setdemorandstream(0689432);
387 n_start = 100;
388 NN1 = cell(1, n_start);
389 perfs1 = zeros(1, n_start);
390 val_t1 = zeros(n_start, length(val.x));
391 for i = 1:n_start
392     fprintf('Training %d/%d\n', i, n_start);
393     NN1{i} = train(net, tr.x, tr.y);
394     val_t1(i, :) = NN1{i}(val.x);
395     perfs1(i) = mse(net, val_t1(i, :), val.y);
396 end
397 [m1,i1] = min(perfs1)
398
399 % plot results (1.3.7)
400 figure;
401 plot(tr.x, tr.y, 'o', 'MarkerSize',2,'MarkerFaceColor','blue',...
402     'MarkerEdgeColor','blue');
403 hold on;
404 plot(val.x, val.y, 'o', 'MarkerSize',2,'MarkerFaceColor','red',...
405     'MarkerEdgeColor','red');
406 plot(val.x, val_t1(i1,:), '-', 'LineWidth',2,'Color','green');
407 plot(tr.x, cos(2 * pi * tr.x), 'g-', 'LineWidth',2,'Color','black');
408 hold off;
409 legend('Training Set','Validation Set','Approximated Function','True Function');
410 dim = [0.2 0.6 0.3 0.3];
411 annotation('textbox',dim,'String',sprintf('MSE = %.4f',m1),'FitBoxToText','on');
412
413 %%% Comparison of regularization constant
414 clear
415 setdemorandstream(0689432);
416 % Size of the training data set is set to 150 observations.
417 tr.x = linspace(-1, 1, 150);
418 val.x = linspace(-0.9, 0.9, 150);
419 % Amount of noise (sd=1.2)
420 tr.y = cos(2 * pi * tr.x) + 1.2 * randn(size(tr.x));
421 val.y = cos(2 * pi * val.x) + 1.2 * randn(size(val.x));
422 [sum(tr.y);sum(val.y)]
423 % The remain characteristics are held constant.
424 % Early stopping, no regularization
425 net = fitnet(6, 'trainlm');
426
427 % Modifying regularization constant
428 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
429 net.performParam.regularization = 0
430 %net.performParam.regularization = 0.1
431 %net.performParam.regularization = 0.25
432 %net.performParam.regularization = 0.5
433
434 % Training net
435 setdemorandstream(0689432);
436 n_start = 100;
437 NN1 = cell(1, n_start);
438 perfs1 = zeros(1, n_start);
439 val_t1 = zeros(n_start, length(val.x));
440 for i = 1:n_start
441     fprintf('Training %d/%d\n', i, n_start);
442     NN1{i} = train(net, tr.x, tr.y);
443     val_t1(i, :) = NN1{i}(val.x);
444     perfs1(i) = mse(net, val_t1(i, :), val.y);
445 end
446 [m1,i1] = min(perfs1)
447
448 % plot results (1.3.7)
449 figure;
450 plot(tr.x, tr.y, 'o', 'MarkerSize',2,'MarkerFaceColor','blue',...
451     'MarkerEdgeColor','blue');
452 hold on;
453 plot(val.x, val.y, 'o', 'MarkerSize',2,'MarkerFaceColor','red',...
454     'MarkerEdgeColor','red');
455 plot(val.x, val_t1(i1,:), '-', 'LineWidth',2,'Color','green');
456 plot(tr.x, cos(2 * pi * tr.x), 'g-', 'LineWidth',2,'Color','black');
457 hold off;
458 legend('Training Set','Validation Set','Approximated Function','True Function');
459 dim = [0.2 0.6 0.3 0.3];
460 annotation('textbox',dim,'String',sprintf('MSE = %.4f',m1),'FitBoxToText','on');
461
462 %%% Comparison of initial weights
463 clear
464 setdemorandstream(0689432);
465 % Size of the training data set is set to 150 observations.
466 tr.x = linspace(-1, 1, 150);
467 val.x = linspace(-0.9, 0.9, 150);
468 % Amount of noise (sd=1.2)
469 tr.y = cos(2 * pi * tr.x) + 1.2 * randn(size(tr.x));
470 val.y = cos(2 * pi * val.x) + 1.2 * randn(size(val.x));
471 [sum(tr.y);sum(val.y)]
472 % The remain characteristics are held constant.
473 % Early stopping, no regularization
474 net = fitnet(6, 'trainlm');
475

```

```

476 % Training net
477 setdemorandstream(0689432);
478 % UNCOMMENT THE SEED THAT YOU WANT TO USE
479 %setdemorandstream(13245);
480 %setdemorandstream(67890);
481 %setdemorandstream(111111);
482 n_start = 100;
483 NN1 = cell(1, n_start);
484 perfs1 = zeros(1, n_start);
485 val_t1 = zeros(n_start, length(val.x));
486 for i = 1:n_start
487     fprintf('Training %d/%d\n', i, n_start);
488     NN1{i} = train(net, tr.x, tr.y);
489     val_t1(i, :) = NN1{i}(val.x);
490     perfs1(i) = mse(net, val_t1(i, :), val.y);
491 end
492 [m1,i1] = min(perfs1)
493
494 % plot results (1.3.7)
495 figure;
496 plot(tr.x, tr.y, 'o', 'MarkerSize',2, 'MarkerFaceColor', 'blue', ...
497      'MarkerEdgeColor', 'blue');
498 hold on;
499 plot(val.x, val.y, 'o', 'MarkerSize',2, 'MarkerFaceColor', 'red', ...
500      'MarkerEdgeColor', 'red');
501 plot(val.x, val_t1(i1,:), '-', 'LineWidth',2, 'Color', 'green');
502 plot(tr.x, cos(2 * pi * tr.x), 'g-', 'LineWidth',2, 'Color', 'black');
503 hold off;
504 legend('Training Set', 'Validation Set', 'Approximated Function', 'True Function');
505 dim = [0.2 0.6 0.3 0.3];
506 annotation('textbox',dim, 'String', sprintf('MSE = %4f',m1), 'FitBoxToText', 'on');
507
508 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
509 %% Curse of dimensionality
510 clear
511 setdemorandstream(0689432);
512 % m = 1
513 % training set
514 tr.x = linspace(-5,5,100)
515 tr.y = (sin(pi*tr.x))./(pi*tr.x)
516 % test set
517 test.x = linspace(-4.9, 4.9,100)
518 test.y = (sin(pi*test.x))./(pi*test.x)
519
520 % Net1.1: algorithms
521 setdemorandstream(0689432);
522 net1.1 = fitnet(5);
523 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
524 % net1.1 = fitnet(5,'trainbr'); % Bayesian Regularization
525 % net1.1 = fitnet(5,'trainbfg'); % Quasi-Newton
526 % net1.1 = fitnet(5,'trainscg'); % Scaled Conjugate Gradient
527
528 net1.1 = configure(net1.1, tr.x, tr.y);
529 net1.1.inputs{1}.processFcns = {};
530 net1.1.outputs{2}.processFcns = {};
531 n_start = 100;
532 NN1.1 = cell(1, n_start);
533 y_est1.1 = zeros(n_start, length(test.x));
534 perfs1.1 = zeros(1, n_start);
535 for i = 1:n_start
536     fprintf('Training %d/%d\n', i, n_start);
537     NN1.1{i} = train(net1.1, tr.x, tr.y);
538     y_est1.1(i, :) = NN1.1{i}(test.x);
539     perfs1.1(i) = mse(net1.1, y_est1.1(i,:), test.y);
540 end
541 [m1.1,i1.1] = min(perfs1.1)
542
543 % plot results
544 figure;
545 plot(tr.x, tr.y, 'g-', 'LineWidth',2, 'Color', 'black');
546 hold on;
547 plot(test.x, y_est1.1(i1.1,:), '-', 'LineWidth',2, 'Color', 'green');
548 hold off;
549 legend('True Function', 'Approximated Function');
550 dim = [0.2 0.6 0.3 0.3];
551 annotation('textbox',dim, 'String', 'MSE = 1.8436e-04', 'FitBoxToText', 'on');
552
553 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
554 % m = 2
555 %clear
556 tr.x = repelem(-5:0.1:5,length(-5:0.1:5))
557 tr.y = repmat(-5:0.1:5,1,length(-5:0.1:5))
558 tr_aux = sqrt(tr.x.^2 + tr.y.^2)
559 tr.z = (sin(pi*tr_aux))./(pi*tr_aux)
560 % all to matrix
561 tr.xmat = vec2mat(tr.x,101)
562 tr.ymat = vec2mat(tr.y,101)
563 tr.zmat = vec2mat(tr.z,101)
564 figure; surf(tr.xmat, tr.ymat, tr.zmat)
565 tr_data = [tr.x;tr.y]
566
567 % Test set
568 test.x= repelem(-4.9:0.125:4.9,length(-4.9:0.125:4.9))
569 test.y = repmat(-4.9:0.125:4.9,1,length(-4.9:0.125:4.9))
570 test_aux = sqrt(test.x.^2 + test.y.^2)

```

```

572 test.z = (sin(pi*test.aux))./(pi*test.aux)
573 % all to matrix
574 test.xmat= vec2mat(test.x,79)
575 test.ymat= vec2mat(test.y,79)
576 test.zmat = vec2mat(test.z,79)
577 figure;surf(test.xmat,test.ymat,test.zmat)
578 test.data = [test.x;test.y]
579
580 % Train neural network
581 % net2_1: LM algorithm
582 setdemorandstream(0689432);
583 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
584 net2_1 = fitnet(50);
585 % net2_1 = fitnet(50,'trainbr');
586 % net2_1 = fitnet(50,'trainbfg');
587 % net2_1 = fitnet(50,'trainscg');
588
589 net2_1 = configure(net2_1, tr_data, tr_z);
590 net2_1.inputs{1}.processFcns = {};
591 net2_1.outputs{2}.processFcns = {};
592 %view(net2_1)
593 n_start = 20;
594 NN2_1 = cell(1, n_start);
595 z_est2_1 = zeros(n_start, length(test.z));
596 perfs2_1 = zeros(1, n_start);
597 for i = 1:n_start
598     fprintf('Training %d/%d\n', i, n_start);
599     NN2_1{i} = train(net2_1, tr_data, tr_z);
600     z_est2_1(i, :) = NN2_1{i}(test.data);
601     perfs2_1(i) = mse(net2_1, z_est2_1(i,:), test.z);
602 end
603 [m2_1,i2_1] = min(perfs2_1)
604
605 % Transform data to matrix
606 z_estmat2_1 = vec2mat(z_est2_1(i2_1,:),79)
607 figure;surf(tr.xmat,tr.ymat,tr.zmat);
608 figure;surf(test.xmat,test.ymat,z_estmat2_1);
609
610
611 %%%%%%%%%%
612 % m = 5
613 %clear
614 tr.x1 = repmat(repelem(-5:2.5:5,length(-5:2.5:5)^4),1,length(-5:2.5:5))
615 tr.x2 = repmat(repelem(-5:2.5:5,length(-5:2.5:5)^3),1,length(-5:2.5:5)^2)
616 tr.x3 = repmat(repelem(-5:2.5:5,length(-5:2.5:5)^2),1,length(-5:2.5:5)^3)
617 tr.x4 = repmat(repelem(-5:2.5:5,length(-5:2.5:5)),1,length(-5:2.5:5)^4)
618 tr.x5 = repmat(-5:2.5:5,1,length(-5:2.5:5).^5)
619 %[tr.x1;tr.x2;tr.x3; tr.x4 ;tr.x5]
620 tr.aux = sqrt(tr.x1.^2 + tr.x2.^2 + tr.x3.^2 + tr.x4.^2 + tr.x5.^2)
621 tr.z = (sin(pi*tr.aux))./(pi*tr.aux)
622 tr.data = [tr.x1;tr.x2;tr.x3; tr.x4 ;tr.x5]
623
624 % Test set
625 test.x1 = repmat(repelem(-4.75:2.5:4.75,length(-4.75:2.5:4.75)^4),1,length(-4.75:2.5:4.75))
626 test.x2 = repmat(repelem(-4.75:2.5:4.75,length(-4.75:2.5:4.75)^3),1,length(-4.75:2.5:4.75)^2)
627 test.x3 = repmat(repelem(-4.75:2.5:4.75,length(-4.75:2.5:4.75)^2),1,length(-4.75:2.5:4.75)^3)
628 test.x4 = repmat(repelem(-4.75:2.5:4.75,length(-4.75:2.5:4.75)),1,length(-4.75:2.5:4.75)^4)
629 test.x5 = repmat(-4.75:2.5:4.75,1,length(-4.75:2.5:4.75).^5)
630 test.data = [test.x1;test.x2;test.x3;test.x4;test.x5]
631 test.aux = sqrt(test.x1.^2 + test.x2.^2 + test.x3.^2 + test.x4.^2 + test.x5.^2)
632 test.z = (sin(pi*test.aux))./(pi*test.aux)
633
634 % Train neural network
635 % net3_1: LM algorithm
636 setdemorandstream(0689432);
637 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
638 net3_1 = fitnet(100);
639 %net3_2 = fitnet(100,'trainbr');
640 %net3_3 = fitnet(100,'trainbfg');
641 %net3_4 = fitnet(100,'trainscg');
642 net3_1 = configure(net3_1, tr_data, tr_z);
643 net3_1.inputs{1}.processFcns = {};
644 net3_1.outputs{2}.processFcns = {};
645 n_start = 10;
646 NN3_1 = cell(1, n_start);
647 z_est3_1 = zeros(n_start, length(test.z));
648 perfs3_1 = zeros(1, n_start);
649 for i = 1:n_start
650     fprintf('Training %d/%d\n', i, n_start);
651     NN3_1{i} = train(net3_1, tr_data, tr_z);
652     z_est3_1(i, :) = NN3_1{i}(test.data);
653     perfs3_1(i) = mse(net3_1, z_est3_1(i,:), test.z);
654 end
655 [m3_1,i3_1] = min(perfs3_1)

```

• Script Exercise Session 2

```

1 % Setting working directory
2 cd '...'
3 %% Santa Fe laser data — time-series prediction
4 clear
5 % Importing data
6 lasertrain = importdata('lasertrain.dat');
7 laserpred = importdata('laserpred.dat');
8

```

```

9 % Plotting data
10 figure;plot(lasertrain);
11 figure;plot(laserpred);
12
13 % Transforming the format
14 lasertrain = tonndata(lasertrain,false,false)
15 laserpred = tonndata(laserpred,false,false)
16
17 %%% Comparing lags
18 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
19 net1 = narnet(1:2,10);
20 %net1 = narnet(1:7,10);
21 %net1 = narnet(1:14,10);
22 %net1 = narnet(1:21,10);
23
24 % Training net1
25 setdemorandstream(0689432);
26 n_start = 100;
27 NN1_1 = cell(1, n_start);
28 perfs1_1 = zeros(1, n_start);
29 test_t1_1(n_start).output = zeros(2, length(laserpred));
30 for i = 1:n_start
31     fprintf('Training %d/%d\n', i, n_start);
32     [Xs,Xi,Ai,Ts] = preparets(net1,{},{},lasertrain);
33     NN1_1{i} = closeloop(train(net1,Xs,Ts,Xi,Ai));
34     yini = lasertrain(end-NN1_1{i}.numLayerDelays+1:end); % Last values from training data
35     [Xs,Xi,Ai] = preparets(NN1_1{i},{},{},[yini laserpred]);
36 % predict on test data
37 test_t1_1(i).output = NN1_1{i}(Xs,Xi,Ai)
38 perfs1_1(i) = perform(NN1_1{i}, laserpred, test_t1_1(i).output)
39 end
40 [m1_1,i1_1] = min(perfs1_1)
41
42 figure;
43 plot(fromnndata(test_t1_1(i1_1).output, true, false, false), '-', 'LineWidth',2, 'Color', 'green');
44 hold on
45 plot(fromnndata(laserpred, true, false, false), 'g-', 'LineWidth',2, 'Color', 'black');
46 hold off;
47 legend('Approximated Function', 'True Function');
48 dim = [0.2 0.6 0.3 0.3];
49 annotation('textbox',dim, 'String', sprintf('MSE = %.4f',m1_1), 'FitBoxToText', 'on');
50 % White noise
51 E = gsubtract(laserpred, test_t1_1(i1_1).output);
52 ploterrcorr(E)
53
54 %%% Comparing hidden units
55 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
56 net1 = narnet(1:14,10);
57 %net1 = narnet(1:14,20);
58 %net1 = narnet(1:14,30);
59 %net1 = narnet(1:14,50);
60
61 % Training net1
62 setdemorandstream(0689432);
63 n_start = 100;
64 NN2_1 = cell(1, n_start);
65 perfs2_1 = zeros(1, n_start);
66 test_t2_1(n_start).output = zeros(2, length(laserpred));
67 for i = 1:n_start
68     fprintf('Training %d/%d\n', i, n_start);
69     [Xs,Xi,Ai,Ts] = preparets(net1,{},{},lasertrain);
70     NN2_1{i} = closeloop(train(net1,Xs,Ts,Xi,Ai));
71     yini = lasertrain(end-NN2_1{i}.numLayerDelays+1:end); % Last values from training data
72     [Xs,Xi,Ai] = preparets(NN2_1{i},{},{},[yini laserpred]);
73 % predict on test data
74 test_t2_1(i).output = NN2_1{i}(Xs,Xi,Ai)
75 perfs2_1(i) = perform(NN2_1{i}, laserpred, test_t2_1(i).output)
76 end
77 [m2_1,i2_1] = min(perfs2_1)
78
79 figure;
80 plot(fromnndata(test_t2_1(i2_1).output, true, false, false), '-', 'LineWidth',2, 'Color', 'green');
81 hold on
82 plot(fromnndata(laserpred, true, false, false), 'g-', 'LineWidth',2, 'Color', 'black');
83 hold off;
84 legend('Approximated Function', 'True Function');
85 dim = [0.2 0.6 0.3 0.3];
86 annotation('textbox',dim, 'String', sprintf('MSE = %.4f',m2_1), 'FitBoxToText', 'on');
87 % White noise
88 E = gsubtract(laserpred, test_t2_1(i2_1).output);
89 ploterrcorr(E)
90
91 %%% Comparing training algorithms
92 net1 = narnet(1:14,50);
93
94 % UNCOMMENT THE ALGORITHM YOU WANT TO USE
95 net1.trainFcn = 'trainlm';
96 %net1.trainFcn = 'traingd';
97 %net1.trainFcn = 'trainscg';
98 %net1.trainFcn = 'trainbfg';
99
100 % Training net1
101 setdemorandstream(0689432);
102 n_start = 100;
103 NN3_1 = cell(1, n_start);
104 perfs3_1 = zeros(1, n_start);

```



```

105 test_t3_1(n_start).output = zeros(2, length(laserpred));
106 for i = 1:n_start
107     fprintf('Training %d/%d\n', i, n_start);
108     [Xs,Xi,Ai,Ts] = preparets(net1,{},{},lasertrain);
109     NN3_1{i} = closeloop(train(net1,Xs,Ts,Xi,Ai));
110     yini = lasertrain(end-NN3_1{i}.numLayerDelays+1:end); % Last values from training data
111     [Xs,Xi,Ai] = preparets(NN3_1{i},{},{},[yini laserpred]);
112 % predict on test data
113     test_t3_1(i).output = NN3_1{i}(Xs,Xi,Ai)
114     perfs3_1(i) = perform(NN3_1{i}, laserpred, test_t3_1(i).output)
115 end
116 [m3_1,i3_1] = min(perfs3_1)
117
118 figure;
119 plot(fromnnndata(test_t3_1(i3_1).output, true, false, false), '-', 'LineWidth',2, 'Color', 'green');
120 hold on
121 plot(fromnnndata(laserpred, true, false, false), 'g-', 'LineWidth',2, 'Color', 'black');
122 hold off;
123 legend('Approximated Function', 'True Function');
124 dim = [0.2 0.6 0.3 0.3];
125 annotation('textbox',dim, 'String', sprintf('MSE = %.4f',m3_1), 'FitBoxToText', 'on');
126 % White noise
127 E = gsubtract(laserpred, test_t3_1(i3_1).output);
128 ploterrcorr(E)
129
130 %%% Comparing training algorithms
131 net1 = narnet(1:14,50);
132 % UNCOMMENT THE REGULARIZATION YOU WANT TO USE
133 net1.performParam.regularization = 0
134 %net1.performParam.regularization = 0.1
135 %net1.performParam.regularization = 0.25
136 %net1.performParam.regularization = 0.5
137
138 % Training net1
139 setdemorandstream(0689432);
140 n_start = 100;
141 NN4_1 = cell(1, n_start);
142 perfs4_1 = zeros(1, n_start);
143 test_t4_1(n_start).output = zeros(2, length(laserpred));
144 for i = 1:n_start
145     fprintf('Training %d/%d\n', i, n_start);
146     [Xs,Xi,Ai,Ts] = preparets(net1,{},{},lasertrain);
147     NN4_1{i} = closeloop(train(net1,Xs,Ts,Xi,Ai));
148     yini = lasertrain(end-NN4_1{i}.numLayerDelays+1:end); % Last values from training data
149     [Xs,Xi,Ai] = preparets(NN4_1{i},{},{},[yini laserpred]);
150 % predict on test data
151     test_t4_1(i).output = NN4_1{i}(Xs,Xi,Ai)
152     perfs4_1(i) = perform(NN4_1{i}, laserpred, test_t4_1(i).output)
153 end
154 [m4_1,i4_1] = min(perfs4_1)
155
156 figure;
157 plot(fromnnndata(test_t4_1(i4_1).output, true, false, false), '-', 'LineWidth',2, 'Color', 'green');
158 hold on
159 plot(fromnnndata(laserpred, true, false, false), 'g-', 'LineWidth',2, 'Color', 'black');
160 hold off;
161 legend('Approximated Function', 'True Function');
162 dim = [0.2 0.6 0.3 0.3];
163 annotation('textbox',dim, 'String', sprintf('MSE = %.4f',m4_1), 'FitBoxToText', 'on');
164 % White noise
165 E = gsubtract(laserpred, test_t4_1(i4_1).output);
166 ploterrcorr(E)
167
168 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
169 %%% Alfabet recognition
170 clear
171 openExample('nnet/apport1')
172
173 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
174 %%% Breast Cancer Wisconsin — classification problem
175 clear
176 filebreast = 'C:\Users\danie\Documents\Daniel Gil\KULeuven\Stage 2\Term 1\Data Mining and Neural Networks\Exercise
Sessions\Datasets\bcw.mat'
177 load(filebreast)
178
179 aux(Y==0) = 1
180 aux(Y==1) = 0;
181 YY= [Y; aux]
182 %vec2ind(YY) 1 is cancer and 2 is no cancer
183 clearvars aux
184
185 % division of data
186 % 90% for designing the neural net
187 % 10% for testing
188 Q = size(X,2) %569 columns
189 Q1 = floor(Q * 0.80); % 455 columns
190 Q2 = Q - Q1; % 114 columns
191 setdemorandstream(0689432)
192 ind = randperm(Q);
193 ind1 = ind(1:Q1);
194 ind2 = ind(Q1 + (1:Q2));
195 tr_x = X(:, ind1);
196 tr_y = YY(:, ind1);
197 test_x = X(:, ind2);
198 test_y = YY(:, ind2);
199 [sum(tr_x(1,:));sum(tr_y(1,:));sum(test_x(1,:));sum(test_y(1,:))]

```

```

200
201 %%% Comparing hidden units
202 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
203 net1 = patternnet(2);
204 %net1 = patternnet(5);
205 %net1 = patternnet(10);
206 %net1 = patternnet(15);
207
208 % Training net1
209 setdemorandstream(0689432);
210 n_start = 100;
211 NN1 = cell(1, n_start);
212 perfs1 = zeros(1, n_start);
213 test_t1(n_start).output = zeros(2, Q2);
214 for i = 1:n_start
215     fprintf('Training %d/%d\n', i, n_start);
216     NN1{i} = train(net1, tr_x, tr_y);
217     test_t1(i).output = NN1{i}(test_x);
218     perfs1(i) = perform(net1, test_y, test_t1(i).output);
219 end
220 [m1, i1] = min(perfs1)
221
222 [c1, cm1] = confusion(test_y, test_t1(i1).output)
223 fprintf('Percentage Correct Classification : %f%%\n', 100*(1-c1));
224 fprintf('Percentage Incorrect Classification : %f%%\n', 100*c1);
225 figure; plotconfusion(test_y, test_t1(i1).output)
226 figure; plotroc(test_y, test_t1(i1).output)
227 dim = [0.7 0.2 0.3 0.3];
228 corr = 100*(1-c1);
229 incorr = 100*c1;
230 str = {sprintf('CE = %4f', m1), sprintf('% Correct = %2f', corr), sprintf('% Correct = %2f', incorr)};
231 annotation('textbox', dim, 'String', str, 'FitBoxToText', 'on');
232
233 %%% Comparing algorithms
234 % UNCOMMENT THE NET THAT YOU WANT TO TRAIN
235 net1 = patternnet(5, 'trainscg');
236 %net1 = patternnet(5, 'trainlm'); % Levenberg
237 %net1 = patternnet(5, 'traingd'); % Backpropagation = Gradient descent
238 %net1 = patternnet(5, 'trainbfg'); % Quasi-Newton
239
240 % Training net1
241 setdemorandstream(0689432);
242 n_start = 100;
243 NN1 = cell(1, n_start);
244 perfs1 = zeros(1, n_start);
245 test_t1(n_start).output = zeros(2, Q2);
246 for i = 1:n_start
247     fprintf('Training %d/%d\n', i, n_start);
248     NN1{i} = train(net1, tr_x, tr_y);
249     test_t1(i).output = NN1{i}(test_x);
250     perfs1(i) = perform(net1, test_y, test_t1(i).output);
251 end
252 [m1, i1] = min(perfs1)
253
254 [c1, cm1] = confusion(test_y, test_t1(i1).output)
255 fprintf('Percentage Correct Classification : %f%%\n', 100*(1-c1));
256 fprintf('Percentage Incorrect Classification : %f%%\n', 100*c1);
257 figure; plotconfusion(test_y, test_t1(i1).output)
258 figure; plotroc(test_y, test_t1(i1).output)
259 dim = [0.7 0.2 0.3 0.3];
260 corr = 100*(1-c1);
261 incorr = 100*c1;
262 str = {sprintf('CE = %4f', m1), sprintf('% Correct = %2f', corr), sprintf('% Correct = %2f', incorr)};
263 annotation('textbox', dim, 'String', str, 'FitBoxToText', 'on');
264
265
266 %%% Comparing regularization
267 net1 = patternnet(5, 'trainscg');
268
269 % Modifying regularization constant
270 % UNCOMMENT THE REGULARIZATION YOU WANT TO USE
271 net1.performParam.regularization = 0
272 %net1.performParam.regularization = 0.01
273 %net1.performParam.regularization = 0.1
274 %net1.performParam.regularization = 0.2
275
276 % Training net1
277 setdemorandstream(0689432);
278 n_start = 100;
279 NN1 = cell(1, n_start);
280 perfs1 = zeros(1, n_start);
281 test_t1(n_start).output = zeros(2, Q2);
282 for i = 1:n_start
283     fprintf('Training %d/%d\n', i, n_start);
284     NN1{i} = train(net1, tr_x, tr_y);
285     test_t1(i).output = NN1{i}(test_x);
286     perfs1(i) = perform(net1, test_y, test_t1(i).output);
287 end
288 [m1, i1] = min(perfs1)
289
290 [c1, cm1] = confusion(test_y, test_t1(i1).output)
291 fprintf('Percentage Correct Classification : %f%%\n', 100*(1-c1));
292 fprintf('Percentage Incorrect Classification : %f%%\n', 100*c1);
293 figure; plotconfusion(test_y, test_t1(i1).output)
294 figure; plotroc(test_y, test_t1(i1).output)
295 dim = [0.7 0.2 0.3 0.3];

```

```

296 corr = 100*(1-c1);
297 incorr = 100-c1;
298 str = {sprintf('CE = %.4f',m1),sprintf('% Correct = %.2f',corr),sprintf('% Correct = %.2f',incorr)};
299 annotation('textbox',dim,'String',str,'FitBoxToText','on');

```

• Script Exercise Session 3

```

1 % Setting working directory
2 cd '...'
3
4 %% Dimensionality reduction by PCA analysis
5 clear
6 a) Investigate PCA analysis in order to achieve a dimensionality reduction
7 load cho_dataset
8
9 % Standardize the variables
10 [pn, std_p] = mapstd(choInputs);
11 [tn, std_t] = mapstd(choTargets);
12 % PCA
13 [pp, pca_p] = processpca(pn, 'maxfrac', 0.001);
14 [m, n] = size(pp)
15 % To get scores after pca (pp object) just do pca.p.transform*pn
16
17 % (b) For the case of 21 inputs, define a training, validation and test set and apply the Levenberg–Marquardt
    algorithm:
18 % Set indices for test, validation and training sets
19 Test_ix = 2:4:n;
20 Val_ix = 4:4:n;
21 Train_ix = [1:4:n 3:4:n];
22 % Configure a network
23 % LM algorithm
24 net1 = fitnet(5);
25 net1.divideFcn = 'divideind';
26 net1.divideParam = struct('trainInd', Train_ix, ...
27     'valInd', Val_ix, ...
28     'testInd', Test_ix);
29
30 % Training net1
31 setdemorandstream(0689432);
32 n_start = 100;
33 NN1 = cell(1, n_start);
34 mse_train1 = zeros(1, n_start);
35 mse_test1 = zeros(1, n_start);
36 yhat_train1(n_start).output = zeros(n_start, length(tn(:, Train_ix)));
37 yhat_test1(n_start).output = zeros(n_start, length(tn(:, Test_ix)));
38
39 for i = 1:n_start
40     fprintf('Training %d/%d\n', i, n_start);
41     NN1{i} = train(net1, pn, tn);
42     yhat_train1(i).output = NN1{i}(pn(:, Train_ix));
43     yhat_test1(i).output = NN1{i}(pn(:, Test_ix));
44     mse_train1(i) = mse(net1, tn(:, Train_ix), yhat_train1(i).output);
45     mse_test1(i) = mse(net1, tn(:, Test_ix), yhat_test1(i).output);
46 end
47 [m1,i1] = min(mse_test1)
48
49 % Investigate whether the performance can be improved by means of Bayesian regularization (trainbr). Compare the
    results on test data.
50 net2 = fitnet(5,'trainbr');
51 net2.divideFcn = 'divideind';
52 net2.divideParam = struct('trainInd', Train_ix, ...
53     'valInd', Val_ix, ...
54     'testInd', Test_ix);
55
56 % Training net2
57 setdemorandstream(0689432);
58 n_start = 100;
59 NN2 = cell(1, n_start);
60 mse_train2 = zeros(1, n_start);
61 mse_test2 = zeros(1, n_start);
62 yhat_train2(n_start).output = zeros(n_start, length(tn(:, Train_ix)));
63 yhat_test2(n_start).output = zeros(n_start, length(tn(:, Test_ix)));
64
65 for i = 1:n_start
66     fprintf('Training %d/%d\n', i, n_start);
67     NN2{i} = train(net2, pn, tn);
68     yhat_train2(i).output = NN2{i}(pn(:, Train_ix));
69     yhat_test2(i).output = NN2{i}(pn(:, Test_ix));
70     mse_train2(i) = mse(net2, tn(:, Train_ix), yhat_train2(i).output);
71     mse_test2(i) = mse(net2, tn(:, Test_ix), yhat_test2(i).output);
72 end
73 [m2,i2] = min(mse_test2)
74
75 mse_test1(i1)
76 mse_test2(i2)
77 mse_train1(i1)
78 mse_train2(i2)
79
80 % (c) Compare the training results of the case of 21 inputs (original inputs) and 4 inputs (after dimensionality
    reduction by PCA) by applying trainbr. Which choice would you make between the two options? Motivate your choice.
81 net3 = fitnet(5,'trainbr');
82 net3.divideFcn = 'divideind';
83 net3.divideParam = struct('trainInd', Train_ix, ...
84     'valInd', Val_ix, ...
85     'testInd', Test_ix);

```

```

86
87 % Training net2
88 setdemorandstream(0689432);
89 n_start = 100;
90 NN3 = cell(1, n_start);
91 mse_train3 = zeros(1, n_start);
92 mse_test3 = zeros(1, n_start);
93 yhat_train3(n_start).output = zeros(n_start, length(tn(:, Train_ix)));
94 yhat_test3(n_start).output = zeros(n_start, length(tn(:, Test_ix)));
95
96 for i = 1:n_start
97     fprintf('Training %d/%d\n', i, n_start);
98     NN3{i} = train(net3, pp, tn);
99     yhat_train3(i).output = NN3{i}(pp(:, Train_ix));
100    yhat_test3(i).output = NN3{i}(pp(:, Test_ix));
101    mse_train3(i) = mse(net3, tn(:, Train_ix), yhat_train3(i).output);
102    mse_test3(i) = mse(net3, tn(:, Test_ix), yhat_test3(i).output);
103 end
104 [m3, i3] = min(mse_test3)
105
106 [mse_train1(i1) mse_test1(i1); mse_train2(i2) mse_test2(i2); mse_train3(i3) mse_test3(i2)]
107
108 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
109 %% (2) Input selection by Automatic Relevance Determination (ARD)
110 demard
111 demev1
112
113 clear
114 fileion = 'ionstart.mat'
115 load(fileion)
116 YY=[Y>=0];
117
118 randn('state', 0689432);
119 rand('state', 0689432);
120
121 % Set up network parameters.
122 nin = 33; % Number of inputs.
123 nhidden = 2; % Number of hidden units.
124 nout = 1; % Number of outputs.
125 aw1 = 0.01*ones(1, nin); % First-layer ARD hyperparameters.
126 ab1 = 0.01; % Hyperparameter for hidden unit biases.
127 aw2 = 0.01; % Hyperparameter for second-layer weights.
128 ab2 = 0.01; % Hyperparameter for output unit biases.
129 beta = 50.0; % Coefficient of data error.
130
131 % Create and initialize network.
132 prior = mlpprior(nin, nhidden, nout, aw1, ab1, aw2, ab2);
133 net = mlp(nin, nhidden, nout, 'logistic', prior, beta);
134
135 % Set up vector of options for the optimiser.
136 nouter = 5; % Number of outer loops
137 ninner = 10; % Number of inner loops
138 options = zeros(1,18); % Default options vector.
139 options(1) = 1; % This provides display of error values.
140 options(2) = 1.0e-7; % This ensures that convergence must occur
141 options(3) = 1.0e-7;
142 options(14) = 500; % Number of training cycles in inner loop.
143
144 % Train using scaled conjugate gradients, re-estimating alpha and beta.
145 for k = 1:nouter
146     net = netopt(net, options, X, YY, 'scg');
147     [net, gamma] = evidence(net, X, YY, ninner);
148     fprintf(1, '\n\nRe-estimation cycle %d:\n', k);
149     disp('The first 33 alphas are the hyperparameters for the corresponding');
150     disp('input to hidden unit weights. The remainder are the hyperparameters');
151     disp('for the hidden unit biases, second layer weights and output unit')
152     disp('biases, respectively.')
153     fprintf(1, ' alpha = %8.5f\n', net.alpha);
154     fprintf(1, ' beta = %8.5f\n', net.beta);
155     fprintf(1, ' gamma = %8.5f\n\n', gamma);
156     %disp('')
157     %disp('Press any key to continue.')
158     %pause
159 end
160
161 % Hyperparameters
162 net.alpha
163 % Weights
164 net2.w1
165
166 %% Second neural net without some variables
167 X2 = X
168 X2(:, [8 15 17 19 28]) = [];
169 YY=[Y>=0];
170
171 randn('state', 0689432);
172 rand('state', 0689432);
173
174 % Set up network parameters.
175 nin = 28; % Number of inputs.
176 nhidden = 2; % Number of hidden units.
177 nout = 1; % Number of outputs.
178 aw1 = 0.01*ones(1, nin); % First-layer ARD hyperparameters.
179 ab1 = 0.01; % Hyperparameter for hidden unit biases.
180 aw2 = 0.01; % Hyperparameter for second-layer weights.
181 ab2 = 0.01; % Hyperparameter for output unit biases.

```

```

182 beta = 50.0; % Coefficient of data error.
183
184 % Create and initialize network.
185 prior2 = mlpprior(nin, nhidden, nout, aw1, ab1, aw2, ab2);
186 net2 = mlp(nin, nhidden, nout, 'logistic', prior2, beta);
187
188 % Set up vector of options for the optimiser.
189 nouter = 5; % Number of outer loops
190 ninner = 10; % Number of inner loops
191 options = zeros(1,18); % Default options vector.
192 options(1) = 1; % This provides display of error values.
193 options(2) = 1.0e-7; % This ensures that convergence must occur
194 options(3) = 1.0e-7;
195 options(14) = 500; % Number of training cycles in inner loop.
196
197 % Train using scaled conjugate gradients, re-estimating alpha and beta.
198 for k = 1:nouter
199     net2 = netopt(net2, options, X2, YY, 'scg');
200     [net2, gamma] = evidence(net2, X2, YY, ninner);
201     fprintf(1, '\n\nRe-estimation cycle %d:\n', k);
202     disp('The first 28 alphas are the hyperparameters for the corresponding');
203     disp('input to hidden unit weights. The remainder are the hyperparameters');
204     disp('for the hidden unit biases, second layer weights and output unit')
205     disp('biases, respectively.')
206     fprintf(1, ' alpha = %8.5f\n', net2.alpha);
207     fprintf(1, ' beta = %8.5f\n', net2.beta);
208     fprintf(1, ' gamma = %8.5f\n\n', gamma);
209     %disp(' ')
210     %disp('Press any key to continue.')
211     %pause
212 end
213 % Hyperparameters
214 net2.alpha
215 % Weights
216 net2.w1

```