

# Project 2: Classification and Regression, from linear and logistic regression to neural networks

A. Mattioli<sup>1</sup>, A. Pera<sup>1</sup>, and J. Štorek<sup>2</sup>

<sup>1</sup>University of Milano Bicocca

<sup>2</sup>Czech Technical University in Prague

November 2019

## Abstract

This report aims to provide insights into performance of different methods on both categorical and linear data. Two datasets have been studied regarding classification for categorical data and linear regression for the Franke's function data. Algorithms such as gradient descent – both simple and stochastic – and neural networks have been applied and compared to scikit-learn results. The obtained results agree with scikit-learn ones which generally perform better. Out of studied methods, neural networks emerge to be the most adequate predictive technique for both classification and regression problems.

## 1 Introduction

The main purpose of the project is to explore classification, as prediction of qualitative responses through logistic regression, and to re-evaluate linear regression problems, already addressed in project 1, by introducing multi-layer neural networks.

In particular, in the first part of the analysis, logistic regression is implemented by using a gradient descent algorithm on two different data sets: a classification case, in which the default payment is predicted, and the so-called Wisconsin Cancer data, a scikit-learn's dataset consisting of a binary classification problem of benign or malignant tumors.

A neural network algorithm has been applied on the same two datasets in order to find the optimal weights and biases and to analyze which combination of regularization parameters and learning rates leads to the best performance measure.

In the second part of the study, the neural network algorithm is adapted to a linear regression problem. For this part, Franke function from project 1 [1] has been used and results have been compared with the ones obtained previously.

A background theory of the used techniques can be found in the methods section. A description of the implemented code of the study is contained in the implementation section. The analysis of results can be read in the results and discussion section.

## 2 Methods

### 2.1 Logistic regression

Logistic regression is a widely-used classifier which allows to predict qualitative responses. In particular, logistic regression models the probability of Y to fall into a particular category rather than modelling directly the response itself.

The probability that a data point belongs to a category is given by the Sigmoid or logit function

$$p(t) = \frac{1}{1 + \exp(-t)} = \frac{\exp t}{1 + \exp t} \quad , \quad (1)$$

which represents the likelihood of an event.

Logistic regression models are fit by applying the Maximum Likelihood Estimation (MLE) criterion. In a binary case with independent data points, the likelihood expression is

$$P(D|\vec{\beta}) = \prod_{i=1}^n [p(y_i = 1|x_i, \vec{\beta})]^{y_i} [1 - p(y_i = 1|x_i, \vec{\beta})]^{1-y_i} \quad (2)$$

from which it is possible to obtain the negative log-likelihood and the cost function, the so-called Cross Entropy, as

$$C(\vec{\beta}) = - \sum_{i=1}^n (y_i(\beta_0 + \beta_1 x_i) - \log(1 + \exp(\beta_0 + \beta_1 x_i))) \quad (3)$$

which is easier to minimize in order to find the parameter estimates.

The binary logistic regression has multiple-class and multinomial extensions in which the gradient descent method is typically used in order to find the optimal parameters' values.

### 2.2 Softmax function

The Softmax is a function that takes as input a vector of the results of K linear functions and returns a probability distribution.

$$\sigma(z)_i = \frac{\exp z_i}{\sum_{j=1}^K \exp z_j} \quad (4)$$

for  $i = 1, \dots, K$  and  $z = (z_1, \dots, z_k)$ .

The softmax function is used in multiclass classification methods as multinomial logistic regression, discriminant analysis and neural networks to compute the probability of a unit to be in a certain class.

## 2.3 Simple gradient descent

Gradient descent is an iterative optimization algorithm for finding the minimum of a function and estimating the optimal values of specific parameters.

When fitting a model to a data set, the fit is actually performed by finding the values of  $\beta$  that minimize the cost function. Ideally, they should be determined analytically but when that is not possible the gradient descent can be used.

The basic idea of the algorithm is that function  $F(x)$  decreases faster in the direction of the negative gradient  $-\nabla F(x)$  and that the old value of  $x_k$  is replaced by the value  $x_{k+1}$  shifted closer to the minimum:  $x_{k+1} = x_k - \gamma_k \nabla F(x)$ , where the parameter  $\gamma_k > 0$  is called learning rate. The number of iterations can be then limited by reaching a particular absolute value of the gradient  $\|\nabla F(x)\|$  which corresponds to the distance from the minimum of  $F(x)$ . The first  $x_0$  can be chosen randomly. For  $\gamma_k$  small,  $F(x_{k+1}) \leq F(x_k)$ . In this way, starting from an initial random guess  $x_0$ , for a sufficiently small  $\gamma_k$  the movement is always towards smaller values.

## 2.4 Stochastic gradient descent

Stochastic gradient descent is a stochastic approximation of the simple gradient descent method, in the sense that in the process it replaces the gradient calculated on the entire data with an estimate of it calculated on a subset of the data.

The main idea behind the method is that the gradient can be computed as a sum of gradients, obtained considering the cost function as a sum over data points

$$C(\beta) = \sum_{i=1}^n c_i(\vec{x}_i, \beta) \quad (5)$$

from which

$$\nabla_{\beta} C(\beta) = \sum_i^n \nabla_{\beta} c_i(\vec{x}_i, \beta) \quad (6)$$

In order to approximate the gradient, the sum over all data points is replaced with a sum over the data points in one of the subsets, called minibatches, randomly picked at each gradient descent step.

Typically, a dataset is divided into  $n/M$  minibatches (where  $n$  is the total number of data points and  $M$  is size of minibatches) and the process iterates over all of them for a chosen number of epochs.

## 2.5 Neural Networks

Neural networks are computational systems that consist of layers of connected units sharing mathematical functions. Neural networks can learn to perform tasks by considering examples, generally without being programmed with any task-specific rule.

They can be classified in feed-forward neural networks, convolutional and recurrent ones. Only the first type, in which the information moves forward through the layers, is used for this project.

Each neural network is made of an input layer, one or more hidden layers (single or multi-layer perceptron) and an output layer. Each layer can contain a different number of neurons which need to be activated by a so-called activation function.

The algorithm is based on the output of the activation function

$$y = f\left(\sum_{i=1}^n w_i x_i + b_i\right) = f(z) \quad (7)$$

in which the inputs  $x_i$  are the outputs of previous layers,  $b_i$  is the bias and  $w_i$  are weights. In this way, each layer receives a weighted sum of the outputs of all units (or neurons) in the previous layer.

From a general point of view, all the layers could have different activation functions and their own number of units. The complete functional form for the output of the  $l^{\text{th}}$  layer is

$$y_i^l = f^l \left[ \sum_{j=1}^{N_{l-1}} w_j^l f^{l-1} \left( \sum_{k=1}^{N_{l-2}} w_k^{l-1} f^{l-2} \left( \dots f^1 \left( \sum_{n=1}^{N_0} w_n^1 x_n + b_n^1 \right) \dots \right) + b_k^{l-1} \right) + b_j^l \right] \quad (8)$$

where the only independent variables are the input  $x_n$ .

The most typical functions used as activation functions are the logistic Sigmoid, the hyperbolic tangent function and the softmax function, which is typically used in the last layer for classification problems.

Neural Networks also allow to use the final output in order to modify and adjust the weights, so that the errors are as small as possible, through a back propagation algorithm implemented by using gradient descent updates. The gradient which starts the algorithm depends on the derivatives of the cost function at the output layer

$$\frac{\partial C(\hat{W}^L)}{\partial w_{jk}^L} = \delta_j^L a_k^{L-1} \quad (9)$$

where  $a_k^{L-1}$  is the output of the activation function of the last hidden layer and

$$\delta_j^L = \frac{\partial C}{\partial b_j^L} \quad (10)$$

is the error of the output layer itself, equal to the rate of change of the cost function  $C$  as a function of the bias  $b$ .

In order to propagate backwards, the error in the last layer is used to compute the error in the previous layer and so on as

$$\delta_j^{l-1} = \sum_k \delta_k^l w_{kj}^l f'(z_j^{l-1}) \quad (11)$$

The weights and biases are updated through a gradient descent procedure as

$$w_{jk}^l = w_{jk}^l - \eta \delta_j^l a_k^{l-1} \quad (12)$$

$$b_j^l = b_j^l - \eta \delta_j^l \quad (13)$$

where  $\eta$  is the learning rate. Equal sign has been used in manner of programming language.

## 2.6 Performance measurements

In the machine learning context, checking the performance of a classification model is an important task. In order to evaluate and compare the performance of different methods, several metrics can be used.

Firstly, considering the number of correctly guessed targets over the total amount of targets, the accuracy score metric is defined as

$$A(t_i, y_i) = \frac{\sum_{i=1}^n I(t_i = y_i)}{n} \quad (14)$$

where  $n$  is the total number of targets and  $I$  is indicator function which returns 1 for  $t_i = y_i$  and 0 otherwise. Due to their definition, accuracy score values lay in the interval  $[0,1]$ , where 1 represents the accuracy score of the best possible classification method.

Another method which can be applied is called Area Under Curve (AUC), which computes the area under a Receiving Operating Characteristics (ROC) curve. ROC is a probability curve used to understand the extent to which the model is capable to distinguish between different classes. An example can be seen in Figure 7. This curve plots two parameters: the true positive rate (TPR) and the false positive rate (FPR). The first one is the ratio of true positive to all positives, and the second is the ratio of false positives to all negatives (binary classification). AUC ranges in value from 0 to 1. A model whose predictions are 100% wrong has an AUC of 0.0; one whose predictions are 100% correct has an AUC of 1.0, while an AUC of 0.5 reflects random predictions.

Finally, confusion matrices can be produced in order to study the amount of correctly guessed targets after the use of a specific predictive algorithm. A confusion matrix is a table which allows to show the way the classification method is confused when making predictions and it is strictly connected to the concept of accuracy score. In particular, the principal diagonal of the table contains the number of correctly predicted targets in all the classes studied, while the antidiagonal contains the number of wrong guesses.

For linear regression, instead, in order to quantify the capacity of the methods on fitting the data, two indices are used: the Mean Squared Error (MSE) and the  $R^2$  (variance score) statistics.

The MSE calculates the mean value of the squared difference between data  $\vec{y}$  and prediction  $\vec{\tilde{y}}$ :

$$MSE(\vec{y}, \vec{\tilde{y}}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2. \quad (15)$$

$R^2$  statistic provides an alternative method to express the quality of a fit by evaluating fraction of residual (RSS) and total sum of squares (TSS). It is calculated as

$$R^2(\vec{y}, \vec{\tilde{y}}) = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2}{\sum_{i=0}^{n-1} (y_i - \bar{y})^2}, \quad (16)$$

where  $\bar{y} = \frac{1}{n} \sum_{i=0}^{n-1} y_i$ .

### 3 Implementation

#### 3.1 Taiwan credit card dataset

The Taiwan credit card data dataset describes customer's default payments in Taiwan and contains a response binary variable (default payment: Yes=1, No=0) and 23 explanatory ones, among which 9 are categorical variables – sex, education, marriage, repayment status and default payment. Table 9 summarizes all the recorded variables with their possible values.

#### 3.2 Winsconsin breast cancer diagnostic dataset

The dataset, stored in the sklearn's repository, contains 569 instances and 30 numeric explanatory variables which describe different diagnostic features of benign and malignant tumors.

#### 3.3 Reading and pre-processing the categorical data - credit card case

Firstly, data has been read from the excel source file (see GitHub link [1]) and the rows containing values which have been found to be outside of the possible choices for every categorical predictor have been deleted using `pandas.DataFrame.isin` functionality. This was not the case for repayment status because of vast fraction of badly labeled data.

Hence, a closer look at repayment status data has been taken, noticing the presence of the value "-2" which does not match the right set of recorded values for those categorical variables. It has been decided to set all values "-2" equal to "0" in order to consider them as irrelevant during the analysis.

For repayment status, zero value is also considered as badly labeled and after hotencoding is transformed into zero. This is done manually by storing positions of hotencoded zeros in `list_bad` array. Values at these positions in design matrix are then set to zeros, i.e. repayment status for zero values is encoded only into zero vector which implies that it does not involve fitting procedure.

Secondly, the function `delete_zeros` (see GitHub link [1]) is applied in order to discard all the rows which contain more than 10 zero values in the columns regarding the amount of bill statement and previous payments because they do not contain any useful information. In particular, for each row the number of zeros contained in the previously stated columns is counted and a list is created aiming to track the lines which are to discard. Finally, the list of rows to be deleted is removed from the dataframe, with a loss of more than 1000 lines.

Thirdly, data from numerical predictors has been standardized using the `standardize` function (see GitHub link), which subtracts the mean from the value and divides it by the standard deviation for each column.

Features and targets have been split in order to be able to apply the one-hot-encoding procedure to each categorical variable by scikit-learn's `OneHotEncoder` and `ColumnTransformer` functionalities, obtaining a feature matrix which contains, for each one-hot-encoded variable, a number of columns corresponding to the number of possible categories, which consists of only ones and zeros.

The design matrix is then built from the feature matrix. In order to split the data between the set used for training (train set) and the one applied in the evaluation process (test set), scikit-learn's `train_test_split` function has been applied by setting 0.8 as training share.

### 3.4 Simple Gradient Descent

The gradient descent procedure is developed by the function `our_gradient_descent_beta` which requires the design matrix, the target vector, the learning rate and a small number epsilon as an input and returns vector beta as an output. In particular, the function firstly sets a random vector of length equal to the number of parameters to estimate as

```
betas=np.ones(X_1.shape[1])*0.001
```

then a variable `gradient_norm`, set equal to a small number ( $2 * \epsilon$  specifically) and a null vector of probabilities with length equal to the number of observed cases. Afterwards a loop, which stops whenever the norm of the gradient is less or equal to  $\epsilon$ , is set. Inside the loop, another loop over the different data lines is performed, computing the probability associated with each case according to the logistic regression probability formula

```
p[i] = (np.exp((X_1[i,:]).dot(betas)))/(1+np.exp((X_1[i,:]).dot(betas)))
```

for the i-th subject of the dataset. After completing the inner loop, it is possible to compute the gradient, update the parameters and check the gradient norm. The number of total iterations is also stored, in order to be able to refer to the

number of `while` loop iterations. The function returns the parameters which are then optimized through the gradient descent method itself and can be used in order to compute predicted values for the targets.

### 3.5 Stochastic Gradient Descent

The stochastic gradient descent is applied through the function `stochastic_GD` (see GitHub link), which requires the design matrix, the targets vector, the learning rate, the number of inputs (in terms of rows of the design matrix), the mini-batch size and the number of epochs as arguments. Firstly, the function sets a random vector of parameters which will be updated afterwards and computes the number of mini-batches by dividing the total number of inputs by the batch size.

Secondly, a first loop is set in order to go over all the epochs and another one is created to go over all the mini-batches. In the latter, a number `k` is randomly chosen among all possible mini-batches and both the design matrix and the target vector are sliced according to `k`, obtaining a subset for both of them. A null vector of probabilities is then set and filled with the actual probabilities computed for each element of the subset of the design matrix according to the logistic regression formula for probabilities. After that, the gradient is computed and used to update the parameters betas, which are returned as the result of the function after the loop over the epochs is concluded.

### 3.6 Performance measurements

In the classification cases, accuracy score, AUC and confusion matrices are used to evaluate the performance of different methods. The `accuracy_score` function is implemented in order to compute the accuracy of predictions made using the parameter estimates defined by the gradient descent procedures. In particular, a loop over all values in the probability vector is set aiming to count the number of predictions which correspond exactly to the target values. Finally, the accuracy is computed as the number of right predictions over the total amount of them according to equation (14). The accuracy is then also computed using scikit-learn's functionality `accuracy_score` and similar results are obtained.

The area under curve function is computed using sklearn's functionality `roc_auc_score`, while also storing results for the ROC curve through `roc_curve` functionality.

A confusion matrix is also produced for each of the methods studied, by applying the sklearn's `confusion_matrix`.

In terms of linear regression analysis, mean square error and  $R^2$  metrics are computed by using sklearn's functions `mean_squared_error` and `r2_score`.



### 3.7 Neural Network with feed forward and back propagation algorithms for classification

The neural network has been implemented using both back propagation and feed forward algorithms. Firstly, it has been decided to use three hidden layers, with 50, 40 and 30 neurons each respectively. The number of categories is set equal to two, therefore there is the need to one-hot encode the response vector. This is processed using the `to_categorical_numpy` function (see GitHub link), which allows to one-hot-encode both the train and test  $y$  subsets.

Secondly, weights and biases for the three hidden layers and the output layer are initialized as random vectors.

Combinations of different learning rates (eta) and the regularization parameters (lmbd), where values are set as

```
eta_vals = [0.0000001,0.000001,0.00001,0.0001,0.001]
lmbd_vals = [0.0001,0.001,0.01,0.1,1,10]
```

for the credit card data and as

```
eta_vals = [1e-5,1e-4,1e-3,1e-2,1e-1]
lmbd_vals = [1e-5,1e-4,1e-3,1e-2,1e-1,1]
```

for the breast cancer data, have been studied.

By running a loop over different iteration (with maximum number set to 100) it is possible to update weights and biases through the neural network procedure itself and to find out which of the combination gives the best results in terms of accuracy score and AUC values.

The back propagation algorithm is firstly applied using the `backpropagation` function, which requires the training subsets of both the design matrix and the one-hot-encoded response vector. The function immediately recalls the `feed_forward_train` process, which computes weighted sums of inputs for each layer and respective activations using the sigmoid activation function

```
def sigmoid(x):
    return 1/(1 + np.exp(-x)).
```

Weighted sums and activations are passed to next layer until the output layer is reached. The three activations and the vector of probabilities (output) are then returned as an output of `feed_forward_train` and used in the back propagation process.

In particular, the back propagation algorithm computes the error of each layer and the respective gradient for both weights and biases, starting from the output layer and going back to the first hidden one, considering the Cross Entropy as the cost function to be minimized. The back propagation function then returns all different computed gradients.

Next step is then to regularize all gradients through the parameter lambda and update weights and biases using the learning rate and the computed gradients, applying a gradient descent process.

After a set number of iterations of the loop, accuracy score can be computed and predictions for the response vector obtained by applying the function `predict`, which recalls a feed forward process in order to get a matrix probabilities and then evaluates the maximum value within each column, returning its index as a result. The validity of the algorithm itself is checked by using sklearn's `MLPClassifier` and computing both accuracy score and AUC on the test data.

### 3.8 Neural Network with feed forward and back propagation algorithms for linear regression

The neural network algorithm, applying both back propagation and feed forward, has been used to study a linear regression case from the previous project, in particular regarding the analysis of the Franke's function data.

After setting the design matrix, all the available data is split in test and training subsets, in order to perform more accurate analysis.

The applied neural network algorithm is the same as the one used for the classification case, except of the cost function which is no longer the Cross Entropy but becomes the Mean Square Error (MSE) which leads to a different definition of the error for the output layer. Another difference consists of the output produced in the feed forward algorithm: the softmax function is no longer needed in order to get probabilities in the output layer since the output is now given by the sigmoid function applied to the weighted sum of inputs to the output layer and depends on the output weights, output bias and activation in the previous hidden layer. As a result of this slightly different approach, predictions now depend only on the output produced by the feed forward algorithm.

Results are studied for different combinations of learning rates and regularization parameters as well as in the classification case.

## 4 Results and discussion

Breast cancer data, credit card data and Franke function's data have been studied by different methods such as simple and stochastic gradient descent and neural networks. The validity of all the results has been evaluated by accuracy score, Area Under Curve and confusion matrices in the classification case, by  $R^2$  and MSE statistics in the regression case.

A first insight of the classification problem datasets is given by the heatmaps Figure 1 and Figure 2, representing the correlations between features of both the credit card data and the breast cancer one.

In both cases, some variables seem to be highly correlated and even though a dimensionality reduction approach could have been taken, it has been preferred to consider all the variables during the analysis in order to maintain the original problem frame.

#### 4.1 Logistic regression for breast cancer data set

Method	Accuracy score
Simple gradient descent	0.9300699300699301
Stochastic gradient descent	0.9370629370629371
Sklearn accuracy score	0.958041958041958
Neural network: $\lambda = 0.01, \eta = 0.01$	0.9790209790209791
Neural network sklearn: $\lambda = 0.01, \eta = 0.01$	0.972027972027972

Table 1: Summary of accuracy score for different methods for breast cancer data.

Method	AUC
Simple gradient descent	0.9289308176100629
Stochastic gradient descent	0.9422431865828093
Neural network: $\lambda = 0.01, \eta = 0.01$	0.979454926624738
Neural network sklearn: $\lambda = 0.01, \eta = 0.01$	0.970020964360587

Table 2: Summary of Area Under Curve (AUC) of different used method for breast cancer data.

Results of all used methods for the breast cancer data, in terms of accuracy and AUC, are summarized in Tables 1 and 2.

It is possible to state that neural network is the algorithm which performs in the most satisfying way on the data analyzed. In particular the best accuracy, according to scikit-learn’s functionality `MPLClassifier`, is reached by using  $\lambda = 0.01$  and  $\eta = 0.01$  as parameters. The our own neural network code returns 0.979 as the highest accuracy computed, using the same parameters combination as scikit-learn’s.

Choosing the learning rate  $\eta$  is one of the most crucial parts of building a neural network, since a very small parameter may result in a long training process with long convergence time, whereas a very large value could lead to the learning of a sub-optimal set of weights in an extremely fast way or to an unstable training process. Another important decision to make is regarding the choice of the regularization parameter  $\lambda$ , which allows to put a constraint on the size of the weights aiming to contain their growth and to avoid overfitting. In particular, larger weight values are more penalized if  $\lambda$  is large. A deeper analysis of how differences among several combinations of learning rates and regularization parameters affect the accuracy score values can be found in Figures 3 and 4 in the appendix. Interestingly, it has been noticed that having the learning rate fixed, changes in the regularization parameter result in very small differences in terms of the performance metrics studied, underlining the key role that the learning rate plays in determining the best method to be used.

Accuracy is a very intuitive method used to evaluate the performance of a model, however it is reliable only when the number of false positive and false

negative values in a dataset is almost the same. Having to deal with biased data, it seems reasonable to compute an alternative statistics which aims to evaluate the performance of the algorithms implemented. In particular, the AUC (Area Under Curve) statistics is chosen to determine which method performs the best. In the case of breast cancer data, AUC agrees with accuracy score suggesting the neural network algorithm as the best for the data analyzed in terms of predictions given.

Confusion matrices showed in Figures 5 and 6 in the Appendix are displayed in order to give an insight of how well the different methods applied perform on the classification of the data itself. The dataset used is quite small, however it is possible to see that the neural network algorithm is the one giving the highest rate of correctly guessed targets, by looking at the percentages in the anti-diagonal elements.

Lastly, visualisation of AUC has been made by Receiver Operating Characteristics (ROC) in Figure 7.

## 4.2 Logistic regression on credit card data

Method	Accuracy score
Simple gradient descent	0.8238606369875066
Stochastic gradient descent	0.8215731127925392
Sklearn accuracy score	0.8272039415801513
Neural network: $\lambda = 10, \eta = 0.00001$	0.7858525426711244
Neural network sklearn: $\lambda = 0.0001, \eta = 0.0001$	0.8272039415801513

Table 3: Summary of accuracy score of different methods for credit card data.

Method	AUC
Simple gradient descent	0.6600631977964128
Stochastic gradient descent	0.6524705934208744
Neural network: $\lambda = 10, \eta = 0.0001$	0.6591697411268396
Neural network sklearn: $\lambda = 0.0001, \eta = 0.0001$	0.683187409779005

Table 4: Summary of Area Under Curve (AUC) for different used method for credit card data.

Results of all used methods in terms of accuracy and AUC for the credit card data are summarized in Tables 3 and 4.

In this case, it seems that the neural network algorithm does not perform as well as it does for the breast cancer problem. In particular, looking at the results obtained by using sklearn it seems that the algorithm itself performs quite well, although the outcome is not totally reliable due to convergence issues which are encountered during the iterations. The best conclusions, according to the accuracy score, are obtained by setting  $\lambda=10$  and  $\eta=0.00001$  with scikit-learn,

$\lambda=0.0001$  and  $\eta=0.0001$  with the own code (further studying of parameters can be seen in Figures 10 and 11 in the appendix).

In general, all the values obtained for accuracy scores with different methods are similar. It can be interesting to study the performance of these different methods using some other kind of statistics in order to see whether results would differ more or not. Aiming to achieve a result which is reliable in the case of biased data, AUC index is computed. Using this alternative criterion, it is possible to state that neural network is the method which performs the best on the data in terms of its classification, in particular with the combination  $\lambda=0.0001$  and  $\eta=0.0001$  with sklearn and  $\lambda=10$  and  $\eta=0.0001$  with the own code. As noticed for the breast cancer case, by looking at accuracy score and AUC values over different parameters combinations, it is clear that the learning rate has an important role in determining the best method. In fact, given a specific value for that parameter, the regularization term  $\eta$  does not lead to significant changes in the performance measurement values.

Since the dataset studied is larger than the breast cancer one, confusion matrices (Tables 8 and 9 in the appendix) are harder to interpret, although they show that all methods tend to perform in a similar way, even though the neural network seems to be the one which can classify the targets with  $y=1$  better.

### 4.3 Regression for Franke's function

	OLS	Ridge	Lasso
MSE	0.0020510889	0.0020673085	0.0087683865
$R^2$	0.9737030822	0.9734951306	0.8937880494

Table 5: Statistics obtained with our own code from project 1.

	OLS	Ridge	Lasso
MSE	0.0022951092	0.0020454841	0.0089400624
$R^2$	0.9699985517	0.9737485463	0.8914546373

Table 6: Statistics obtained with cross validation our own code from project 1.

	OLS	Ridge	Lasso
MSE	0.0020510889	0.0020673088	0.0087683865
$R^2$	0.9737030822	0.9734951268	0.8937880494

Table 7: Statistics obtained with scikit-learn functions from project 1.

In Tables 5 to 7 results from our previous studies of Franke's function with Ordinary Least Squares, Ridge and Lasso regression methods are summarized.

	MSE	$R^2$
NN own code: $\lambda = 0.00001, \eta = 1$	0.017093883312645722	0.9014514754185878
NN sklearn: $\lambda = 0.0001, \eta = 0.1$	0.0001340666188039936	0.9992270880035212

*Table 8: Statistics obtained with the neural network algorithm*

Aiming to check whether there existed more effective methods for prediction in a linear regression context, the analysis of the data produced by the so-called Franke function has been extended to a more complex method, which is neural networks.

Particularly, the neural network algorithm implemented to perform a prediction analysis on the regression problem studied has been evaluated in terms of reliability by studying the convergence to 0 of the MSE for the training data over 100 iterations, simultaneously checking the progressive increase of the  $R^2$  metric towards 1. The described behaviours can be observed in Figure 12.

By comparing previous results with the ones displayed in table 8, it can be stated that neural network tends to perform better than the other regression methods studied, particularly looking at sklearn’s results. In fact, both the MSE and the  $R^2$  statistics denote a better predictive approach on the data analyzed. The results obtained using the own code, however, tend to be less precise than the ones given by sklearn’s functionalities, since the latter use more efficient methods for weight optimization and prediction, such as the so-called adam gradient descent optimizer. Hence, sklearn’s outcomes can be considered as more reliable and they can be used to define the algorithm which works the best on the given linear regression problem, which in this case is the neural network.

Different combination of learning rate and regularization parameters have been studied (see table in the Appendix) in order to evaluate performance metrics over different models, resulting in  $\lambda=0.00001$  and  $\eta=1$  in the case of the own neural network algorithm and in  $\lambda=0.0001$  and  $\eta=0.1$  for sklearn as the best combination of parameters to use.

## 5 Conclusions

Different methods, such as simple gradient descent, stochastic gradient descent and neural networks, have been applied in order to evaluate the predictive performance on three different datasets – two classification and one linear regression problems. All the considered algorithms behave similarly over all kinds of analyzed problems, recognizing the neural network as the best method in terms of performance compared to the others. In the linear regression case, neural networks bring better results than the regression methods studied in the previous project, such as OLS, Ridge and Lasso. In the both classification problems, neural networks overcome gradient descent techniques in terms of prediction ac-

curacy, even though the results they give for credit card data are less impressive than the ones obtained for the breast cancer data case.

In comparison to our own algorithms, scikit-learn's functionalities performed always better except of neural network breast cancer data case. However, results obtained with our own algorithms agree with the scikit-learn's ones within 7% error, in most cases is the agreement even better (error within 4%).

Gradient descent has been found to be a simple method sensitive to learning rate  $\gamma$  choice. However, obtaining a precise result is quite time demanding. This can be improved by using stochastic gradient descent which utilizes idea of batches and avoids falling into local minimum.

On the other hand, neural networks are an excellent method to apply in order to get predictions for given data, although they consume high amounts of memory and the back-propagation algorithm inside them leads to the black-box problem, which is defined as the series of difficulties humans encounter when trying to understand the actual meaning of values in the learning process of the algorithm itself.

Gradient descent (GD) is a simple and fairly well studied algorithm, since it goes fast to near-optimum results, but it is relatively slow when close to the minimum. This issue is partially improved by using stochastic gradient descent (SGD), which updates the weights element-wise instead of as a whole. SGD is then computationally less demanding than GD since it doesn't run through the whole training set in each iteration. Both GD and SGD perform well on strongly convex problems and less in other situations, because knowing how far to travel along the gradient can result in a considerable complication.

In general, neural networks are an excellent method to apply in order to get predictions for given data, although they consume high amounts of memory and the back-propagation algorithm inside them leads to the black-box problem, which is defined as the series of difficulties humans encounter when trying to understand the actual meaning of values in the learning process of the algorithm itself.

In conclusion, neural network appears to be the best method among the ones analyzed to provide predictions of given responses in problems of different types. An interesting comparison could be given by studying the behaviour of other techniques, such as decision trees and convolutional neural networks, in the given context.

## References

- [1] Project repository accessible in Project2\_wd directory at [https://github.com/storejar/Data\\_analysis\\_and\\_machine\\_learning.git](https://github.com/storejar/Data_analysis_and_machine_learning.git)
- [2] James, G., Witten, D., Hastie, T., Tibshirani, R. (2013). An Introduction to Statistical Learning. Springer Texts in Statistics. Springer New York. <https://doi.org/10.1007/978-1-4614-7138-7>

- [3] Hastie, Trevor, et al. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer, 2017.
- [4] Morten Hjorth-Jensen, (2019), Lectures Notes in FYS-STK4155. Data Analysis and Machine Learning: Optimization and Gradient Methods, <https://compphysics.github.io/MachineLearning/doc/pub/Splines/html/Splines.html>
- [5] Morten Hjorth-Jensen, (2019), Lectures Notes in FYS-STK4155. Data Analysis and Machine Learning: Logistic Regression, <https://compphysics.github.io/MachineLearning/doc/pub/LogReg/html/LogReg.html>
- [6] Morten Hjorth-Jensen, (2019), Lectures Notes in FYS-STK4155. Data Analysis and Machine Learning: Neural networks, from the simple perceptron to deep learning, <https://compphysics.github.io/MachineLearning/doc/pub/NeuralNet/html/NeuralNet.html>

## A Appendix



Label	Description
LIMIT_BAL	Amount of given credit in New Taiwan (NT) dollars
SEX	Gender (1=male, 2=female)
EDUCATION	(1=graduate school, 2=university, 3=high school, 4=others)
MARRIAGE	Marital status (1=married, 2=single, 3=others)
AGE	Age in years
PAY_0	Repayment status in September, 2005 (-1=pay duly, 1=payment delay for one month, 2=payment delay for two months, ... 8=payment delay for eight months, 9=payment delay for nine months and above)
PAY_2	Repayment status in August, 2005 (scale same as above)
PAY_3	Repayment status in July, 2005 (scale same as above)
PAY_4	Repayment status in June, 2005 (scale same as above)
PAY_5	Repayment status in May, 2005 (scale same as above)
PAY_6	Repayment status in April, 2005 (scale same as above)
BILL_AMT1	Amount of bill statement in September, 2005 (NT dollar)
BILL_AMT2	Amount of bill statement in August, 2005 (NT dollar)
BILL_AMT3	Amount of bill statement in July, 2005 (NT dollar)
BILL_AMT4	Amount of bill statement in June, 2005 (NT dollar)
BILL_AMT5	Amount of bill statement in May, 2005 (NT dollar)
BILL_AMT6	Amount of bill statement in April, 2005 (NT dollar)
PAY_AMT1	Amount of previous payment in September, 2005 (NT dollar)
PAY_AMT2	Amount of previous payment in August, 2005 (NT dollar)
PAY_AMT3	Amount of previous payment in July, 2005 (NT dollar)
PAY_AMT4	Amount of previous payment in June, 2005 (NT dollar)
PAY_AMT5	Amount of previous payment in May, 2005 (NT dollar)
PAY_AMT6	Amount of previous payment in April, 2005 (NT dollar)
default.payment.next.month	Default payment in June, 2005 (1=yes, 0=no)

*Table 9: Summary of the all data variables.*

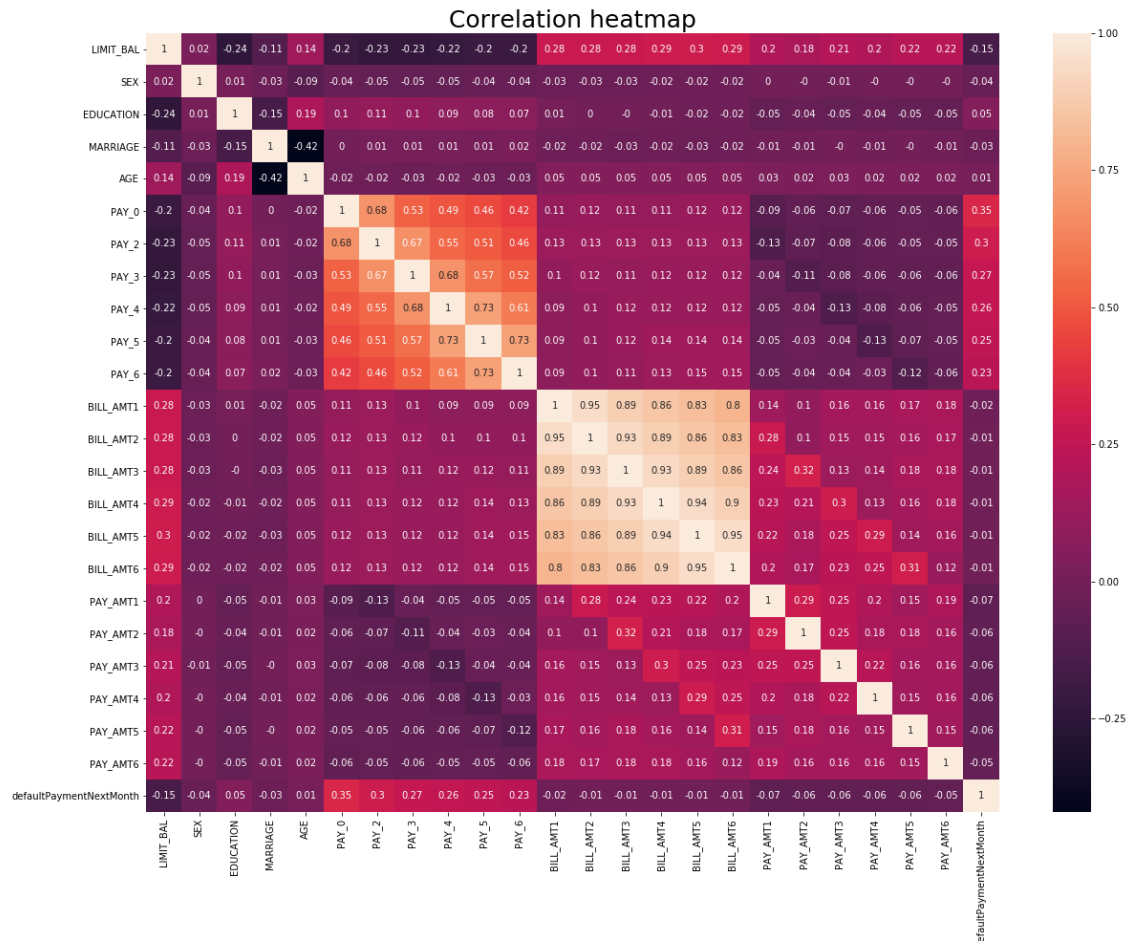


Figure 1: Correlation matrix heatmap for credit card data.

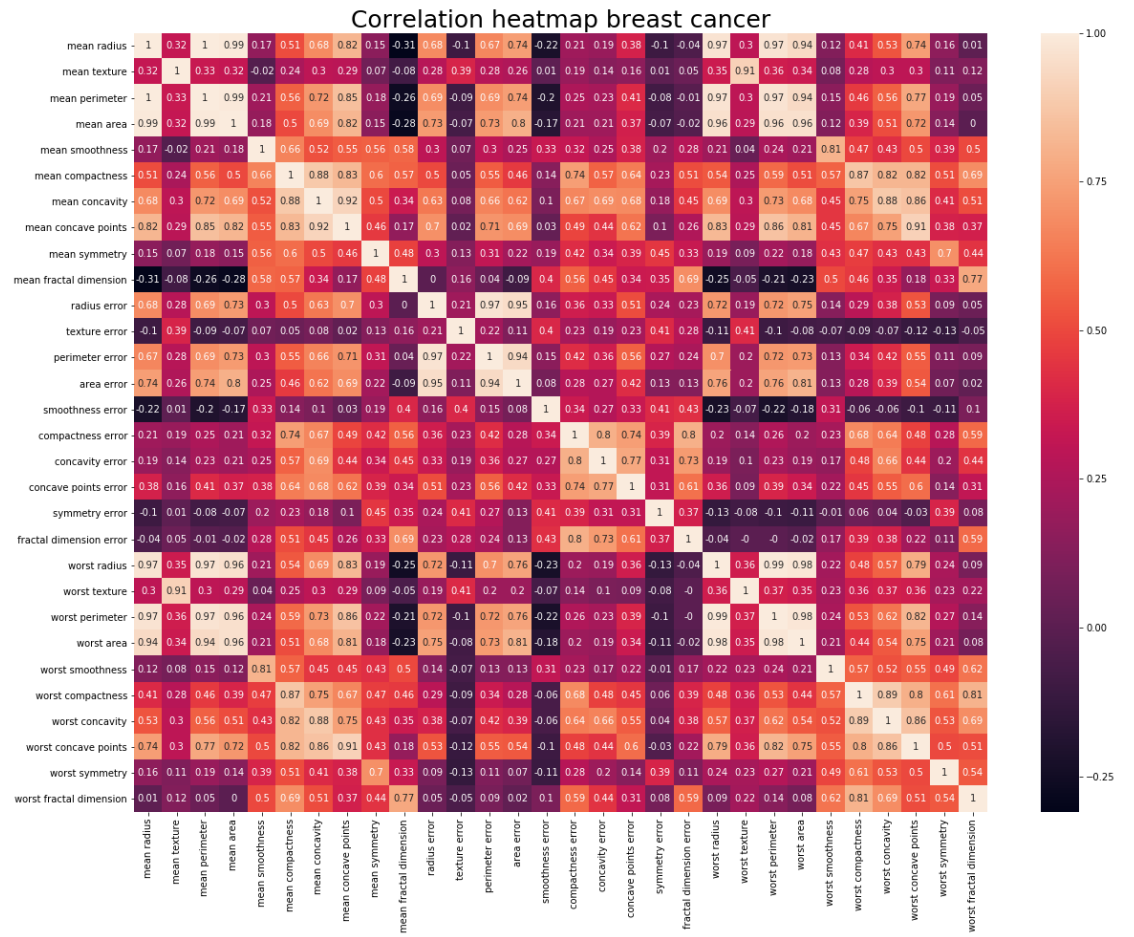


Figure 2: Correlation matrix Breast Cancer data set

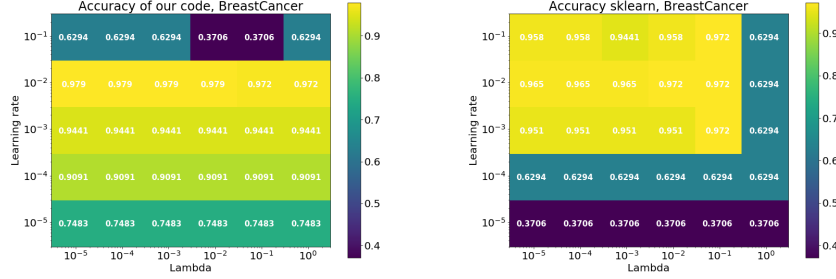


Figure 3: Accuracy dependence on learning rate  $\eta$  and lambda  $\lambda$  for breast cancer data left for our algorithm and right for scikit-learn.

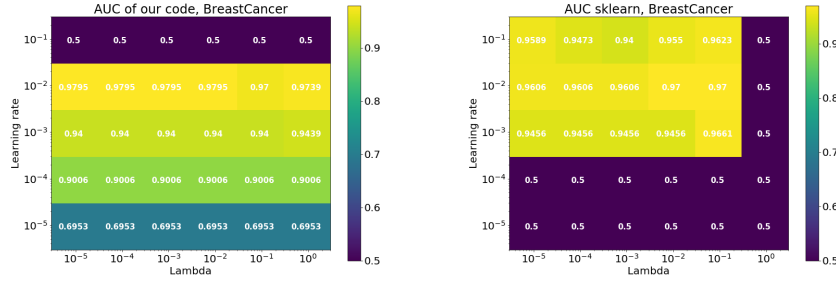


Figure 4: Area Under Curve dependence on learning rate  $\eta$  and lambda  $\lambda$  for breast cancer data left for our algorithm and right for scikit-learn.

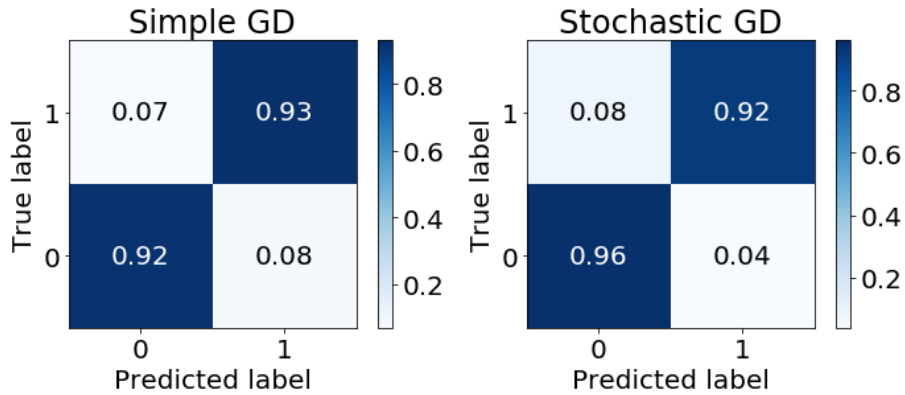


Figure 5: Confusion matrices for breast cancer data for our simple gradient descent (left) and stochastic gradient descent (right).

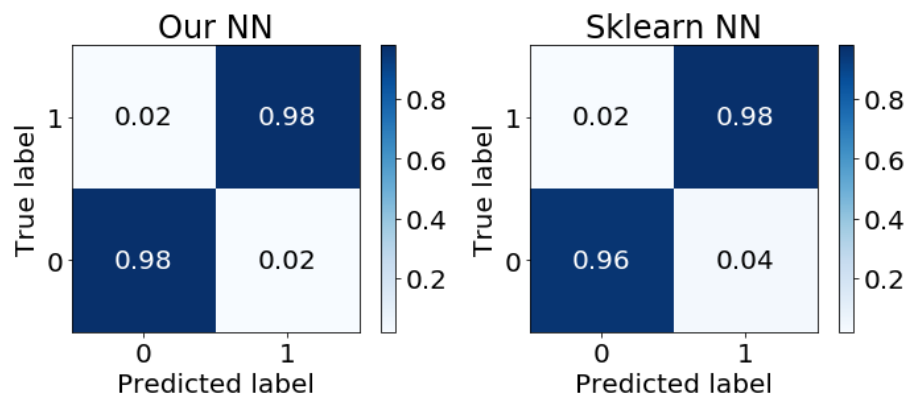


Figure 6: Confusion matrices for breast cancer data for Neural Networks with our own code on the left and from scikit-learn on the right.

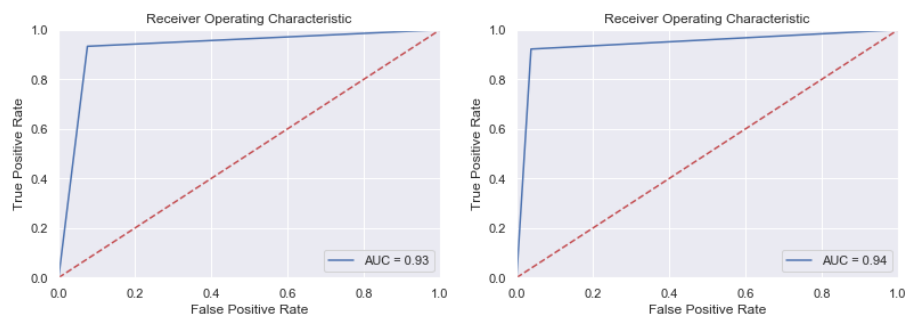


Figure 7: Receiver Operating Characteristic for simple gradient descent on the left and for stochastic gradient descent on the right on breast cancer data

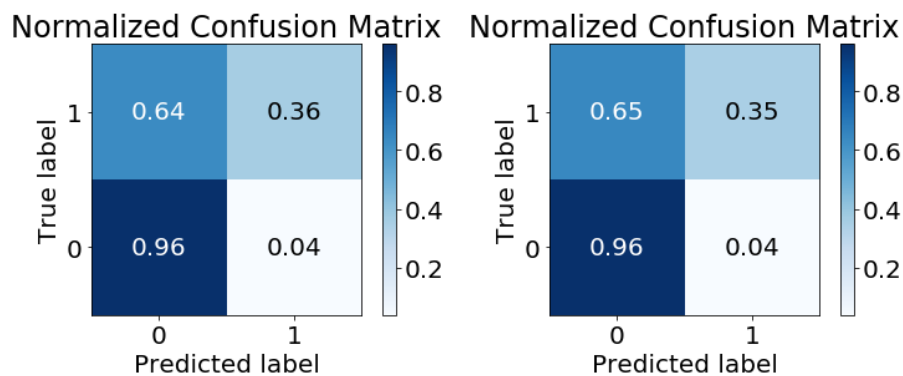


Figure 8: Confusion matrices on credit card data respectively for simple gradient descent on the left and stochastic gradient descent on the right

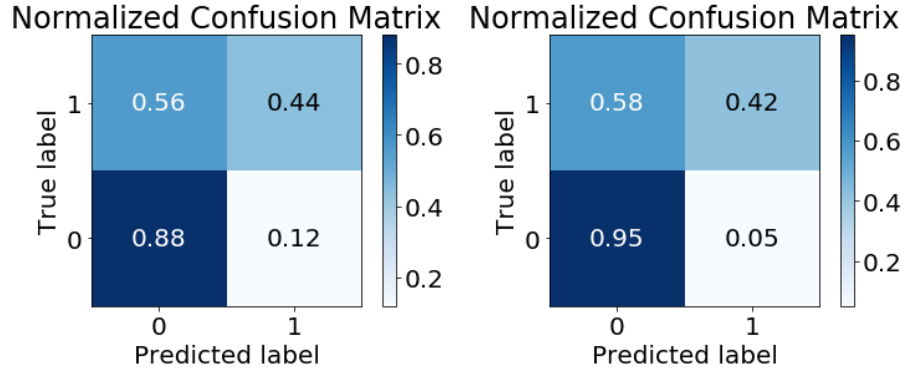


Figure 9: Confusion matrices on credit card data respectively for Neural Networks with own code on the left and from sklearn on the right

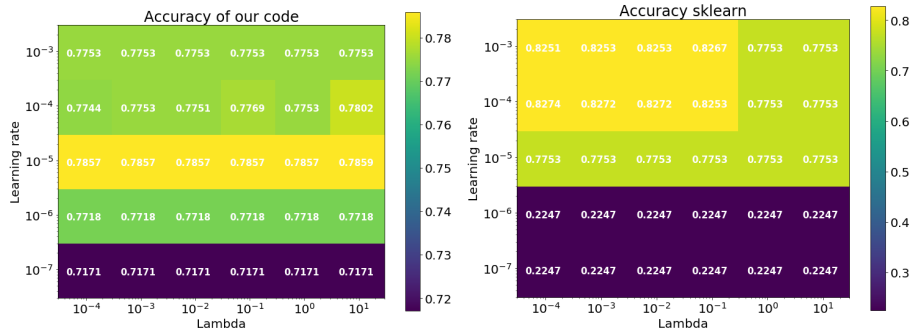


Figure 10: Accuracy dependence on learning rate  $\eta$  and lambda  $\lambda$  for credit card data left for our algorithm and right for scikit-learn.

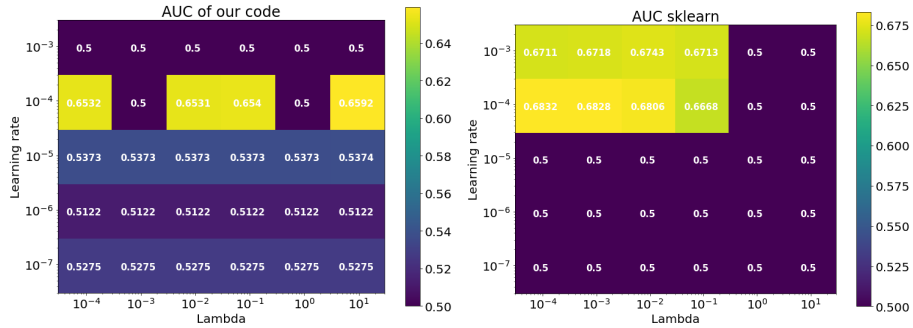


Figure 11: Area Under Curve dependence on learning rate  $\eta$  and lambda  $\lambda$  for credit card data left for our algorithm and right for scikit-learn.

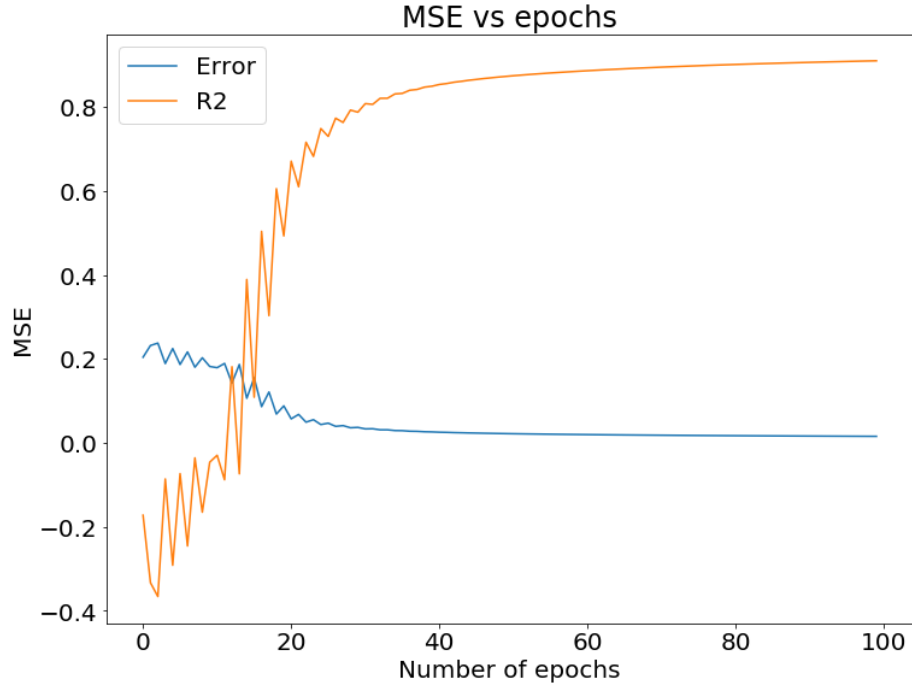


Figure 12: MSE and  $R^2$  computed with own Neural networks code and plotted against the number of epochs, for Franke's function.

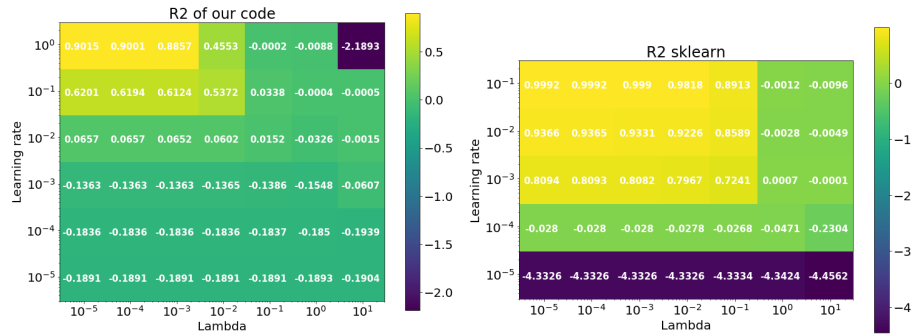


Figure 13:  $R^2$  dependence on learning rate  $\eta$  and lambda  $\lambda$  for Franke's function, left for our algorithm and right for scikit-learn.