

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

## Project 2 FYS-STK 3155

Liang Jia and Emil B. Hågan  
20.11.2021



# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Theory</b>	<b>2</b>
2.1	Stochastic Gradient Descent . . . . .	2
2.2	Logistic Regression . . . . .	3
2.3	Neural Network . . . . .	4
<b>3</b>	<b>Method and data</b>	<b>6</b>
3.1	Stochastic Gradient Descent . . . . .	6
3.2	Logistic Regression . . . . .	7
3.3	Neural Network . . . . .	8
<b>4</b>	<b>Results</b>	<b>10</b>
4.1	Stochastic Gradient Descent . . . . .	10
4.2	Logistic Regression . . . . .	11
4.3	FFNN with different activation functions . . . . .	12
4.4	NN with Keras . . . . .	15
<b>5</b>	<b>Discussion</b>	<b>16</b>
5.1	Stochastic Gradient Descent . . . . .	16
5.2	Logistic Regression . . . . .	16
5.3	FFNN with different activation functions . . . . .	16
5.4	NN with Keras . . . . .	17
<b>6</b>	<b>Conclusion</b>	<b>17</b>
<b>A</b>	<b>Appendix A CODE:</b>	<b>18</b>
A.1	Stochastic Gradient Decent . . . . .	18
A.2	FFNN simple case . . . . .	26
A.3	FFNN final code . . . . .	29
A.4	Breast cancer data implementation code . . . . .	35
A.5	Keras NN Cancer data code . . . . .	38
<b>B</b>	<b>Appendix B REFERENCES:</b>	<b>39</b>

## Abstract

As technology developed and our perception of reality and the universe changed our methods of describing the world changed with it. With the new age and the development of the computer we changed how we do calculations to a numerical perspective. With this tool we could reach new heights. One of these tools are *Neural Networks* (NN). With this method we can analyze big data sets, transcribe speech to text, filter noise from a recording. The aim of this article was to study regression and classification problems where we develop our own Feed-Forward Neural Network (FFNN). Where we first looked at Stochastic Gradient Descent with added momentum. Then creating our own FFNN, where we looked at different activation functions. Further we implement a classification analysis where we change our cost function to calculate an accuracy score for. For the final part we looked at logistic regression code and compare this method with our FFNN. The aim is to see how well our FFNN preforms and the capabilities of a NN. The data set we studied was the Wisconsin cancer data from pythons package *sci-kit learn*. The advantages of the Sigmoid function is that it preforms well in its bounds and overall not terrible outside this area, the same can be said for the ReLU and the linear functions. Although the linear function might lose some of the data and miss some of the information. The worse functions where the Leaky ReLU and tanh, which preformed so bad for this NN that we should not use them for this application. By testing different activation function in the output layer we also saw that the best preforming function was the softmax, where it gave consistent results. The softmax is normally used for where we have an output in multi-class classification problems. The implementation of logistic regression on the same data got quite similar performance, with right hyper parameters. They both can reach about 97% accuracy on test data and more than 99% on train data. The best activation functions for classification cases will be the softmax and sigmoid and for regression the linear would work best.

## 1 Introduction

We have always wanted answers to the big questions, as technology developed and our perception of reality and the universe changed our methods of describing the world changed with it. With the new age and the development of the computer we changed how we do calculations to a numerical perspective. With this tool we could reach new heights and one of these tools are *Neural Networks* (NN). With this method we can analyze big data sets, transcribe speech to text, filter noise from a recording as Nvidia is doing with "NVIDIA RTX Voice" and something you might not think of at first, but they way NETFLIX recommends entertainment to a specific viewer.

The aim of this article is to study regression and classification problems where we develop our own Feed-Forward Neural Network (FFNN). This might seem like a big task to start with so we will first look at Stochastic Gradient Descent (SGD) with added momentum. Then we will look at creating our own FFNN where we first create an implementation for a simple case and then generalize it further. With our FFNN we will also look at different activation functions such as Sigmoid, RELU and leaky RELU. Further we will implement a classification analysis where we change our cost function to calculate an accuracy score for *TRUE* or *FALSE*. For the final part we will look at logistic regression code and compare this method with our FFNN. The aim is to see how well our FFNN preforms and the capabilities of a Neural Network. The data set we studied was the Wisconsin cancer data from pythons package *sci-kit learn*, which is already well structured and ready for the implementation.

The code for the different parts can be found at our [GitHub](#)

## 2 Theory

### 2.1 Stochastic Gradient Descent

In this first section we will look at the SGD which derives from Gradient Descent (GD), but addresses some of the issues with GD. Some of these issues with the GD is that with this method we locate the local minima of our function, now the problem here is that we might not know if this minima is the one we want or another local minima within the function we study. When we for example look at terrain data we might get some issues with this. The next issue we will face with the GD is that the initial conditions matter heavily on the result. This problem is because with different initial conditions we might end up at another local minima and GD is also very sensitive to the learning rate, if this learning rate is small the training will take longer and will have need of more computational power. Therefore this heavily influence the computational time of the method and GD is computationally expensive to calculate for data sets which are large.

However SGD addresses these issues and is therefore much better. The SGD comes from studying the cost function and minimize it, the function can be expressed as a sum over  $n$  data points of  $\{\mathbf{x}_i\}_{i=1}^n$ ,

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

To compute the SGD we can define it as a sum over  $i$  gradients as such:

$$\nabla_{\beta} C(\beta) = \sum_{i=1}^n \nabla_{\beta} c_i(\mathbf{x}_i, \beta). \quad (2)$$

We also want to further develop this method and add momentum to the SGD. With the SGD we don't calculate the exact derivative of our loss function, but we estimate it in small batches. Therefore we are not always moving in the right direction of our minima, which comes from "noise" in our derivatives. The added momentum will help negate this issue. Another reason for adding momentum is in the ravines. This ravine is where the surface curve very sharply in one demotion, compared to any other. These commonly appeared close to these minimas and the SGD has issues with locating them. This is illustrated for normal SGD with figure 1 and with momentum in figure 1. We can look at momentum as a way to keeping memory of which direction we are moving in.

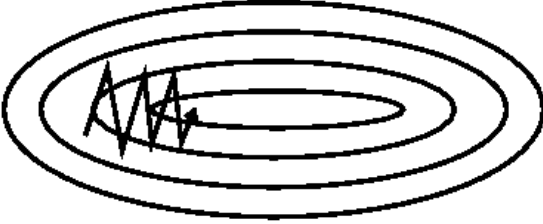


Figure 1: Shows SGD illustration with ravines from Bushaev, V. (2017, December 5) [5].

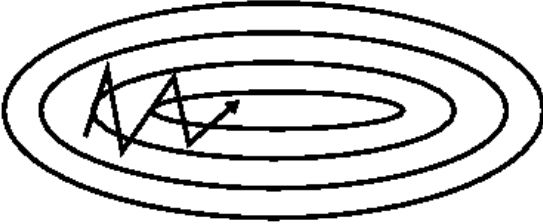


Figure 2: Shows SGD with momentum illustration with ravines from Bushaev, V. (2017, December 5) [5].

Now lets express SGD with momentum and a typical implementation of the method:

$$\mathbf{v}_t = \gamma \mathbf{v}_{t-1} + \eta_t \nabla_{\theta} E(\theta_t + \gamma \mathbf{v}_{t-1}) \quad (3)$$

$$\theta_{t+1} = \theta_t - \mathbf{v}_t, \quad (4)$$

here we introduce a momentum parameter  $\gamma$  which is defined within  $0 \leq \gamma \leq 1$ . The gradient is here also over all the different mini batches we just don't denote it for simplicity. The equation is updated with a Nesterov Acceleration Gradient or NAG rule and allows for a larger learning rate than that of a standard gradient decent with momentum. Lastly for this section we define our learning rate (LR) or schedule as:

$$lr = \frac{t_0}{t + t_1}, \quad (5)$$

where  $t_0$  and  $t_1$  is the number of mini-batches and epochs respectively. The only variable in the LR is  $t$  and it is this parameter that is sent to the equation to update the LR. Note that we in the code don't calculate the gradients or derivatives, but we use *autograd* which is a python library that does this for us.

## 2.2 Logistic Regression

In addition to SGD and NN, we will also implement the logistic regression algorithm, which is used to model probability of a certain category. Logistic regression is mostly applied in binary problems, true or false, positive or negative etc. Furthermore, logistic regression will guide us to the next part of neural network.

Since we aimed to use logistic regression for the binary classification, it is very important to encode the binary categories into 0 and 1. We chose the probability higher than or equal to 0.5 as 1 and otherwise 0.

$$y_i = \begin{bmatrix} 0 & no \\ 1 & yes \end{bmatrix} = \begin{bmatrix} 0 & p < 0.5 \\ 1 & p \geq 0.5 \end{bmatrix} \quad (6)$$

As known, the linear regression equation can be written as

$$y = \beta_0 + X_1\beta_1 + X_2\beta_2 + \dots + X_n\beta_n \quad (7)$$

Given that the probability of a data  $X_i$  belonging to a certain category can be calculated by logit/Sigmoid function as follow:

$$p(t) = \frac{1}{1 + \exp^{-t}} = \frac{\exp^t}{1 + \exp^t}, \quad (8)$$

where  $1 - p(t) = p(t)$

After the combination of this two equation we will take two parameters as illustration, we get the probabilities for each category by:

$$p(y_i = 1|X_i, \beta) = \frac{\exp^{\beta_0 + X_1\beta_1}}{1 + \exp^{\beta_0 + X_1\beta_1}} \quad (9)$$

$$p(y_i = 0|X_i, \beta) = 1 - p(y_i = 1|X_i, \beta) \quad (10)$$

The maximum likelihood estimation (MLE) was deployed to estimate the parameters of an assumed probability distribution with given a data set. Then our target is to maximize the probability of the given data set. The approximate likelihood of every individual probabilities with respect to a certain outcome  $y_i$  which can be obtained from:

$$p(D|\beta) = \prod_{i=1}^n [p(y_i = 1|X_i, \beta)]^{y_i} [1 - p(y_i = 1|X_i, \beta)]^{1-y_i} \quad (11)$$

The cost/loss function was obtained by negatively log-likelihood, which is known as the cross entropy in the statistics.

$$C(\beta) = - \sum_{i=1}^n (y_i(\beta_0 + X_1\beta_1) - \log(1 + \exp^{\beta_0 + X_1\beta_1})) \quad (12)$$

In addition the loss from the data set, the l2 regularization loss was also taken into consideration to avoid overfitting of the model on training data. Besides, the l2 regularization was applied in the code of Neural Network as well.

$$C_{reg} = \lambda \sum_{i=ij} w_{ij}^2 \quad (13)$$

In the training of logistic regression, we used simply SGD algorithm to get the optimized parameters and the minimized cost at the same time.

Finally, we will evaluate the performance of our model based on the following equation:

$$Accuracy = \frac{\sum_{i=1}^n I(\tilde{y}_i = y_i)}{n} \quad (14)$$

We will also need the R2 score which is defined as:

$$R^2(\hat{y}\hat{y}) = 1 - \frac{\sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2}{\sum_{i=0}^{n-1} (y_i - \bar{y})^2}, \quad (15)$$

and the mean value of  $\tilde{y}$  is given by:

$$\bar{y} = \frac{1}{n} \sum_{i=0}^{n-1} y_i. \quad (16)$$

## 2.3 Neural Network

We have now implemented and understood some linear models from Ordinary Least Square (OLS), Ridge, Lasso to the SGD. Now we want to look at Artificial Neural Network (ANN). There are a few different variants of NN. Here we will only focus on the first and simplest one, Feed Forward Neural Network (FFNN), where we have a visual illustration in figure 3 describing how a NN works.

We can define the weights and biases as  $P = \{P_{hidden}, P_{outout}\}$  for each neuron and we have  $N_{hidden}$  amount of neurons of features in the hidden layer. This  $P_{hidden}$  will be a matrix with  $N_{features} \times N_{hidden}$  dimensions, where  $N_{features}$  are number of features of the input layer. We define our  $P_{hidden}$  such that we have the bias in the first column and weights in the other for each hidden neurons. For the output layer the  $P_{outout}$  will be a  $N_{hidden} \times N_{outout}$  matrix, with the same shape as  $P_{hidden}$  for each neuron. We should define the  $N_{outout}$  layer based on our targets to be predicted. We then let the function we want to study be  $F(N(x, P)) = xN(x, N)$  and  $A(x)$  be where our function is  $i = 0$ . Then we want to find the cost function our NN must solve:

$$\begin{aligned} & \min \left\{ (F'(N(x, P))) - (-\gamma F(N(x, P)))^2 \right\}, \\ & \min \{ (F'(N(x, \{P_{hidden}, P_{outout}\}))) \\ & - (-\gamma F(N(x, \{P_{hidden}, P_{outout}\})))^2 \}. \end{aligned} \quad (17)$$

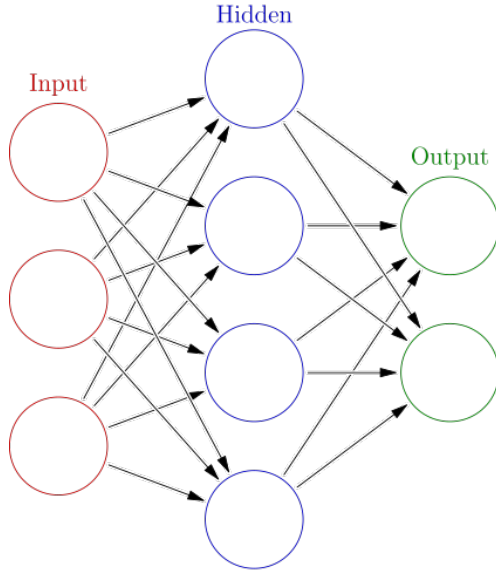


Figure 3: Shows a simple FFNN structure from Wikipedia [2]

We don't just minimize for one value of  $x$ , but also for every  $N$  iteration as well and then the total error becomes derived from equation 18:

$$\min \left\{ \frac{1}{N} \sum_{i=1}^N (F'(N(x, P))) - (-\gamma F(N(x, P)))^2 \right\} \quad (18)$$

then we let  $x$  be a vector with  $x_i$  elements and also described in terms of  $P_{hidden}$  and  $P_{output}$ :

$$\min_{P_{hidden}, P_{output}} = C(x, \{P_{hidden}, P_{output}\}). \quad (19)$$

The cost function we chose is the Mean Square Error (MSE) for the regression problem, which is defined as:

$$MSE(\hat{y}\hat{y}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 \quad (20)$$

For the loss function we use cross entropy to solve the classification problem, which is defined in equation 12.

Now we want to move on and describe the way we obtain  $z$  from the hidden biases and weights. We can

define  $z_{i,j}^{hidden}$ , where  $j$  is the neuron of the input layer, as:

$$\begin{aligned} z_{i,j}^{hidden} &= b_i^{hidden} + w_i^{hidden} x_j \\ &= (b_i^{hidden} + w_i^{hidden}) \begin{pmatrix} 1 \\ x_j \end{pmatrix}, \end{aligned} \quad (21)$$

furthermore we can define for weighting at  $i$ -th iteration of the hidden neurons as:

$$\begin{aligned} z_i^{hidden} &= (b_i^{hidden} + w_i^{hidden} x_1 + b_i^{hidden} \\ &\quad + w_i^{hidden} x_2 + \dots + b_i^{hidden} + w_i^{hidden} x_N) \\ &= (b_i^{hidden} \quad w_i^{hidden}) \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_N \end{pmatrix} \\ &= p_{i,hidden}^T X, \end{aligned} \quad (22)$$

We have that  $p_{i,hidden}^T$  is each row in  $P_{hidden}$ . Now that we have found  $z_i^{hidden}$  for every iteration of  $i$  we can send this to the activation function of our choice  $a_i(z)$  and we will have the output of this function  $x_i^{hidden}$  equal to:

$$x_i^{hidden} = f(z_i^{hidden}) \quad (23)$$

The output is then sent to the output layer which is comprised of  $N$  neurons and combines for each of the neurons in our hidden layer. Lastly the output layer combines the weights and biases of the output layer and we now have the output of  $z_i^{output}$ :

$$z_i^{output} = (b_i^{output} \quad w_i^{output}) \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_N \end{pmatrix} \quad (24)$$

In this NN we also wanted to add back propagation and feed forward which we will define now and we will start with back propagation. We have to choose a optimization method to minimize the cost function and we will use SGD with momentum, which is described above, for this purpose.

The activation functions we want to study are the Sigmoid, ReLU, Leaky ReLU, the tanh function and

a linear function which is just the same as its inputs. They are defined as such:

$$\begin{aligned} \text{Sigmoid} &= \frac{1}{1 + e^{-x}} \\ \text{ReLU} &= \max(0, x) \\ \text{Leaky ReLU} &= \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{otherwise} \end{cases} \\ \text{tanh} &= \frac{e^x - e^{-x}}{e^x + e^{-x}} \end{aligned} \quad (25)$$

One might also want to add an accuracy part to see how well we predict the values. We want to study the accuracy of the FFNN and it is described as the sum of correctly guessed targets divided by the total number of the set. We then have the function:

$$\text{Accuracy} = \frac{\sum_{i=1}^n I(t_i = y_i)}{n} \quad (26)$$

Where  $t_i$  is the wanted target and  $y_i$  is the one from our FFNN and  $n$  is the total number of the set of targets.

### 3 Method and data

#### 3.1 Stochastic Gradient Descent

We started by looking at some example code of GD, then started to code with the equations above for SGD and momentum 1 and 4 respectively. Bellow is an outtake of our implementation of just the SGD with momentum calculation. Then we use this code with conjunction with our code for project one with some small adjustments. We have changed the inversion calculation the Ordinary Least Squares (OLS) and Ridge regression method with our SGD. We do this simply by calling our SGD in the class and changing our x and y input to that of the OLS and Ridge regression. We now get the theta by calling our SGD *"StochasticGradientDecent(x, y).SGD()"*.

```

1  def SGD(self):
2      X = np.c_[np.ones((len(self.x_full), 1),
3      self.x_full)]
4      #size_matrix = x.shape[0]
5      self.theta = np.random.randn(X.shape
6      [1], 1) #Initilize theta for matrix shape.
7
8      #xy = np.c_[x.reshape(size_matrix, -1)
9      , y.reshape(size_matrix, 1)]
10
11     #Main SGD loop for epochs of
12     minibatches

```

```

9      for epoc in range(self.n_epoc):
10         #Second SGD loop with random
11         choice of k
12         for k in range(self.m):
13             random_index = self.M*np.
14             random.randint(self.m)
15             xi = X[random_index:
16             random_index+self.M]
17             yi = self.y_full[random_index:
18             random_index+self.M]
19
20             eta = self.ls(epoc*self.m+k) #
21             Calling function to cal. eta
22
23             #self.v_ = gamma*self.v_ + eta
24             *gradient(x_iter, y_iter, self.theta -
25             gamma*self.v_) #Cal. v where gradient is
26             from autograd
27             place_hold = self.theta + self
28             .gamma*self.v_
29             x_grad = egrad(self.gradient,
30             2) #Gradient with respect to theta
31             self.v_ = self.gamma*self.v_ +
32             eta * x_grad(xi, yi, place_hold) #Cal. v
33             where gradient is from autograd, self.
34             gradient(xi, yi, self.theta)
35             self.theta = self.theta - self
36             .v_ #Theta +1 from this itteration of
37             theta and v
38
39     return self.theta

```

To collect the data we do two different things, firstly we write the R2 score and the MSE (Mean Square Error) to a file and plot the MSE and R2 score. We also plot the noise on the MSE as well.

We also do some initial testing to make sure we input the correct parameters and variables. Where we check for the length of x and y, that the epochs are the right size and that  $\gamma$  is in range. This makes it easier to know if something is wrong with the initialization of the code and so we don't have to check for every new line of code each time we do something new.

```

1  def __call__(self):
2
3      #Checks matrix size of rows
4      size_matrix = self.x_full.shape[0]
5      if size_matrix != self.y_full.shape
6      [0]:
7          raise ValueError("'x' and 'y' must
8          have same rows")
9
10     #Check to see if batches are right
11     size
12     self.n_epoc = int(self.n_epoc)
13     if not 0 < self.n_epoc <= size_matrix:
14         raise ValueError("Must have a
15         batch size less or equal to observations
16         and greater than zero.")
17
18     #Checks gamma is in range
19     if not 0 <= self.gamma <= 1:
20         raise ValueError("Gamma must be
21         equal or greater than zero and equal or
22         less than 1.")

```

```

17     #Checks gamma is in range
18     if not 0 <= self.gamma <= 1:
19         raise ValueError("Gamma must be
equal or greater than zero and equal or
less than 1.")

```

## 3.2 Logistic Regression

Since logistic regression is a linear model, we start by initializing the parameters (weights and biases) for the linear equation. For simplicity, we just initialize the weights as a zero matrix with number of features of input data with 1 dimension and bias as 0. This way we could easily test if the training session is changing the parameters.

```

1  def wb_init(self, dim):
2      weights = np.zeros((dim,1))
3      bias = 0
4      return weights, bias

```

The cost function was obtained from combining the Sigmoid function and linear transform of the input data.

```

1  def sigmoid(self, z):
2      s = 1/(1 + np.exp(-z))
3      return s
4
5  def probability(self, weights, X, bias):
6      prob = self.sigmoid(np.dot(weights.T,
7      X) + bias)
8      return prob
9
10 def cost(self, prob, y, num_sample,
11 weights):
12     data_cost = np.mean(-np.sum(y * np.log
13     (prob) + (1-y) * np.log(1-prob)))
14     self.data_cost.append(data_cost)
15
16     reg_cost = 0
17     if self.regularizer_l2 > 0:
18         reg_cost = self.regularizer_l2 *
19         np.sum(weights**2)
20     self.reg_cost.append(reg_cost)
21
22     #print("data_cost", data_cost)
23     #print("reg_cost", reg_cost)
24
25     cost = data_cost + reg_cost
26     #print("cost", cost)
27
28     return cost

```

For the SGD algorithm, we just used a simple version of SGD, where only learning rate was taken into consideration. The useful information like costs, parameters ect. was recorded in the dictionary for further calculation.

```

1  def sgd(self, weights, bias, X, y,
2      print_cost = False):
3      costs = []
4
5      for i in range(self.n_iter):
6
7          grads, cost = self.model_status(
8          weights, bias, X, y)

```

```

7          grad_weights = grads["grad_weights
8          "]
9          grad_bias = grads["grad_bias"]
10
11         #updates
12         weights = weights - self.lr *
13         grad_weights
14         bias = bias - self.lr * grad_bias
15
16         #record cost
17         costs.append(cost)
18
19         #print cost every 100 training
20         epoch
21         if print_cost:
22             print ("Cost after iteration %
23             i: %f" %(i, cost))
24
25         parameters = {
26             "weights": weights,
27             "bias": bias
28         }
29
30         gradients = {
31             "grad_weights": grad_weights,
32             "grad_bias": grad_bias
33         }
34
35         return parameters, gradients, costs

```

The training function is calling the initialization function to get weights and bias, the SGD function is to get optimized parameters using train data, and test the model on our test data.

```

1  def train(self, X_train, y_train, X_test,
2      y_test, print_cost = False, print_score =
3      False):
4      #initilize parameters with normal
5      distribution
6      weights, bias = self.wb_init(X_train.
7      shape[0])
8
9      #update parameters with sgd algorithm
10     parameters, gradients, costs = self.
11     sgd(weights, bias, X_train, y_train,
12     print_cost)
13     self.weights = parameters["weights"]
14     self.bias = parameters["bias"]
15
16     #make predictions
17     y_pred_test = self.predict(X_test)
18     y_pred_train = self.predict(X_train)
19     #print(y_pred_train, y_train)
20
21     train_acc = self.accuracy(y_pred_train,
22     y_train)
23     test_acc = self.accuracy(y_pred_test,
24     y_test)
25
26     if print_score:
27         print ("train accuracy: %f" %(
28         train_acc))
29         print ("test accuracy: %f" %(
30         test_acc))
31
32     output = {
33         "costs": costs,
34         "data_cost": self.data_cost,
35         "reg_cost": self.reg_cost,
36         "y_pred_train": y_pred_train,
37         "y_pred_test": y_pred_test,
38         "weights": self.weights,
39         "bias": self.bias,

```



```

29         "learning rate": self.lr,
30         "train accuracy": train_acc,
31         "test accuracy": test_acc
32     }
33
34     return output

```

The evaluation of logistic model is based the theory part of logistic regression, and calculated as percent and programmed as follow:

```

1     def accuracy(self, y_pred, y):
2         accuracy = 100 - np.mean(np.abs(y_pred
3             - y)) * 100
4         return accuracy

```

### 3.3 Neural Network

The aim of this part is to write a FFNN and we start by writing some code for the simplest case where we have one hidden layer to get a better understanding of how NN works. We start by setting up our biases and weights, where the weights are normally distributed and we add a small number 0.01, so if they are zero they get a value which is nonzero (we might get complications if we get a divide by zero). Bellow is a snip-it of the code for the defining of the weights and an biases.

```

1     def crt_b_w(self):
2         # weights and bias in our hidden
3         #Note addind +0.01 so that if we have
4         zero its changed to a low value
5         self.h_weights = 0.01 + np.random.
6         normal(self.features, self.hidden_neurons)
7         #with normal distribution
8         self.h_bias = np.zeros(self.
9             hidden_neurons)
10
11         # weights and bias in our output
12         self.out_weights = 0.01 + np.random.
13         normal(self.hidden_neurons, self.
14             categories)
15         self.out_bias = np.zeros(self.
16             categories)

```

For the back propagation for this the simple case we update the weights and biases in the hidden and output layer. Where the output layer is updated with the activation function and error in the hidden layer for the weights and the error in the output for the biases. The hidden layer is updated for the weights with the  $X^T$  and the error in the hidden layer and the bias is the error in the hidden layer. We also update both the weights depending on if  $\lambda > 0$ , where we multiply  $\lambda$  with the weights and sum them. The last step is to update the weights and biases in the hidden and output layer with  $\eta$  which is our learning rate 5 and then subtract that from the original weights and biases in both layers. We then have the code:

```

1     def backprop(self):
2         self.error_in_out = self.probability -
3             self.Y_full

```

```

3         self.error_in_hidden = np.matmul(self.
4             error_in_out, self.out_weights.T) * self.
5             activation_hidden*(1-self.
6                 activation_hidden)
7
8         self.grad_weight_out = np.matmul(self.
9             activation_hidden.T, self.error_in_hidden)
10        self.grad_bias_out = np.sum(self.
11            error_in_out, axis=0)
12
13        self.grad_weight_hidden = np.matmul(
14            self.X_full.T, self.error_in_hidden)
15        self.grad_bias_hidden = np.sum(self.
16            error_in_hidden, axis=0)
17
18        if self.lmbda > 0:
19            self.grad_weight_out += self.lmbda
20            * self.out_weights
21            self.grad_weight_hidden += self.
22                lmbda * self.h_weights
23
24        self.out_weights -= self.eta*self.
25            grad_weight_out
26        self.out_bias -= self.eta*self.
27            grad_bias_out
28
29        self.h_weights -= self.eta*self.
30            grad_weight_hidden
31        self.h_bias -= self.eta*self.
32            grad_bias_hidden

```

Then we have the feed forward method which is quiet simply taken from then theory in equation 23 and added to code:

```

1     def ff(self):
2         #Feed forward for network saved
3         globally in class
4         self.z_hidden = np.matmul(self.X_full,
5             self.h_weights) + self.h_bias
6
7         activation_hidden = self.
8             activation_func_hidden(self.z_hidden)
9
10        self.z_out = np.matmul(
11            activation_hidden, self.out_weights) +
12            self.out_bias
13        self.a_expect = self.
14            activation_func_out(self.z_out)
15        self.probability = self.a_expect/np.
16            sum(self.a_expect, axis=1, keepdim=True)

```

Finally we add a step to train the function which will take random data points and run both the feed forward and back propagation. We only need to do this one to get the parameters and then we can simply run just for the outputs to get the error, probability and so on.

```

1     def train_function(self):
2         indec = np.arange(self.inputs)
3
4         for i in range(self.epochs):
5             for l in range(self.iter):
6                 data_points=np.random.choice(
7                     indec, size=self.batch_sz, replace=False)

```

```

8         self.X_full = self.X_full[
data_points]
9         self.Y_full = self.Y_full[
data_points]
10
11         self.ff()
12         self.backprop()

```

After looking at this simple case of a FFNN we needed to make it generalized so we could change the number of wanted hidden layers. We didn't know where to start, but referenced to the book "Neural Networks from Scratch in Python" (Kinsley, H. & Kukiela, D. 2020. [4]) and the lecture notes [1], which are excellent at describing how to set up a NN and for it to have a customizable number of hidden layers. For this we changed the method quite a bit. We now need to add a forward, backward and a predict to all the activation functions so that for each hidden layer we adjust for the activation function. For the layer where we define the weights and biases we also need a forward and backward function to update them accordingly. We start by doing the same as for the simple case, we initialize the weights and biases the same way, but add a regularization to use in the backward function later. The implementation we ended up using is as follows:

```

1 class Layer():
2     def __init__(self, n_features, neurons,
3         lmbd = 0):
4         # Initialize weights and biases
5         self.n_features = n_features
6         self.neurons = neurons
7         self.weights = 0.01 * np.random.randn(
8             self.n_features, self.neurons) #with
9             normal distribution
10        self.bias = np.zeros((1,self.neurons))
11
12        # Set strength of regularization, the
13        # regularizer should be greater than or
14        # equal to 0
15        self.weight_regularizer_l2 = lmbd
16
17    def forward(self, inputs, train):
18        self.input = inputs
19
20        #Calculate output value from previous
21        #layer's inputs, weights and biases
22        self.out = np.dot(inputs, self.weights
23        ) + self.bias
24
25    def back(self, d_val):
26        # gradients on weights and biases
27        self.grad_weights = np.dot(self.input.
28        T, d_val)
29        self.grad_bias = np.sum(d_val, axis=0,
30        keepdims=True)
31
32        #L2 regularization on weights
33        if self.weight_regularizer_l2 > 0:
34            self.grad_weights += 2 * self.
35            weight_regularizer_l2 * self.weights
36
37        #Gradient on inputs
38        self.grad_input = np.dot(d_val, self.
39        weights.T)

```

For the activation functions we can look at one example, lets look at the Sigmoid to accommodate for different hidden layers we have to calculate for the forward, backward and the predicted value. Almost every part has this method since we need to update the for each of the hidden layers back and fourth. The forward will just be the sigmoid function of the input, but the backward is *evaluated value \* (1 - output) \* output*. Lastly the predicted part is  $(output > 0.5) * 1$  and the function as code looks like this (we do this for each of the different activation functions):

```

1 class Activ_Sigmoid():
2     def forward(self, inputs, train):
3         self.input = inputs
4
5         self.out = 1/(1 + np.exp(-inputs))
6     def back(self, d_val):
7         self.grad_input = d_val * (1 - self.
8         out) * self.out
9
10    def predict(self, out):
11        return (out > 0.5) * 1

```

One of the biggest changes is that we split the different parts into their own Classes so that we can utilize them freely as we want and input them with the different layers. Another change is that we now use attributes to classify if we have different parameters or weights for example, by classifying the different inputs of the classes. After changing the code to the more advanced version we also add a method part where we run through the different classes and do the train and finalization. The training is done similarly to the part above we run through the epochs and train our model depending on the activation function chosen and the weights and biases, our implementation of the train is as follows (just a small part of it):

```

1     def train(self, X, y, *, n_epoc = 1,
2         validation_data = None, print_epoch =
3         False):
4
5         self.accuracy.init(y)
6
7         for epochs in range(1, n_epoc+1):
8
9             out = self.forward(X, train=True)
10
11             loss_dat, reg_loss = self.loss.cal
12             (out, y, regularization = True)
13             loss = loss_dat + reg_loss
14
15             predict = self.activ_out.predict(
16             out) ???
17             accuracy = self.accuracy.cal(
18             predict, y)
19
20             self.back(out, y)
21
22             self.optimiz.pre_up_par()
23             for layer in self.tlayer:
24                 self.optimiz.up_par(layer)
25             self.optimiz.post_up_par()

```

We then have the finalization part which runs through the layers by first taking the first layer where  $i == 0$  and set the previous and the next object, then we run through for  $i < \text{layercount} - 1$  for the previous and the next object and lastly save the last objects. Then we check if layer has attribute "weights" and if it does we add it to the trainable layers (we don't need to check for biases) and update the loss object for the trainable layers. We then have the code:

```

1      for i in range(layer_iter):
2
3          if i == 0:
4              self.layer[i].prev = self.
5              layer_inp
6              self.layer[i].next = self.
7              layer[i+1]
8
9          elif i < layer_iter - 1:
10             self.layer[i].prev = self.
11             layer[i-1]
12             self.layer[i].next = self.
13             layer[i+1]
14
15         else:
16             self.layer[i].prev = self.
17             layer[i-1]
18             self.layer[i].next = self.loss
19             self.activ_out = self.layer[i]
20
21         if hasattr(self.layer[i], "weights
22         "):
23             self.tlayer.append(self.layer[
24             i])
25
26         # Update loss object with trainable
27         layers
28         self.loss.remember_training_layer(self
29         .tlayer)

```

## 4 Results

### 4.1 Stochastic Gradient Descent

With our SGD we tested the implementation by changing the inversion calculation from project of the OLS and Ridge regression. We have the resulting images for the SGD in figure 4 and 5 and for the results from project 1 with python's inversion we have the figures 6 and 7.

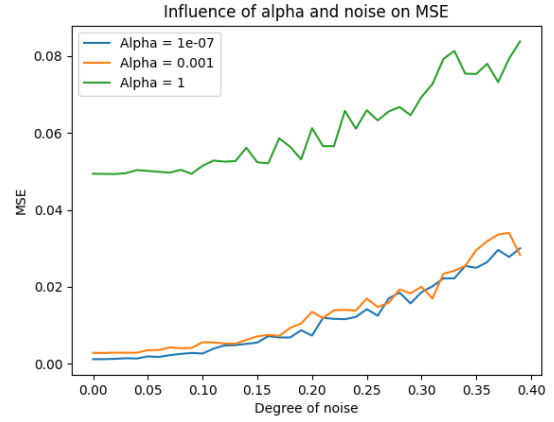


Figure 4: Show the influence of alpha and noise on our MSE with SGD.

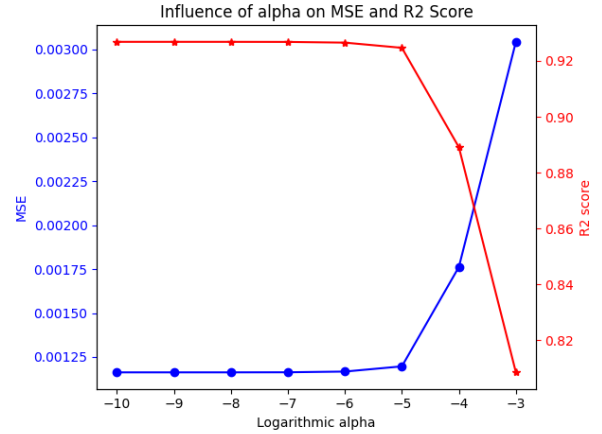


Figure 5: Show the influence of  $\alpha$  on MSE Vs. R2 score with SGD.

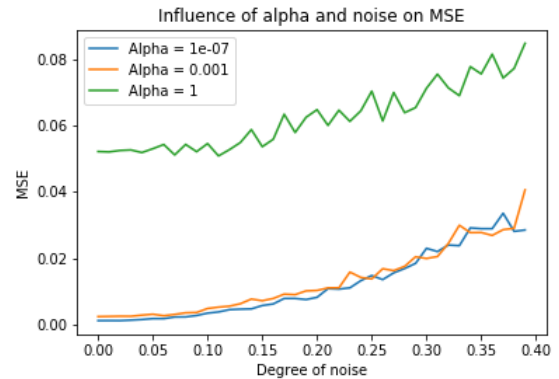


Figure 6: Show the influence of alpha and noise on our MSE with python inversion.

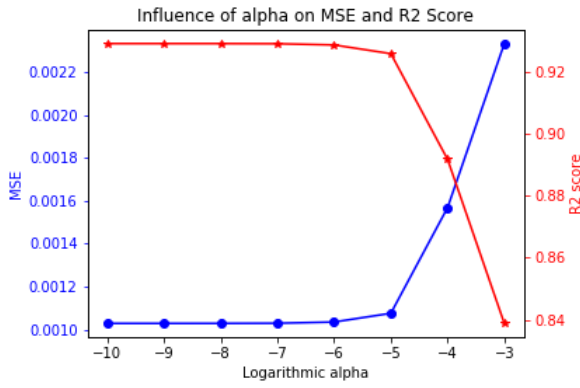


Figure 7: Show the influence of  $\alpha$  on MSE Vs. R2 score with python inversion.

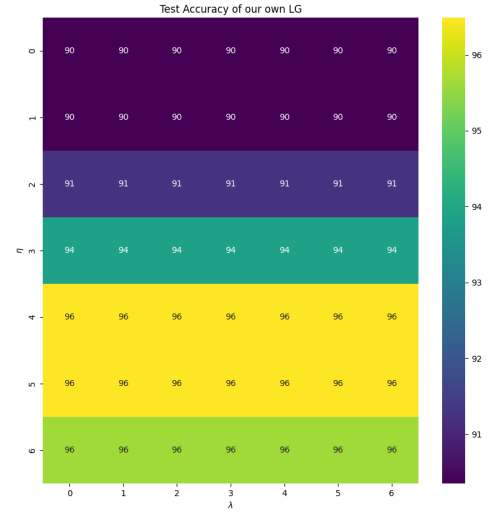


Figure 9: Show the test accuracy of our own LG.

## 4.2 Logistic Regression

Our logistic regression code was tested with Wisconsin breast cancer data set. The results of training cost was plotted against the increasing epochs (figure 8). In addition, the grid search method was used to find the best combination of the learning rate  $\eta$  and regularization parameter  $\lambda$  with respect to accuracy on training data (figure 9) and accuracy on test data (figure 10).

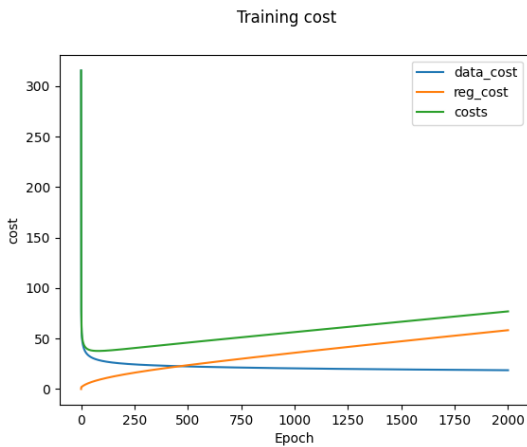


Figure 8: Show the training cost.

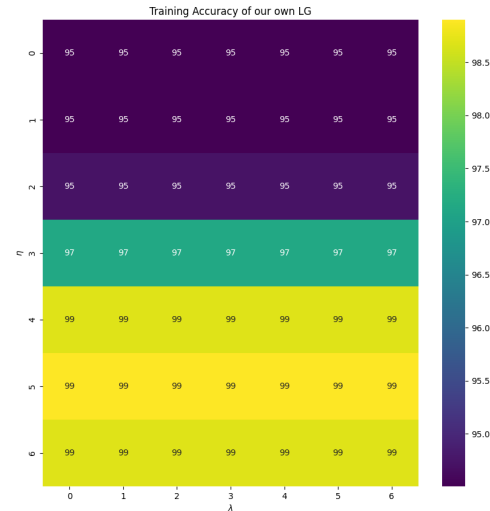


Figure 10: Show the train accuracy of our own LG.

The SGD classifier from sci-kit learn package was used to compare the performance of our own SGD code. The SGD optimizer is the only part that has been replaced with. The similar grid search method is employed as well shown in figure 11 and figure 12.

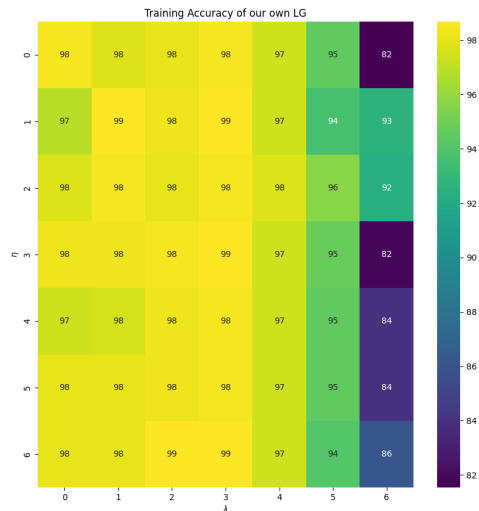


Figure 11: Show the train accuracy of our own LG with SGD classifier from sklearn package.

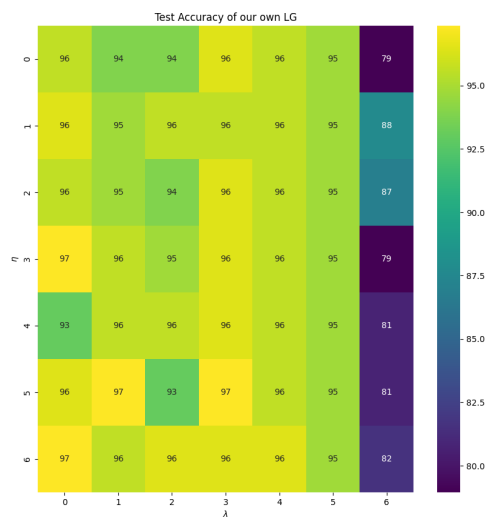


Figure 12: Show the train accuracy of our own LG with SGD classifier from sklearn package.

### 4.3 FFNN with different activation functions

We firstly looked at the influence of number of neurons in the hidden layer on train and test accuracy and result is shown in figure 13.

racy of train and test dataset with different number of neurons in the hidden

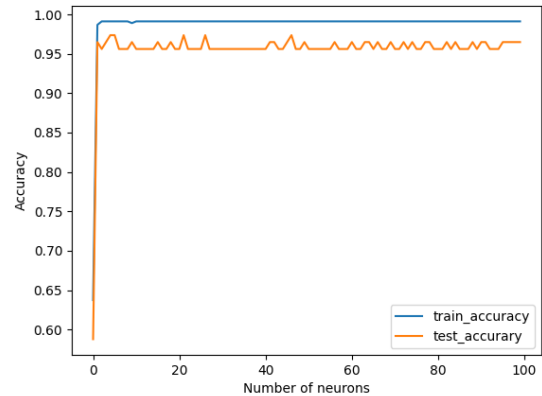


Figure 13: Show the train and test accuracy with different number of hidden neurons.

We tested our FFNN for different the activation functions and made a heat map for the accuracy of the different functions where we plotted both the test accuracy and the train accuracy to see the difference. Firstly we tested the Sigmoid function which is shown in figure 14 and 15 for the test and train. With the heat map it is easy to see the difference before and after and we can see that the accuracy in train is much higher than that of the test.

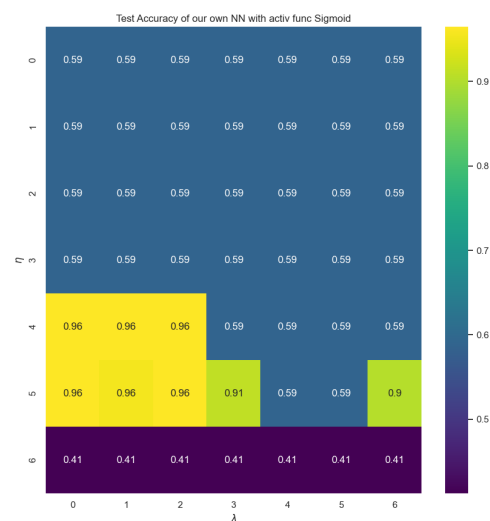


Figure 14: Show the test accuracy of the Sigmoid activation function in then hidden layers.

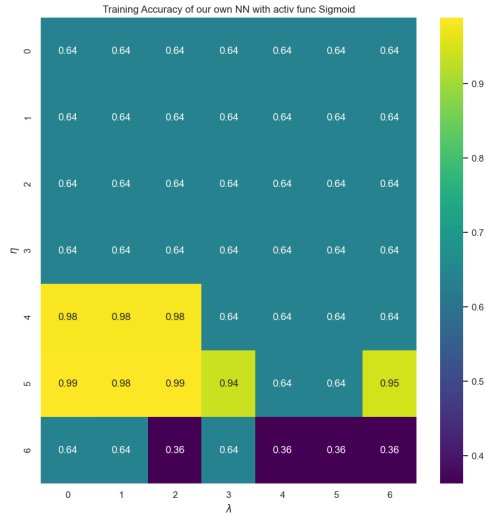


Figure 15: Show the train accuracy of the Sigmoid activation function in then hidden layers.

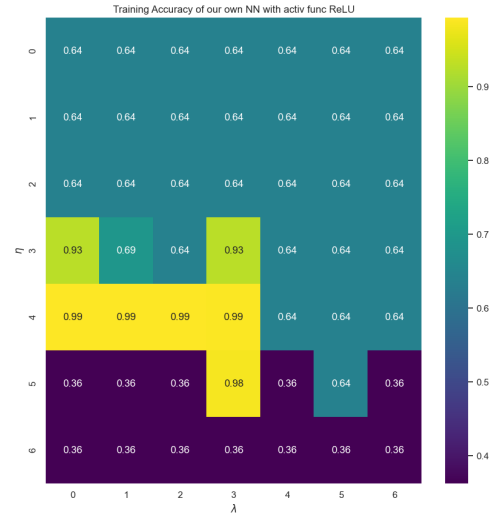


Figure 17: Show the train accuracy of the ReLU activation function in then hidden layers.

Then we tested for the ReLU activation function, which is presented in figure 16 and 17. Where we have an increase in accuracy from the test to the train. We can see that for both the function tested we have a pattern for the accuracy.

We then tested the activation function Leaky ReLU and we have the resulting figures 18 and 19. Where we have quite low accuracy in both cases and it does not increase more than 5% at best and it actually decreases for bigger  $\eta$ .

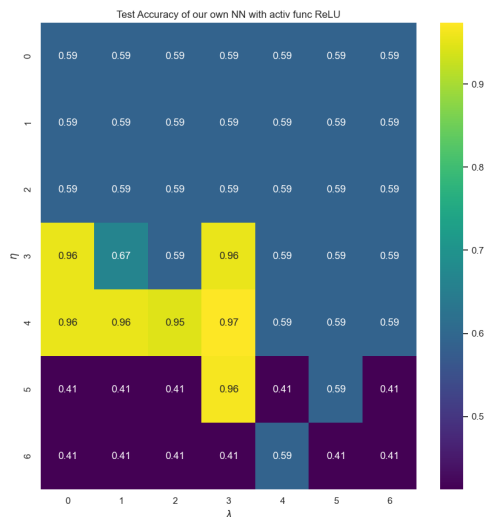


Figure 16: Show the test accuracy of the ReLU activation function in then hidden layers.

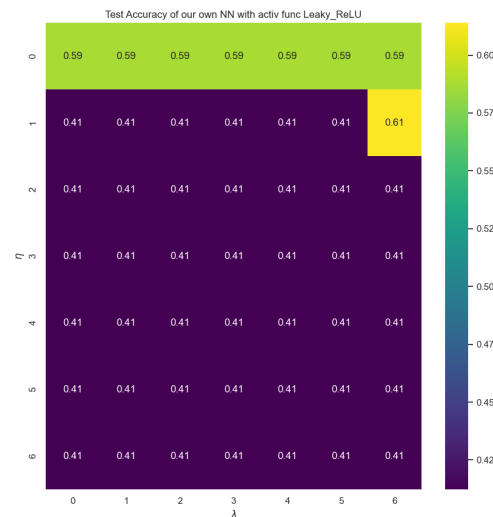


Figure 18: Show the test accuracy of the Leaky ReLU activation function in then hidden layers.

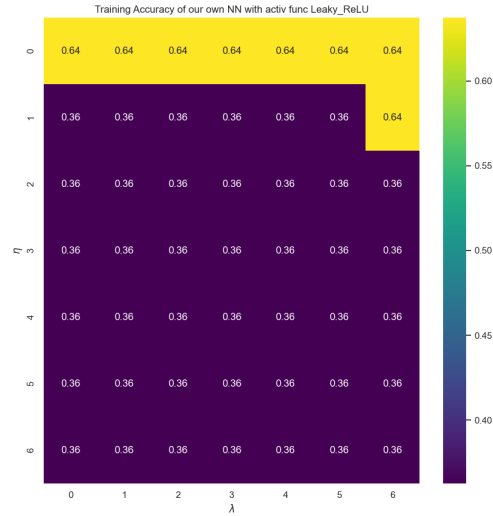


Figure 19: Show the train accuracy of the Leaky ReLU activation function in then hidden layers.

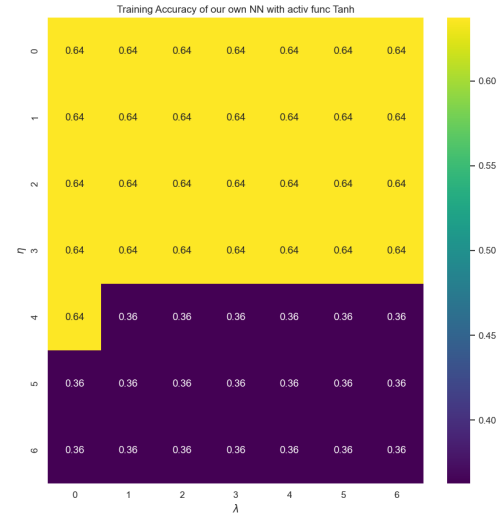


Figure 21: Show the train accuracy of the tanh activation function in then hidden layers.

We ended up testing for two more activation functions the tanh and a linear function. For the tanh we have the resulting figures 20 and 21. This activation function gave low accuracy's for all values of  $\eta$  and  $\lambda$ , it is better for  $\eta$  from  $10^{-5}$  to  $10^{-1}$ . Then for the linear function we have the resulting images 22 and 23. With this activation function we have a wider verity of parameters to choose from and have optimal accuracy's.

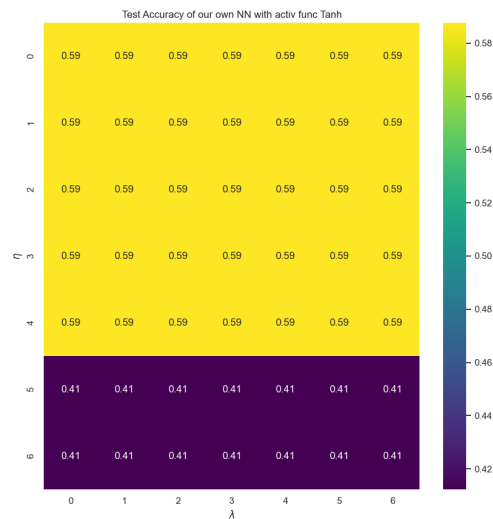


Figure 20: Show the test accuracy of the tanh activation function in then hidden layers.

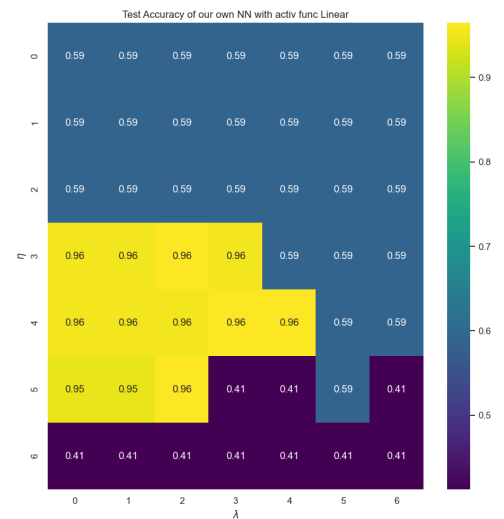


Figure 22: Show the test accuracy with a linear activation function in then hidden layers.

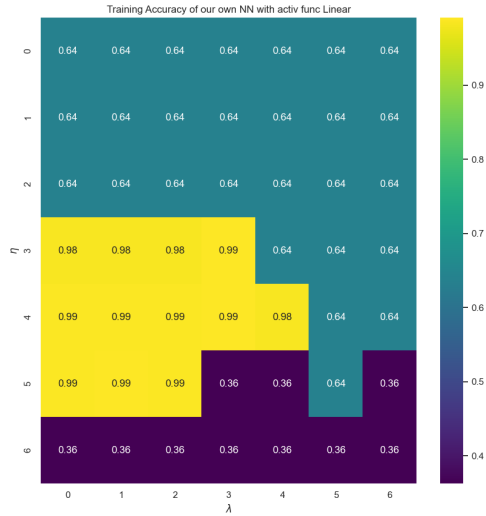


Figure 23: Show the train accuracy with a linear activation function in then hidden layers.

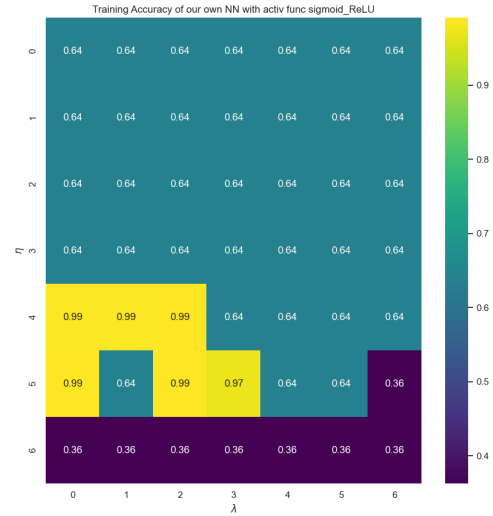


Figure 24: Show the train accuracy with a Sigmoid and ReLU activation functions

The accuracy for all of the activation functions except the tanh function had a final accuracy from 96% to 98% which is very high and acceptable result for the output of the NN.

We used softmax for all of the results above, but we wanted to see what happen if we changed this up and used for example the ReLU or the tanh function instead. Simply changing the functions gave us the resulting figures for Sigmoid and ReLU 24 and for the Sigmoid and tanh functions we have figure 25.

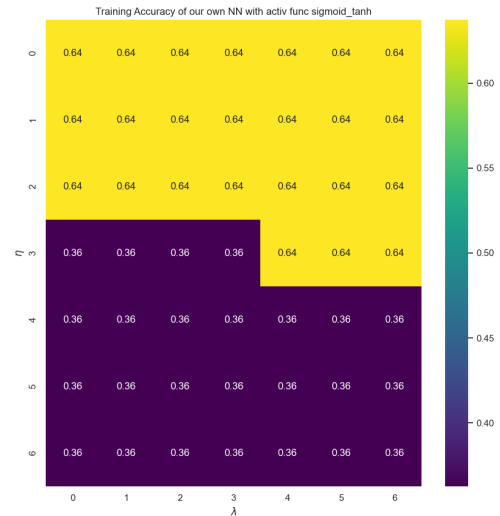


Figure 25: Show the train accuracy with a Sigmoid and tanh activation functions.

#### 4.4 NN with Keras

We also wanted to test our code with that of one of the already existing library's NN and chose Keras. We can simply just send our data to the library to get the same results as our own code. Where we tested for one activation function and chose to test with the Sigmoid function. From this we have the figures 26 and 27. Where we can see that we have a much larger area



for choosing  $\eta$  and  $\lambda$ , but the accuracy is still as good as what we got from our FFNN.

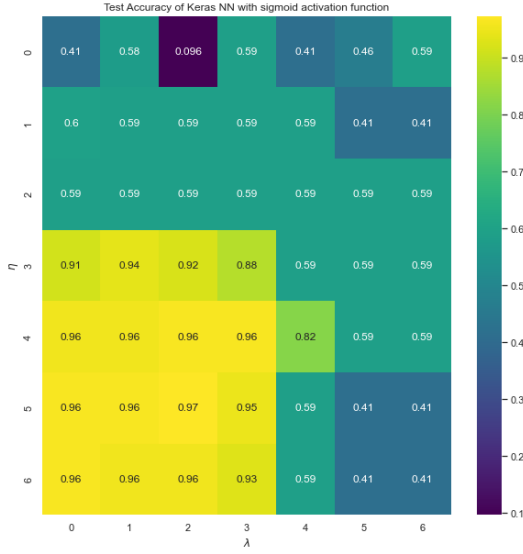


Figure 26: Show the test accuracy with Keras with Sigmoid as the activation function.

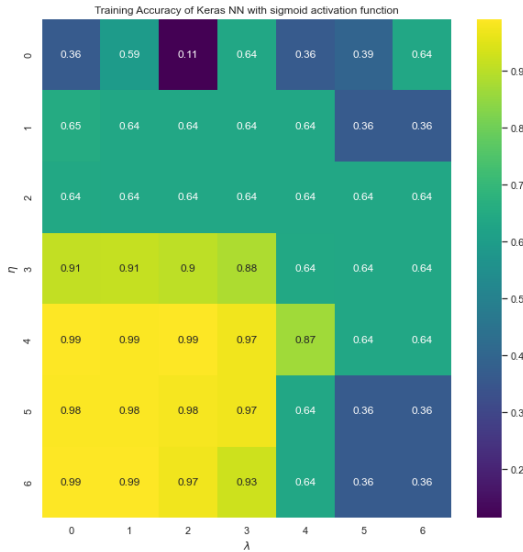


Figure 27: Show the train accuracy with Keras with Sigmoid as the activation function.

## 5 Discussion

### 5.1 Stochastic Gradient Descent

From the resulting images we can clearly see that changing the inversion calculation with our own SGD there is not a huge difference. We have that the MSE and R2 score staid almost the same and the difference is negligible, the same can be said for the influence of

noise and  $\alpha$ , but the most interesting part is the MSE and R2 score.

### 5.2 Logistic Regression

The training cost from figure 8 shows how the cost from our data set and regularization cost change with the increasing epochs. However, the regularization cost becomes higher and higher with the increasing epoch, which we expected to be plateau from a certain epoch. In addition, the train and test accuracy upon different learning rate  $\eta$  and regularization parameter  $\lambda$  indicates that the model is sensitive to the learning rate, but not responsive to the regularization parameter. Therefore, we speculate that the SGD algorithm may have some defect, resulting in the continuous changing of parameters for every epoch.

In order to validate if the SGD algorithm has problems, we chose to use the SGD classifier from sci-kit learn package in replacement of our own SGD algorithm used in logistic regression. In the resulting SGD classifier it seems that the regularization parameter  $\lambda$  is able to influence the accuracy on both train and test data.

Ignoring of the regularization parameter, we could get the best performance with learning rate of 0.1 and 1.

### 5.3 FFNN with different activation functions

The number of neurons does not seem to influence both train and test accuracy that much as long as the number of neurons are not too small. This means that we could improve the runtime hugely without sacrificing the performance by selecting the proper number of neurons in the hidden layer.

The different activation functions differed from each other and we can see that for the three first, we have some difference in where they are optimal for the different values of the hyper-parameters  $\eta$  and  $\lambda$ . We got the biggest area for these parameters by using the linear function as the activation function. Surprisingly the ReLU function preformed much better than the Leaky ReLU. We would expect the leaky version to preform much better, since we don't have to face the "dying ReLU" problem and could catch some characteristics when we have values for ReLU which are always under zero. We would expect that the backward would preform on the learning rate, but it preforms worse. From the Sigmoid function we can see that it works well for our data where we have either "True" or "False", but it is not a zero-centric function which might give us some problems. This function is also quite computationally expensive because of its exponential in its nature, this is the

same for tanh. The ReLU on the other hand is not as computationally expensive.

We also changed the softmax function with both ReLU and tanh to see if they would perform better. The resulting images shows that they did change the outcome, but not as one would expect. We might expect that the change would impact each other positively where they have low accuracy's, whereas making it perform inadequately instead.

## 5.4 NN with Keras

When we tested the Keras package the resulting heat maps where we have a large area for the choosing the hyper-parameters  $\eta$  and  $\lambda$ . There are a lot of values that give an accuracy of over 90%. The difference in Keras and our own code is not huge if we look at the Sigmoid results we have a lot of similarities, but just a smaller area for the  $\eta$  and  $\lambda$  values. Then our code performs well, but we have some improvements if we want to make it even better. In the Keras training results there is an anomaly for  $\eta = 10^{-5}$  and  $\lambda = 10^{-3}$ , where we have an exceptionally small accuracy. This might come from the model breaking down with the Sigmoid function and resulting in an abnormally small accuracy which is not good.

## 6 Conclusion

The resulting accuracy's from our model and that of already well made implementations we can see that we get quite good accuracy's as long as we are careful about choosing the right hyper parameters. We get an accuracy well above 90%. The advantages of the Sigmoid function is that it performs well in its bounds and overall not terrible outside this area, the same can be said for the ReLU and the linear functions. Although the linear function might lose some of the data and miss some of the information. The worse functions where the Leaky ReLU and tanh, which performed so bad for this NN that we should not use them for this application. By testing different activation function in the output layer we also saw that the best performing function was the softmax, where it gave consistent results. The softmax is normally used for where we have an output in multi-class classification problems.

The implementation of logistic regression on the same data got quite similar performance, with right hyper parameters. They both can reach about 97% accuracy on test data and more than 99% on train data. The runtime of both models is highly dependent on the number of iteration. It would be interesting to test the runtime for both models with their best combination of hyper parameters for the further steps of

implementing these two models. The best activation functions for classification cases will be the softmax and sigmoid and for regression the linear would work best.

## A Appendix A CODE:

*The code and results can be found following this link to our [GitHub](#)*

### A.1 Stochastic Gradient Decent

```
1 import numpy as npn
2 import autograd.numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.linear_model import SGDRegressor
5 from autograd import grad
6 from autograd import elementwise_grad as egrad # for functions that vectorize over inputs
7 from autograd import holomorphic_grad as hgrad
8 from sklearn import linear_model
9 #from AnalysisFunctions import *
10 from sklearn.preprocessing import PolynomialFeatures
11 import matplotlib.pyplot as plt
12 from matplotlib import cm
13 from imageio import imread
14
15
16 """
17 Analasys functions from project 1
18 """
19
20 def FrankeFunction(x,y, noise = 0):
21     term1 = 0.75*np.exp(-(0.25*(9*x-2)**2) - 0.25*((9*y-2)**2))
22     term2 = 0.75*np.exp(-((9*x+1)**2)/49.0 - 0.1*(9*y+1))
23     term3 = 0.5*np.exp(-(9*x-7)**2/4.0 - 0.25*((9*y-3)**2))
24     term4 = -0.2*np.exp(-(9*x-4)**2 - (9*y-7)**2)
25     return (term1 + term2 + term3 + term4 + noise*np.random.randn(len(x)))
26
27
28 def R2(zReal, zPredicted):
29     """
30     :param zReal: actual z-values, size (n, 1)
31     :param zPredicted: predicted z-values, size (n, 1)
32     :return: R2-score
33     """
34     R2 = 1 - (np.sum((zReal - zPredicted)**2)/np.sum((zReal - np.mean(zReal))**2))
35     return R2
36
37 def MeanSquaredError(zReal, zPredicted):
38     """
39     :param zReal: actual z-values, size (n, 1)
40     :param zPredicted: predicted z-values, size (n, 1)
41     :return: Mean squared error
42     """
43     MSE = np.sum((zReal - zPredicted)**2)/len(z)
44     return MSE
45
46 def betaCI_OLS(zReal, beta_mean, X):
47     """
48     :param zReal: actual z-values, size (n, 1)
49     :param beta_mean: mean of beta
50     :param X: dataset
51     Compute a 90% confidence interval for the beta coefficients
52     """
53
54     # Calculate variance squared in the error
55     z_hat = X.dot(beta_mean)
56     N, P = np.shape(X)
57     sigma2 = (np.sum(np.power((zReal-z_hat), 2)))/N
58
59     # Calculate the variance squared of the beta coefficients
60     var_beta = np.diag(sigma2*np.linalg.inv((X.T.dot(X))))
61
62     # The square root of var_beta is the standard error. Confidence intervals are calculated as mean
63     # +/- Z*SE
64     ci_minus = beta_mean - 1.645*var_beta
65     ci_plus = beta_mean + 1.645*var_beta
```

```

65     return ci_minus, ci_plus
66
67
68
69 def betaCI_Ridge(zReal, beta_mean, X, 1):
70     """
71     :param zReal: actual z-values, size (n, 1)
72     :param beta_mean: mean of beta
73     :param X: dataset
74     Compute a 90% confidence interval for the beta coefficients - Ridge
75     """
76
77     # Calculate variance squared in the error
78     z_hat = X.dot(beta)
79     N, P = np.shape(X)
80     sigma_2 = (np.sum(np.power((zReal-z_hat), 2)))/N
81
82     # Calculate the variance squared of the beta coefficients
83     XTX= X.T.dot(X)
84     R, R = np.shape(XTX)
85     var_beta = np.diag(sigma_2*np.linalg.inv((XTX + 1*np.identity(R))))
86
87     # The square root of var_beta is the standard error. Confidence intervals are calculated as mean
88     +/- Z*SE
89     ci_minus = beta_mean - 1.645*var_beta
90     ci_plus = beta_mean + 1.645*var_beta
91
92     return ci_minus, ci_plus
93
94 def plotFrankes(x_, y_, z_):
95     """
96     Plot Franke's function
97     """
98
99     fig = plt.figure()
100     ax = fig.gca(projection='3d')
101
102     surf = ax.plot_surface(x_, y_, z_, cmap=cm.coolwarm,
103                           linewidth=0, antialiased=False)
104
105     # Customize the z axis.
106     ax.set_zlim(-0.10, 1.40)
107     ax.zaxis.set_major_locator(LinearLocator(10))
108     ax.zaxis.set_major_formatter(FormatStrFormatter('%0.2f'))
109
110     ax.set_xlabel('X')
111     ax.set_ylabel('Y')
112     ax.set_zlabel('Z - Franke')
113
114     # Add a color bar which maps values to colors.
115     clb = fig.colorbar(surf, shrink=0.5, aspect=5)
116     clb.ax.set_title('Level')
117
118     plt.show()
119
120 #Ordinary Least Squared function
121 def ols(x, y, z, degree = 5):
122     #x: vector of size(n, 1)
123     #y: vector of size(n,1)
124     #z: vector of size(n,1)
125     xyb_ = np.c_[x, y]
126     poly = PolynomialFeatures(degree)
127     xyb = poly.fit_transform(xyb_)
128
129     #Change from inverse to SGD
130     #print("SGD=",np.shape(StochasticGradientDecent(x = (xyb.T.dot(xyb)), y=z).SGD()))
131     #beta = (StochasticGradientDecent(x = (xyb.T.dot(xyb)), y=z).SGD()).dot(xyb.T).dot(z)
132     #beta = StochasticGradientDecent(x = xyb, y=z).SGD().dot(xyb.T).dot(z)
133     beta = StochasticGradientDecent(x, y).SGD().dot(xyb.T).dot(z)
134     #beta = np.linalg.inv(xyb.T.dot(xyb)).dot(xyb.T).dot(z)
135
136     return beta

```

```

137
138 def RidgeRegression(x, y, z, degree=5, l=0.0001):
139     """
140     :param x: numpy vector of size (n, 1)
141     :param y: numpy vector of size (n, 1)
142     :param degree: degree of polynomial fit
143     :param l: Ridge penalty coefficient
144     :return: numpy array with the beta coefficients
145     """
146     # Calculate matrix with x, y - polynomials
147     M_ = np.c_[x, y]
148     poly = PolynomialFeatures(degree)
149     M = poly.fit_transform(M_)
150
151     # Calculate beta
152     A = np.arange(1, degree + 2)
153     rows = np.sum(A)
154     #Change from inverse to SGD
155     #beta = (StochasticGradientDecent(x = (M.T.dot(M) + l * np.identity(rows)), y=y).SGD()).dot(M.T)
156     #.dot(z)
157     beta = (StochasticGradientDecent(x = (M.T.dot(M) + l * np.identity(rows)), y=y).SGD()).dot(M.T)
158     .dot(z)
159
160     return beta
161
162 def Lasso(x, y, z, degree=5, a=1e-06):
163
164     X = np.c_[x, y]
165     poly = PolynomialFeatures(degree=degree)
166     X_ = poly.fit_transform(X)
167
168     clf = linear_model.Lasso(alpha=a, max_iter=5000, fit_intercept=False)
169     clf.fit(X_, z)
170     beta = clf.coef_
171
172     return beta
173
174 def bootstrap(x, y, z, p_degree, method, n_bootstrap=100):
175     # Randomly shuffle data
176     data_set = np.c_[x, y, z]
177     np.random.shuffle(data_set)
178     set_size = round(len(x)/5)
179
180     # Extract test-set, never used in training. About 1/5 of total data
181     x_test = data_set[0:set_size, 0]
182     y_test = data_set[0:set_size, 1]
183     z_test = data_set[0:set_size, 2]
184     test_indices = np.linspace(0, set_size-1, set_size)
185
186     # And define the training set as the rest of the data
187     x_train = np.delete(data_set[:, 0], test_indices)
188     y_train = np.delete(data_set[:, 1], test_indices)
189     z_train = np.delete(data_set[:, 2], test_indices)
190
191     Z_predict = []
192
193     MSE = []
194     R2s = []
195     for i in range(n_bootstrap):
196         x_, y_, z_ = resample(x_train, y_train, z_train)
197
198         if method == 'Ridge':
199             # Ridge regression, save beta values
200             beta = RidgeRegression(x_, y_, z_, degree=p_degree)
201         elif method == 'Lasso':
202             beta = Lasso(x_, y_, z_, degree=p_degree)
203         elif method == 'OLS':
204             beta = ols(x_, y_, z_, degree=p_degree)
205         else:
206             print('ERROR: Cannot recognize method')
207             return 0
208
209     M_ = np.c_[x_test, y_test]

```

```

208     poly = PolynomialFeatures(p_degree)
209     M = poly.fit_transform(M_)
210     z_hat = M.dot(beta)
211
212     Z_predict.append(z_hat)
213
214     # Calculate MSE
215     MSE.append(np.mean((z_test - z_hat)**2))
216     R2s.append(R2(z_test, z_hat))
217
218     # Calculate MSE, Bias and Variance
219     MSE_M = np.mean(MSE)
220     R2_M = np.mean(R2s)
221     bias = np.mean((z_test - np.mean(Z_predict, axis=0, keepdims=True))**2)
222     variance = np.mean(np.var(Z_predict, axis=0, keepdims=True))
223     return MSE_M, R2_M, bias, variance
224
225
226
227 class StochasticGradientDecent(object):
228     """docstring for StochasticGradientDecent."""
229
230     def __init__(self, x, y, n_epoc = 50, M = 10, n=1000, gamma=0.9, dtype = "float64"):
231
232         self.x_full = x
233         self.y_full = y
234         self.n_epoc = n_epoc
235         #size of each minibatch
236         self.M = M
237         self.n = n
238         self.gamma = gamma
239         #Some initial conditions
240         #number of minibatch
241         self.m = int(self.n/self.M)
242         self.t0 = self.M
243         self.t1 = self.n_epoc
244         #theta dimension is based on the number of columns in design matrix
245         self.v_ = 0
246
247     def __call__(self):
248
249         #Checks matrix size of rows
250         size_matrix = self.x_full.shape[0]
251         if size_matrix != self.y_full.shape[0]:
252             raise ValueError("'x' and 'y' must have same rows")
253
254         #Check to see if batches are right size
255         self.n_epoc = int(self.n_epoc)
256         if not 0 < self.n_epoc <= size_matrix:
257             raise ValueError("Must have a batch size less or equal to observations and greater than
258 zero.")
259
260         #Checks gamma is in range
261         if not 0 <= self.gamma <= 1:
262             raise ValueError("Gamma must be equal or greater than zero and equal or less than 1.")
263
264         #Checks gamma is in range
265         if not 0 <= self.gamma <= 1:
266             raise ValueError("Gamma must be equal or greater than zero and equal or less than 1.")
267
268         #gradient
269         def gradient(self, x, y, theta):
270             return (2.0/self.M)*x.T @ ((x @ theta) - y).T
271
272         #This the learning scheduel for eta
273         def ls(self, t):
274             return self.t0/(t+self.t1)
275
276         #The eta values function
277         def eta(self, t):
278             return self.t0**2/(t+self.t1)
279

```

```

280 def SGD(self):
281     X = np.c_[np.ones((len(self.x_full),1)), self.x_full]
282     #size_matrix = x.shape[0]
283     self.theta = np.random.randn(X.shape[1],1) #Initilize theta for matrix shape.
284
285     #xy = np.c_[x.reshape(size_matrix, -1), y.reshape(size_matrix, 1)]
286
287     #Main SGD loop for epochs of minibatches
288     for epoc in range(self.n_epoc):
289         #Second SGD loop with random choice of k
290         for k in range(self.m):
291             random_index = self.M*np.random.randint(self.m)
292             xi = X[random_index:random_index+self.M]
293             yi = self.y_full[random_index:random_index+self.M]
294
295             eta = self.ls(epoc*self.m+k) #Calling function to cal. eta
296
297
298             #self.v_ = gamma*self.v_ + eta*gradient(x_iter, y_iter, self.theta - gamma*self.v_)
299             #Cal. v where gradient is from autograd
300             place_hold = self.theta + self.gamma*self.v_
301             x_grad = egrad(self.gradient, 2) #Gradient with respect to theta
302             self.v_ = self.gamma*self.v_ + eta * x_grad(xi, yi, place_hold) #Cal. v where
303             gradient is from autograd, self.gradient(xi, yi, self.theta)
304             self.theta = self.theta - self.v_ #Theta +1 from this itteration of theta and v
305
306         return self.theta
307
308 """
309 Analysis of a Lasso Regression model of Franke's function from project 1
310 """
311
312 """
313 # Load data
314 X = np.load('data.npy')
315 x = X[:, 0]
316 y = X[:, 1]
317 z = FrankeFunction(x, y)
318
319 alphas = [10**-10, 10**-9, 10**-8, 10**-7, 10**-6, 10**-5, 10**-4, 10**-3]
320 alpha_logs = [-10, -9, -8, -7, -6, -5, -4, -3]
321
322 Bs = []
323 for al in alphas:
324     Bs.append(Lasso(x, y, z, degree=5, a=al))
325
326 # Generate new test data
327 x_test = np.random.rand(200)
328 y_test = np.random.rand(200)
329 z_test = FrankeFunction(x_test, y_test)
330
331 M_ = np.c_[x_test, y_test]
332 poly = PolynomialFeatures(5)
333 M = poly.fit_transform(M_)
334 MSEs = []
335 R2s = []
336 text_file = open("../Results/SGD/Bootstrap_lasso_SGD.txt", "w")
337 for i in range(len(alphas)):
338     z_predict = M.dot(Bs[i])
339     MSE = MeanSquaredError(z_test, z_predict)
340     MSEs.append(MSE)
341     R2_score = R2(z_test, z_predict)
342     R2s.append(R2_score)
343     text_file.write('--- Alpha value: {0} ---\n Mean Squared error: {1:.7f} \n R2 Score: {2:.7f}\n'.
344                     format(alphas[i], MSE, R2_score))
345
346 # make plot
347 fig, ax1 = plt.subplots()
348 ax1.plot(alpha_logs, MSEs, 'bo-')
349 ax1.set_xlabel('Logarithmic alpha')
350 # Make the y-axis label, ticks and tick labels match the line color.

```

```

350 ax1.set_ylabel('MSE', color='b')
351 ax1.tick_params('y', colors='b')
352
353 ax2 = ax1.twinx()
354 ax2.plot(alpha_logs, R2s, 'r*-')
355 ax2.set_ylabel('R2 score', color='r')
356 ax2.tick_params('y', colors='r')
357
358 plt.title('Influence of alpha on MSE and R2 Score')
359 fig.tight_layout()
360 plt.savefig('../Results/SGD/MSE_R2_alpha_SGD.png')
361
362 # Investigate how the alpha values are influenced by noise
363 noise = np.arange(0, 0.4, 0.01)
364 alphas = [10**-7, 10**-3, 1]
365 Bs = []
366
367 # Generate more data to test
368 x_test = np.random.rand(200)
369 y_test = np.random.rand(200)
370 M_ = np.c_[x_test, y_test]
371 poly5 = PolynomialFeatures(5)
372 M = poly5.fit_transform(M_)
373
374 for al in alphas:
375     B = []
376     #print(al)
377     for n in noise:
378         z = FrankeFunction(x, y, noise=n)
379         B.append(Lasso(x, y, z, degree=5, a=al))
380     Bs.append(B)
381
382 lines = []
383 plt.figure()
384 for i in range(len(alphas)):
385     text_file.write('--- alpha value: {} --- \n'.format(alphas[i]))
386     line = []
387     for j in range(len(noise)):
388         z_test = FrankeFunction(x_test, y_test, noise=noise[j])
389         z_predict = M.dot(Bs[i][j])
390         MSE = MeanSquaredError(z_test, z_predict)
391         line.append(MSE)
392         R2_score = R2(z_test, z_predict)
393         text_file.write(' Noise: {0} \n Mean Squared error: {1:.7f} \n R2 Score: {2:.7f}\n'.format(
noise[j], MSE, R2_score))
394     plt.plot(noise, line, label='Alpha = {0}'.format(alphas[i]))
395
396 plt.legend()
397 plt.xlabel('Degree of noise')
398 plt.ylabel('MSE')
399 plt.title('Influence of alpha and noise on MSE')
400 plt.savefig('../Results/SGD/alpha_noise_MSE_SGD.png')
401
402 MSE_1, R2_1, bias_1, variance_1 = bootstrap(x, y, z, method='Lasso', p_degree=5)
403 text_file.write('--- BOOTSTRAP --- \n')
404 text_file.write('MSE: {} \n'.format(MSE_1))
405 text_file.write('R2: {} \n'.format(R2_1))
406 text_file.write('Bias: {} \n'.format(bias_1))
407 text_file.write('Variance: {} \n'.format(variance_1))
408
409 text_file.close()
410 """
411
412
413
414 """
415 Part 6 from project 1
416 """
417
418 # Load the terrain
419 terrain1 = imread('SRTM_data_Norway_2.tif')
420 # Show the terrain
421 plt.figure()

```



```

422 plt.title('Terrain area')
423 plt.imshow(terrain1, cmap='gray')
424 plt.xlabel('X')
425 plt.ylabel('Y')
426 plt.show()
427
428 # Choose a smaller part of the data set
429 terrain = terrain1[500:750, 0:250]
430 # Show the terrain
431 plt.figure()
432 plt.imshow(terrain, cmap='gray')
433 plt.xlabel('X')
434 plt.ylabel('Y')
435 plt.savefig('../Results/SGD/terrain_original.png')
436
437 # Make zero matrix to later fit data
438 num_rows, num_cols = np.shape(terrain)
439 num_observations = num_rows * num_cols
440 X = np.zeros((num_observations, 3))
441
442 # make a matrix with all the values from the data on the form [x y z]
443 index = 0
444 #X = X - np.mean(X)
445 for i in range(0, num_rows):
446     for j in range(0, num_cols):
447         X[index, 0] = i # x
448         X[index, 1] = j # y
449         X[index, 2] = terrain[i, j] # z
450         index += 1
451
452 # OLS example
453 # extract x, y, z
454 xt = X[:,0, np.newaxis]
455 yt = X[:,1, np.newaxis]
456 zt = X[:,2, np.newaxis]
457
458
459
460
461 degree = [2, 4, 6, 8]
462 text_file = open("../Results/SGD/terrain_CI_ols.txt", "w")
463 for d in degree:
464     beta = ols(xt, yt, zt, degree=d)
465
466     M_ = np.c_[xt, yt]
467     poly = PolynomialFeatures(d)
468     M = poly.fit_transform(M_)
469     z_predict = M.dot(beta)
470
471
472     T = np.zeros([num_rows, num_cols])
473     index = 0
474     # create matrix for imshow
475     for i in range(0, num_rows):
476         for j in range(0, num_cols):
477             T[i, j] = (z_predict[index])
478             index += 1
479
480     plt.figure()
481     plt.imshow(T, cmap='gray')
482     plt.xlabel('X')
483     plt.ylabel('Y')
484     plt.savefig('../Results/SGD/terrain_ols_d{}.png'.format(d))
485
486
487     z_test = np.zeros(zt.shape[0])
488     for i in range(zt.shape[0]):
489         z_test[i] = zt[i][0]
490
491     beta_test = np.zeros(beta.shape[0])
492     for i in range(beta.shape[0]):
493         beta_test[i] = zt[i][0]
494

```

```

495     conf1, conf2 = betaCI_OLS(z_test, beta_test, M)
496     #print(conf2.shape)
497     for i in range(len(conf1)):
498         text_file.write('Beta {0}: {1:5f} & [{2:5f}, {3:5f}] \n'.format(i, beta_test[i], conf1[i],
499                               conf2[i]))
500 text_file.close()
501
502 # Evaluate model with bootstrap algorithm
503 text_file = open("../Results/SGD/terrain_mse_ols.txt", "w")
504 MSE, R2, bias, variance = bootstrap(xt, yt, zt, p_degree=8, method='OLS', n_bootstrap=100)
505 text_file.write('MSE: {0:5f} & R2: {1:5f} & bias: {2:5f} & var: {3:5f}'.format(MSE, R2, bias,
506                               variance))
507 text_file.close()
508 #####
509 ##Ridge Regression
510 #####
511 text_file = open("../Results/SGD/terrain_CI_ridge.txt", "w")
512 for d in degree:
513     beta = RidgeRegression(xt, yt, zt, degree=d)
514
515     M_ = np.c_[xt, yt]
516     poly = PolynomialFeatures(d)
517     M = poly.fit_transform(M_)
518     z_predict = M.dot(beta)
519
520     T = np.zeros([num_rows, num_cols])
521     index = 0
522     # create matrix for imshow
523     for i in range(0, num_rows):
524         for j in range(0, num_cols):
525             T[i, j] = (z_predict[index])
526             index += 1
527
528     plt.figure()
529     plt.imshow(T, cmap='gray')
530     plt.xlabel('X')
531     plt.ylabel('Y')
532     plt.savefig('../Results/SGD/terrain_ridge_d{}.png'.format(d))
533
534     z_test = np.zeros(zt.shape[0])
535     for i in range(zt.shape[0]):
536         z_test[i] = zt[i][0]
537
538     beta_test = np.zeros(beta.shape[0])
539     for i in range(beta.shape[0]):
540         beta_test[i] = zt[i][0]
541
542     conf1, conf2 = betaCI_OLS(z_test, beta_test, M)
543     #print(conf2.shape)
544     for i in range(len(conf1)):
545         text_file.write('Beta {0}: {1:5f} & [{2:5f}, {3:5f}] \n'.format(i, beta_test[i], conf1[i],
546                               conf2[i]))
547
548 text_file.close()
549
550 # Evaluate model with bootstrap algorithm
551 text_file = open("../Results/SGD/terrain_mse_ridge.txt", "w")
552 MSE, R2, bias, variance = bootstrap(xt, yt, zt, p_degree=8, method='Ridge', n_bootstrap=100)
553 text_file.write('MSE: {0:5f} & R2: {1:5f} & bias: {2:5f} & var: {3:5f}'.format(MSE, R2, bias,
554                               variance))
555 text_file.close()
556 #####
557 ##Lasso Regression
558 #####
559 text_file = open("../Results/SGD/terrain_CI_lasso.txt", "w")
560 for d in degree:
561     beta = Lasso(xt, yt, zt, degree=d)
562
563

```

```

564 M_ = np.c_[xt, yt]
565 poly = PolynomialFeatures(d)
566 M = poly.fit_transform(M_)
567 z_predict = M.dot(beta)
568
569
570 T = np.zeros([num_rows, num_cols])
571 index = 0
572 # create matrix for imshow
573 for i in range(0, num_rows):
574     for j in range(0, num_cols):
575         T[i, j] = (z_predict[index])
576         index += 1
577 plt.figure()
578 plt.imshow(T, cmap='gray')
579 plt.xlabel('X')
580 plt.ylabel('Y')
581 plt.savefig('../Results/SGD/terrain_lasso_d{}.png'.format(d))
582
583
584 z_test = np.zeros(zt.shape[0])
585 for i in range(zt.shape[0]):
586     z_test[i] = zt[i][0]
587
588 beta_test = np.zeros(beta.shape[0])
589 for i in range(beta.shape[0]):
590     beta_test[i] = zt[i][0]
591
592
593 conf1, conf2 = betaCI_OLS(z_test, beta_test, M)
594 #print(conf2.shape)
595 for i in range(len(conf1)):
596     text_file.write('Beta {0}: {1:5f} & [{2:5f}, {3:5f}] \n'.format(i, beta_test[i], conf1[i],
597 conf2[i]))
598 text_file.close()
599
600 # Evaluate model with bootstrap algorithm
601 text_file = open("../Results/SGD/terrain_mse_lasso.txt", "w")
602 MSE_l, R2_l, bias_l, variance_l = bootstrap(xt, yt, zt, method='Lasso', p_degree=8)
603 text_file.write('MSE: {0:5f} & R2: {1:5f} & bias: {2:5f} & var: {3:5f}'.format(MSE_l, R2_l, bias_l,
604 variance_l))
605 text_file.close()
606
607
608 """
609 n = 1000
610 x = 2*np.random.rand(n,1)
611 y = 4+3*x+np.random.randn(n,1)
612 X = np.c_[np.ones((n,1)), x]
613
614 #print(X)
615
616 if __name__ == "__main__":
617     theta = StochasticGradientDecent(x, y).SGD()
618     print(theta)
619 """

```

## A.2 FFNN simple case

```

1 import autograd.numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import SGDRegressor
4 from autograd import grad
5 from autograd import elementwise_grad as egrad # for functions that vectorize over inputs
6 from sklearn import linear_model
7 from SGD import StochasticGradientDecent
8
9
10 """Feed Forward Neural Network"""
11 class NeuralNetwork(object):
12     def __init__(self, x, y, hidden_neurons=50, categories= 10, n_epochs=10 ,batch_sz = 100, lambda =

```

```

0.001, activation_func_hidden = "sigmoid", activation_func_out = "tanh"):
13
14     #Setting up initial conditions for class
15     self.X_full = x
16     self.Y_full = y
17
18     self.input = x.shape[0]
19     self.feauters = x.shape[1]
20     self.hidden_neurons = hidden_neurons
21     self.categories = categories
22     self.n_epochs = n_epochs
23     self.lmbda = lmbda
24     self.batch_sz = batch_sz
25     self.iter = self.input // self.batch_sz
26     self.eta = (self.n_epochs/2)/(self.n_epochs/2+self.n_epochs) #Learning scheduel
27
28     #Initilize theta from SGD
29     self.theta = StochasticGradientDecent(x, y, n_epoc = self.n_epochs, M = self.categories, n=
np.size(self.input), gamma=0.3).SGD()
30
31     #Allows other activation function for hidden layer
32     if activation_func_hidden == "sigmoid":
33         self.activation_func_hidden = self.sigmoid
34     #Allows other activation function for output layer
35     elif activation_func_out == "tanh":
36         self.activation_func_out = self.tanh
37     elif activation_func_out == "relu":
38         self.activation_func_out = self.re
39
40     #Creating bias and weight by running function
41     self.crt_b_w()
42
43     def __call__(self):
44
45         #Checks matrix size of rows
46         size_matrix = self.x_full.shape[0]
47         if size_matrix != self.y_full.shape[0]:
48             raise ValueError("'x' and 'y' must have same rows")
49
50         #Small tests to check input
51         size_matrix = X.shape[0]
52         self.n_epochs = int(self.n_epoch)
53         if not 0 < self.n_epoch <= size_matrix:
54             raise ValueError("Must have a 'epochs' size less or equal to observations and greater
than zero")
55
56         self.batch_sz = int(self.batch_sz)
57         if self.batch_sz <= 0:
58             raise ValueError("'Batch size' must be greater than 0.")
59
60
61     def crt_b_w(self):
62         # weights and bias in our hidden
63         #Note addind +0.01 so that if we have zero its changed to a low value
64         self.h_weights = 0.01 + np.random.normal(self.features, self.hidden_neurons) #with normal
distribution
65         self.h_bias = np.zeros(self.hidden_neurons)
66
67         # weights and bias in our output
68         self.out_weights = 0.01 + np.random.normal(self.hidden_neurons, self.categories)
69         self.out_bias = np.zeros(self.categories)
70
71     #Sigmoid activation function
72     def sigmoid(self, x):
73         return 1/(1 + np.exp(-x))
74
75     #Tanh activation function
76     def tanh(self, x):
77         return (np.exp(x) - np.exp(-x)) / (np.exp(x) + np.exp(-x))
78
79     #ReLU activation function
80     def relu(self, x):
81

```

```

82     return np.maximum(0,x)
83
84 def cost_MSE(self y_h):
85     n = np.size(y)
86     C = 0
87     for i in range(n):
88         C += (self.Y_full[i] - y_h[i])**2
89     return 1/n * C
90
91 def ff(self):
92     #Feed forward for network saved globally in class
93     self.z_hidden = np.matmul(self.X_full, self.h_weights) + self.h_bias
94
95     activation_hidden = self.activation_func_hidden(self.z_hidden)
96
97     self.z_out = np.matmul(activation_hidden, self.out_weights) + self.out_bias
98     self.a_expect = self.activation_func_out(self.z_out)
99     self.probability = self.a_expect/np.sum(self.a_expect, axis=1, keepdim=True)
100
101
102 def ff_out(self):
103     #feed forward output saved locally in function
104     z_hidden = np.matmul(X, self.h_weights) + self.h_bias
105
106     activation_hidden = self.activation_func_hidden(z_hidden)
107
108     z_out = np.matmul(activation_hidden, self.out_weights) + self.out_bias
109     a_expect = self.activation_func_out(z_out)
110     probability = a_expect/np.sum(a_expect, axis=1, keepdim=True)
111     return probability
112
113 def backprop(self):
114     self.error_in_out = self.probability - self.Y_full
115
116     self.error_in_hidden = np.matmul(self.error_in_out, self.out_weights.T) * self.
activation_hidden*(1-self.activation_hidden)
117
118     self.grad_weight_out = np.matmul(self.activation_hidden.T, self.error_in_hidden)
119     self.grad_bias_out = np.sum(self.error_in_out, axis=0)
120
121
122     self.grad_weight_hidden = np.matmul(self.X_full.T, self.error_in_hidden)
123     self.grad_bias_hidden = np.sum(self.error_in_hidden, axis=0)
124
125     if self.lmbda > 0:
126         self.grad_weight_out += self.lmbda * self.out_weights
127         self.grad_weight_hidden += self.lmbda * self.h_weights
128
129     self.out_weights -= self.eta*self.grad_weight_out
130     self.out_bias -= self.eta*self.grad_bias_out
131
132     self.h_weights -= self.eta*self.grad_weight_hidden
133     self.h_bias -= self.eta*self.grad_bias_hidden
134
135 def predict(self, X):
136     return np.argmax(self.ff_out(X), axis=1)
137
138 def pred_prob(self, X):
139     return self.ff_out(X)
140
141 def train_function(self):
142     indec = np.arange(self.inputs)
143
144     for i in range(self.epochs):
145         for l in range(self.iter):
146             data_points=np.random.choice(indec, size=self.batch_sz, replace=False)
147
148             self.X_full = self.X_full[data_points]
149             self.Y_full = self.Y_full[data_points]
150
151             self.ff()
152             self.backprop()

```

## A.3 FFNN final code

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.datasets import load_breast_cancer
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6
7
8 #define Layers
9 """Initializing for weights an biases for flexible number of layers with number of inputs and amount
   of neurons or hidden layers."""
10 class Layer():
11     def __init__(self, n_features, neurons, lmbd = 0):
12         # Initialize weights and biases
13         self.n_features = n_features
14         self.neurons = neurons
15         self.weights = 0.01 * np.random.randn(self.n_features, self.neurons) #with normal
   distribution
16         self.bias = np.zeros((1,self.neurons))
17
18         # Set strength of regularization, the regularizer should be greater than or equal to 0
19         self.weight_regularizer_l2 = lmbd
20
21     def forward(self, inputs, train):
22         self.input = inputs
23
24         #Calculate output value from previous layer's inputs, weights and biases
25         self.out = np.dot(inputs, self.weights) + self.bias
26
27     def back(self, d_val):
28         # gradients on weights and biases
29         self.grad_weights = np.dot(self.input.T, d_val)
30         self.grad_bias = np.sum(d_val, axis=0, keepdims=True)
31
32         #L2 regularization on weights
33         if self.weight_regularizer_l2 > 0:
34             self.grad_weights += 2 * self.weight_regularizer_l2 * self.weights
35
36         #Gradient on inputs
37         self.grad_input = np.dot(d_val, self.weights.T)
38
39 class Layer_Input:
40     def forward(self, inputs, train):
41         self.out = inputs
42
43
44 '''Define activation function'''
45 class Activ_Sigmoid():
46     def forward(self, inputs, train):
47         self.input = inputs
48
49         self.out = 1/(1 + np.exp(-inputs))
50     def back(self, d_val):
51         self.grad_input = d_val * (1 - self.out) * self.out
52
53     def predict(self, out):
54         return (out > 0.5) * 1
55
56 class Activ_Leaky_ReLU():
57     def forward(self, inputs, train):
58         self.input = inputs
59
60         self.out = np.maximum(0.01 * inputs, inputs)
61
62     def back(self, d_val):
63
64         self.grad_input = d_val.copy()
65
66         #gradient will be 0.01 when inputs are negative
67         self.grad_input[self.input <= 0] = 0.01
68
```

```

69     def predict(self, out):
70         return out
71
72     class Activ_ReLU():
73         def forward(self, inputs, train):
74             self.input = inputs
75
76             self.out = np.maximum(0, inputs)
77
78         def back(self, d_val):
79
80             self.grad_input = d_val.copy()
81
82             #gradient will be 0 when inputs are negative
83             self.grad_input[self.input <= 0] = 0
84
85         def predict(self, out):
86             return out
87
88     class Activ_Linear():
89         def forward(self, inputs, train):
90             self.input = inputs
91             self.out = inputs
92
93         def back(self, d_val):
94             self.grad_input = d_val.copy()
95
96         def predict(self, out):
97             return out
98
99     class Activ_tanh():
100         def forward(self, inputs, train):
101             self.input = inputs
102
103             self.out = (np.exp(inputs) - np.exp(-inputs)) / (np.exp(inputs) + np.exp(-inputs))
104
105         def back(self, d_val):
106             self.grad_input = d_val.copy()
107
108             self.grad_input = 1-self.out**2
109
110         def predict(self, out):
111             return out
112
113     class Activ_Softmax():
114
115         def forward(self, inputs, train):
116             self.input = inputs
117
118             # Get unnormalized probabilities
119             exp_val = np.exp(inputs - np.max(inputs, axis=1, keepdims=True))
120
121             # Normalize them for each sample
122             probabilities = exp_val / np.sum(exp_val, axis=1, keepdims=True)
123
124             self.out = probabilities
125
126         def back(self, d_val):
127
128             # Create uninitialized array
129             self.grad_input = np.empty_like(d_val)
130
131             # Enumerate outputs and gradients
132             for index, (single_out, single_d_val) in enumerate(zip(self.out, d_val)):
133                 # Flatten output array
134                 single_out = single_out.reshape(-1, 1)
135                 # Calculate Jacobian matrix of the output
136                 jacobian_matrix = np.diagflat(single_out) - np.dot(single_out, single_out.T)
137                 # Calculate sample-wise gradient
138                 # and add it to the array of sample gradients
139                 self.grad_input[index] = np.dot(jacobian_matrix, single_d_val)
140
141         def predict(self, out):

```

```

142         return np.argmax(out, axis=1)
143
144     '''Define Optimizer'''
145     class Optim_SGD():
146
147         def __init__(self, momentum = 0, decay = 0, lr = 1):
148             self.momentum = momentum
149             self.lr = lr
150             self.decay = decay
151
152             self.cur_lr = lr
153             self.iter = 0
154
155         def pre_up_par(self):
156             if self.decay:
157                 self.cur_lr = self.lr * (1 / (1 + self.decay * self.iter))
158
159         def up_par(self, layer):
160             if self.momentum:
161                 #If layer dont consist of momentum we create them
162                 if not hasattr(layer, "weight_momentum"):
163                     layer.weight_momentum = np.zeros_like(layer.weights)
164                     layer.bias_momentum = np.zeros_like(layer.bias)
165
166                 #Creating weight updates
167                 weight_up = self.momentum * layer.weight_momentum - self.cur_lr * layer.grad_weights
168                 layer.weight_momentum = weight_up
169
170                 #Creating bias updates
171                 bias_up = self.momentum * layer.bias_momentum - self.cur_lr * layer.grad_bias
172                 layer.bias_momentum = bias_up
173
174             else:
175                 weight_up = -self.cur_lr * layer.grad_weights
176
177                 bias_up = -self.cur_lr * layer.grad_bias
178
179
180             layer.weights += weight_up
181             layer.bias += bias_up
182
183         def post_up_par(self):
184             self.iter += 1
185
186     '''Define Loss'''
187     class Loss:
188
189         def reg_loss(self):
190             #initialize regularization loss
191             reg_loss = 0
192
193             #calculate regularization loss for each the training layer
194             for layer in self.training_layer:
195                 # L2 regularization - weights
196                 if layer.weight_regularizer_l2 > 0:
197                     reg_loss += layer.weight_regularizer_l2 * np.sum(layer.weights * layer.weights)
198
199             return reg_loss
200
201         def remember_training_layer(self, training_layer):
202             self.training_layer = training_layer
203
204         #calculate losses from data and regularization with model output and true values
205         def cal(self, out, y, *, regularization=False):
206             #loss from sample
207             samp_loss = self.forward(out, y)
208             #mean loss
209             loss_dat = np.mean(samp_loss)
210
211             if not regularization:
212                 return loss_dat
213
214             return loss_dat, self.reg_loss()

```



```

215
216
217 # Categorical Cross-entropy loss
218 class Loss_CC(Loss):
219
220     def forward(self, predict, y_real):
221         #size of each batch
222         sample = len(predict)
223
224         # Clip data to avoid denominator of 0
225         predict_clip = np.clip(predict, 1e-8, 1 - 1e-8)
226
227         #Probabilities for target values of categorical labels
228         if len(y_real.shape) == 1:
229             cc = predict_clip[range(sample), y_real]
230
231         elif len(y_real.shape) == 2:
232             cc = np.sum(predict_clip * y_real, axis = 1)
233
234         nll = -np.log(cc)
235         return nll
236
237     def back(self, d_val, y_real):
238
239         sample = len(d_val)
240         lab = len(d_val[0])
241
242         if len(y_real.shape) == 1:
243             y_real = np.eye(lab)[y_real]
244
245         #Claculate and normalize
246         self.grad_input = (-y_real/d_val)/sample
247
248 class Loss_MSE(Loss):
249
250     def forwar(self, predict, y_real):
251         samp_loss = np.mean((y_real - predict)**2, axis=-1)
252
253         return samp_loss
254
255     def backward(self, d_val, y_real):
256         sample = len(d_val)
257
258         out = len(d_val[0])
259
260         #Claculate and normalize
261         self.grad_input = (-2 * (y_real - d_val) / out) / sample
262
263 '''Define Accuracy'''
264 class Accuracy:
265
266     def cal(self, predict, y):
267         # Get comparison results
268         comparisons = self.compare(predict, y)
269         acc = np.mean(comparisons)
270
271         return acc
272
273     def cal_accum(self):
274
275         acc = self.accum_sum / self.accum_count
276
277         return acc
278
279 # Reset variables for accumulated accuracy
280     def reset_var(self):
281
282         self.accum_sum = 0
283         self.accum_count = 0
284
285
286 class Accuracy_Classification(Accuracy):
287

```

```

288 def __init__(self, *, binary=False):
289     # Binary mode?
290     self.binary = binary
291
292 def init(self, y):
293     pass
294
295 # Compares predictions to the ground truth values
296 def compare(self, predict, y):
297
298     if not self.binary and len(y.shape) == 2:
299         y = np.argmax(y, axis = 1)
300
301     return predict == y
302
303
304 class Accuracy_Regression(Accuracy):
305
306     def __init__(self, y):
307         self.precision = None
308
309     def init(self, y, reinit=False):
310         if self.precision is None or reinit:
311             self.precision = np.std(y) / 250
312
313 # Compares predictions to the ground truth values
314 def compare(self, predict, y):
315     return np.absolute(predict - y) < self.precision
316
317 '''Define Model'''
318 class Method():
319
320     def __init__(self):
321         self.layer = []
322         self.train_acc = None
323         self.test_acc = None
324
325     def add_to_list(self, layer):
326         self.layer.append(layer)
327
328     def set_param(self, *, loss, accuracy, optimiz):
329         self.optimiz = optimiz
330         self.accuracy = accuracy
331         self.loss = loss
332
333     def finall(self):
334
335         self.layer_inp = Layer_Input()
336         layer_iter = len(self.layer)
337         self.tlayer = []
338
339         for i in range(layer_iter):
340
341             if i == 0:
342                 self.layer[i].prev = self.layer_inp
343                 self.layer[i].next = self.layer[i+1]
344
345             elif i < layer_iter - 1:
346                 self.layer[i].prev = self.layer[i-1]
347                 self.layer[i].next = self.layer[i+1]
348
349             else:
350                 self.layer[i].prev = self.layer[i-1]
351                 self.layer[i].next = self.loss
352                 self.activ_out = self.layer[i]
353
354             if hasattr(self.layer[i], "weights"):
355                 self.tlayer.append(self.layer[i])
356
357 # Update loss object with trainable layers
358 self.loss.remember_training_layer(self.tlayer)
359
360

```

```

361 def train(self, X, y, *, n_epoc = 1, validation_data = None, print_epoch = False):
362
363     self.accuracy.init(y)
364
365     for epochs in range(1, n_epoc+1):
366
367         out = self.forward(X, train=True)
368
369         loss_dat, reg_loss = self.loss.cal(out, y, regularization = True)
370         loss = loss_dat + reg_loss
371
372         predict = self.activ_out.predict(out) #??
373         accuracy = self.accuracy.cal(predict, y)
374
375         self.back(out, y)
376
377         self.optimiz.pre_up_par()
378         for layer in self.tlayer:
379             self.optimiz.up_par(layer)
380         self.optimiz.post_up_par()
381
382         """
383         if print_epoch:
384             print(f'n_epoc: {epochs}, ' +
385                   f'accuracy: {accuracy:.3f}, ' +
386                   f'loss: {loss:.3f} (' +
387                   f'loss in data: {loss_dat:.3f}, ' +
388                   f'loss in regularization: {reg_loss:.3f}), ' +
389                   f'learningrate: {self.optimiz.cur_lr}')
390         """
391
392         if epochs == n_epoc:
393             self.train_acc = accuracy
394
395
396     if validation_data is not None:
397
398         X_val, y_val = validation_data
399
400         out = self.forward(X_val, train=False)
401
402         loss = self.loss.cal(out, y_val)
403
404         predict = self.activ_out.predict(out)
405
406         accuracy = self.accuracy.cal(predict, y_val)
407
408         self.test_acc = accuracy
409
410         if print_epoch:
411             print(f'validation, ' +
412                   f'acc: {accuracy:.3f}, ' +
413                   f'loss: {loss:.3f}')
414
415     return self.train_acc, self.test_acc
416
417
418 def forward(self, X, train):
419
420     self.layer_inp.forward(X, train)
421
422     for layer in self.layer:
423         layer.forward(layer.prev.out, train)
424
425     return layer.out
426
427 def back(self, out, y):
428
429     self.loss.back(out, y)
430
431     for layer in reversed(self.layer):
432         layer.back(layer.next.grad_input)
433

```

```

434 '''Test Wisconsin Cancer Data'''
435 data = load_breast_cancer()
436 X = data['data']
437 y = data['target']
438
439 #Train test split
440 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 0)
441
442 #Feature Scaling
443 sc = StandardScaler()
444 X_train = sc.fit_transform(X_train)
445 X_test = sc.transform(X_test)
446
447 n_neurons = 100
448 train_accuracy = []
449 test_accuracy = []
450 for n in range(n_neurons):
451
452     # Instantiate the model
453     model = Method()
454     # Add layers
455     model.add_to_list(Layer(30, n, lmbd=5e-4))
456     model.add_to_list(Activ_ReLU())
457     model.add_to_list(Layer(n, 2))
458     model.add_to_list(Activ_Softmax())
459
460     # Set loss, optimizer and accuracy objects
461     model.set_param(
462         loss=Loss_CC(),
463         optimiz=Optim_SGD(lr=0.05, decay=5e-5, momentum=0.9),
464         accuracy=Accuracy_Classification()
465     )
466
467     # Finalize the model
468     model.finall()
469
470     # Train the model
471     train_acc, test_acc = model.train(X_train, y_train, validation_data=(X_test, y_test), n_epoc
472                                     =200, print_epoch = False)
473     train_accuracy.append(train_acc)
474     test_accuracy.append(test_acc)
475
476 x = range(n_neurons)
477 fig = plt.figure()
478 plt.plot(x, train_accuracy, label = "train_accuracy")
479 plt.plot(x, test_accuracy, label = "test_accuracy")
480 fig.suptitle('Accuracy of train and test dataset with different number of neurons in the hidden
481             layer')
482 plt.xlabel('Number of neurons')
483 plt.ylabel('Accuracy')
484 plt.legend()
485 plt.savefig('../Results/acc_n_nruons_NN.png')

```

## A.4 Breast cancer data implementation code

```

1 from NN import *
2 import numpy as np
3 from sklearn.datasets import load_breast_cancer
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 import numba as nb
9 import timeit
10
11 # ensure the same random numbers appear every time
12 np.random.seed(100)
13
14 # load breast cancer data
15 data = load_breast_cancer()
16 X = data['data']
17 y = data['target']

```

```

18
19 #@nb.jit(nopython=True)
20 def run_NN(X, y, active = "ReLU"):
21
22     sns.set()
23     #train test split
24     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 0)
25
26     #Feature Scaling
27     sc = StandardScaler()
28     X_train = sc.fit_transform(X_train)
29     X_test = sc.transform(X_test)
30
31     #set up learning rate eta and regularization lambda for grid search
32     lr = np.logspace(-5, 1, 7)
33     lmbd = np.logspace(-5, 1, 7)
34
35     train_accuracy = np.zeros((len(lr), len(lmbd)))
36     test_accuracy = np.zeros((len(lr), len(lmbd)))
37
38     # grid search
39     for i, eta in enumerate(lr):
40         for j, lm in enumerate(lmbd):
41             # Instantiate the model
42             model = Method()
43
44             if active == "ReLU":
45                 # Add layers ReLU
46                 model.add_to_list(Layer(30, 64, lmbd=lm))
47                 model.add_to_list(Activ_ReLU())
48                 model.add_to_list(Layer(64, 64, lmbd=lm))
49                 model.add_to_list(Activ_ReLU())
50                 model.add_to_list(Layer(64, 2))
51                 model.add_to_list(Activ_Softmax())
52
53             elif active == "Sigmoid":
54                 # Add layers Sigmoid
55                 model.add_to_list(Layer(30, 64, lmbd=lm))
56                 model.add_to_list(Activ_Sigmoid())
57                 model.add_to_list(Layer(64, 64, lmbd=lm))
58                 model.add_to_list(Activ_Sigmoid())
59                 model.add_to_list(Layer(64, 2))
60                 model.add_to_list(Activ_Softmax())
61
62             elif active == "Leaky_ReLU":
63                 # Add layers Leaky_ReLU
64                 model.add_to_list(Layer(30, 64, lmbd=lm))
65                 model.add_to_list(Activ_Leaky_ReLU())
66                 model.add_to_list(Layer(64, 64, lmbd=lm))
67                 model.add_to_list(Activ_Leaky_ReLU())
68                 model.add_to_list(Layer(64, 2))
69                 model.add_to_list(Activ_Softmax())
70
71             elif active == "Linear":
72                 # Add layers Linear
73                 model.add_to_list(Layer(30, 64, lmbd=lm))
74                 model.add_to_list(Activ_Linear())
75                 model.add_to_list(Layer(64, 64, lmbd=lm))
76                 model.add_to_list(Activ_Linear())
77                 model.add_to_list(Layer(64, 2))
78                 model.add_to_list(Activ_Softmax())
79
80             elif active == "Tanh":
81                 # Add layers tanh
82                 model.add_to_list(Layer(30, 64, lmbd=lm))
83                 model.add_to_list(Activ_tanh())
84                 model.add_to_list(Layer(64, 64, lmbd=lm))
85                 model.add_to_list(Activ_tanh())
86                 model.add_to_list(Layer(64, 2))
87                 model.add_to_list(Activ_Softmax())
88
89             elif active == "sigmoid_tanh":
90                 # Add layers sigmoid and tanh

```

```

91         model.add_to_list(Layer(30, 64, lmbd=1m))
92         model.add_to_list(Activ_Sigmoid())
93         model.add_to_list(Layer(64, 64, lmbd=1m))
94         model.add_to_list(Activ_Sigmoid())
95         model.add_to_list(Layer(64, 2))
96         model.add_to_list(Activ_tanh())
97
98     elif active == "sigmoid_ReLU":
99         # Add layers sigmoid and ReLU
100        model.add_to_list(Layer(30, 64, lmbd=1m))
101        model.add_to_list(Activ_Sigmoid())
102        model.add_to_list(Layer(64, 64, lmbd=1m))
103        model.add_to_list(Activ_Sigmoid())
104        model.add_to_list(Layer(64, 2))
105        model.add_to_list(Activ_ReLU())
106
107
108    # Set loss, optimizer and accuracy objects
109    model.set_param(
110        loss=Loss_CC(),
111        optimiz=Optim_SGD(lr=eta, decay=5e-5, momentum=0.9),
112        accuracy=Accuracy_Classification()
113    )
114
115    # Finalize the model
116    model.finall()
117
118    # Train the model
119    train_acc, test_acc = model.train(X_train, y_train, validation_data=(X_test, y_test),
120    n_epoc=200, print_epoch = False)
121
122    train_accuracy[i][j] = train_acc
123    test_accuracy[i][j] = test_acc
124
125    fig, ax = plt.subplots(figsize = (10, 10))
126    sns.heatmap(train_accuracy, annot=True, ax=ax, cmap="viridis")
127    ax.set_title("Training Accuracy of our own NN with activ func " + active )
128    ax.set_ylabel("$\eta$")
129    ax.set_xlabel("$\lambda$")
130    plt.savefig('../Results/eta_lmd_train_acc_' + active + ".png")
131
132    fig, ax = plt.subplots(figsize = (10, 10))
133    sns.heatmap(test_accuracy, annot=True, ax=ax, cmap="viridis")
134    ax.set_title("Test Accuracy of our own NN with activ func " + active)
135    ax.set_ylabel("$\eta$")
136    ax.set_xlabel("$\lambda$")
137    plt.savefig('../Results/eta_lmd_test_acc_' + active + ".png")
138
139    return train_accuracy, test_accuracy
140
141
142 start = timeit.default_timer()
143 run_NN(X, y, active = "Tanh")
144 stop = timeit.default_timer()
145 print('Time: ', stop - start)
146
147
148
149
150 """
151 fig, ax = plt.subplots(figsize = (10, 10))
152 sns.heatmap(train_accuracy, annot=True, ax=ax, cmap="viridis")
153 ax.set_title("Training Accuracy of our own NN")
154 ax.set_ylabel("$\eta$")
155 ax.set_xlabel("$\lambda$")
156 plt.savefig('../Results/eta_lmd_train_acc.png')
157 fig, ax = plt.subplots(figsize = (10, 10))
158 sns.heatmap(test_accuracy, annot=True, ax=ax, cmap="viridis")
159 ax.set_title("Test Accuracy of our own NN")
160 ax.set_ylabel("$\eta$")
161 ax.set_xlabel("$\lambda$")
162 plt.savefig('../Results/eta_lmd_test_acc.png')

```

## A.5 Keras NN Cancer data code

```

1 from tensorflow.keras.layers import Input
2 from tensorflow.keras.models import Sequential          #This allows appending layers to existing models
3 from tensorflow.keras.layers import Dense              #This allows defining the characteristics of a
   particular layer
4 from tensorflow.keras import optimizers                #This allows using whichever optimiser we want (
   sgd,adam,RMSprop)
5 from tensorflow.keras import regularizers              #This allows using whichever regularizer we want
   (l1,l2,l1_l2)
6 from tensorflow.keras.utils import to_categorical      #This allows using categorical cross entropy as
   the cost function
7
8 from sklearn.model_selection import train_test_split
9 from sklearn.datasets import load_breast_cancer
10 from sklearn.preprocessing import StandardScaler
11 import numpy as np
12 from sklearn.metrics import accuracy_score
13 import matplotlib.pyplot as plt
14 import seaborn as sns
15 sns.set()
16
17 #Kera neural network was runned on conda environment
18 #The figure was downloaded and put in our github
19
20 # one-hot in numpy
21 def to_categorical_numpy(integer_vector):
22     n_inputs = len(integer_vector)
23     n_categories = np.max(integer_vector) + 1
24     onehot_vector = np.zeros((n_inputs, n_categories))
25     onehot_vector[range(n_inputs), integer_vector] = 1
26
27     return onehot_vector
28
29 def create_neural_network_keras(n_neurons_layer1, n_neurons_layer2, n_categories, eta, lmbd):
30     model = Sequential()
31     model.add(Dense(n_neurons_layer1, activation='sigmoid', kernel_regularizer=regularizers.l2(lmbd)
32 ))
33     model.add(Dense(n_neurons_layer2, activation='sigmoid', kernel_regularizer=regularizers.l2(lmbd)
34 ))
35     model.add(Dense(n_categories, activation='softmax'))
36
37     sgd = optimizers.SGD(lr=eta)
38     model.compile(loss='categorical_crossentropy', optimizer=sgd, metrics=['accuracy'])
39
40     return model
41
42 data = load_breast_cancer()
43 X = data['data']
44 y = data['target']
45 y = to_categorical(y)
46
47 X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size = 0.2, random_state = 0)
48 #Feature Scaling
49 sc = StandardScaler()
50 X_train = sc.fit_transform(X_train)
51 X_test = sc.transform(X_test)
52
53 DNN_keras = np.zeros((len(eta_vals), len(lmbd_vals)), dtype=object)
54
55 for i, eta in enumerate(eta_vals):
56     for j, lmbd in enumerate(lmbd_vals):
57         DNN = create_neural_network_keras(n_neurons_layer1, n_neurons_layer2, n_categories,
58                                           eta=eta, lmbd=lmbd)
59         DNN.fit(X_train, Y_train, epochs=epochs, batch_size=batch_size, verbose=0)
60         scores = DNN.evaluate(X_test, Y_test)
61
62         DNN_keras[i][j] = DNN
63
64     print("Learning rate = ", eta)

```

```
64 print("Lambda = ", lmbd)
65 print("Test accuracy: %.3f" % scores[1])
66 print()
```

## B Appendix B REFERENCES:

### References

- [1] Lectures and notes from lectures, <https://compphysics.github.io/MachineLearning/doc/web/course.html>
- [2] Wikipedia. (2021, November 11). Artificial Neural Network. Wikipedia. Retrieved November 19, 2021, from [https://en.wikipedia.org/wiki/Artificial\\_neural\\_network](https://en.wikipedia.org/wiki/Artificial_neural_network)
- [3] Kakaraparthi, V. (2019, October 29). Activation functions in neural networks. Medium. Retrieved November 20, 2021, from <https://prateekvishnu.medium.com/activation-functions-in-neural-networks-bf5c542d5fec>.
- [4] Kinsley, H., and Kukiela, D. (2020). Neural Networks From Scratch (Vol. 1). Harrison Kinsley.
- [5] Bushaev, V. (2017, December 5). Stochastic gradient descent with momentum. Medium. Retrieved November 20, 2021, from <https://towardsdatascience.com/stochastic-gradient-descent-with-momentum-a84097641a5d>.