



COMPUTER SCIENCE DEPARTMENT

Practical Deep Learning

Assignment 1 Report

Authors:

Almog Zemach: 205789001

Raz Monsonego: 313292120

May 2022

Contents

1	Part I: The classifier and optimizer	1
1.1	Gradient test for softmax regression	1
1.2	SGD for Least Squares	3
1.3	SGD for softmax	5
2	Part II: The neural network	9
2.1	Jacobian test - regular layer	9
2.2	ResNet Jacobian test	12
2.3	Neural network Gradient test	15
2.4	Learning of the whole neural network	17
2.4.1	Training	17
2.4.2	Swiss Roll data	19
2.4.3	Peaks Data	21
2.4.4	GMM Data	22
2.4.5	Learning with 200 data points	23
3	Conclusions	25
4	Auxiliary code	27

1 Part I: The classifier and optimizer

1.1 Gradient test for softmax regression

Initially, we implement the loss-function using the soft-max regression function, and its gradient with respect to the weights and biases (see code attached).

Then, we make sure we calculated the derivatives correctly by using the gradient-test. The code that performs the gradient test and the results are shown bellow:

Gradient test code:

```
1 def gradient_test(X: np.array, W: np.array, C: np.array, b: np.
2     array, policy):
3     """
4         Gradient test with respect for W.
5         :param X matrix.
6         :param W matrix.
7         :param C matrix.
8         :param b matrix.
9         :param policy. Indicates to which parameter we are doing the
10            test.
11        :return matplotlib graph which shows the gradiant test.
12        """
13    dict_num = {"W": 1, "X": 2, "b": 3}
14    dict_param = {"W": W, "X": X, "b": b}
15    dict_name = {"W": "\$\\delta W$", "X": "\$\\delta X$", "b": "\$\\delta b$"}
16    V = np.random.rand(dict_param[policy].shape[0], dict_param[
17        policy].shape[1])
18    d = (V / np.linalg.norm(V))
19    d_vector = d.reshape(-1, 1)
20    err_1 = []
21    err_2 = []
22    ks = []
23    params = soft_max_regression(X, W, C, b)
24    grad = params[dict_num[policy]].reshape(-1, 1)
25    for k in range(1, 20):
26        epsilon = 0.5 ** k
27        new = dict_param[policy] + epsilon * d
28        dict_args = {"W": (X, new, C, b), "X": (new, W, C, b), "b":
29            (X, W, C, new)}
30        f_x_d, _, _, _ = soft_max_regression(*dict_args[policy])
31        err_1.append(abs(f_x_d - params[0]))
32        err_2.append(abs(f_x_d - params[0] - (epsilon * d_vector.T
33            @ grad)[0][0]))
34        ks.append(k)
35    print_test(ks, err_1, err_2, "Gradiant Test: " + dict_name[
36        policy])
```

Gradient test results (for the soft-max function):

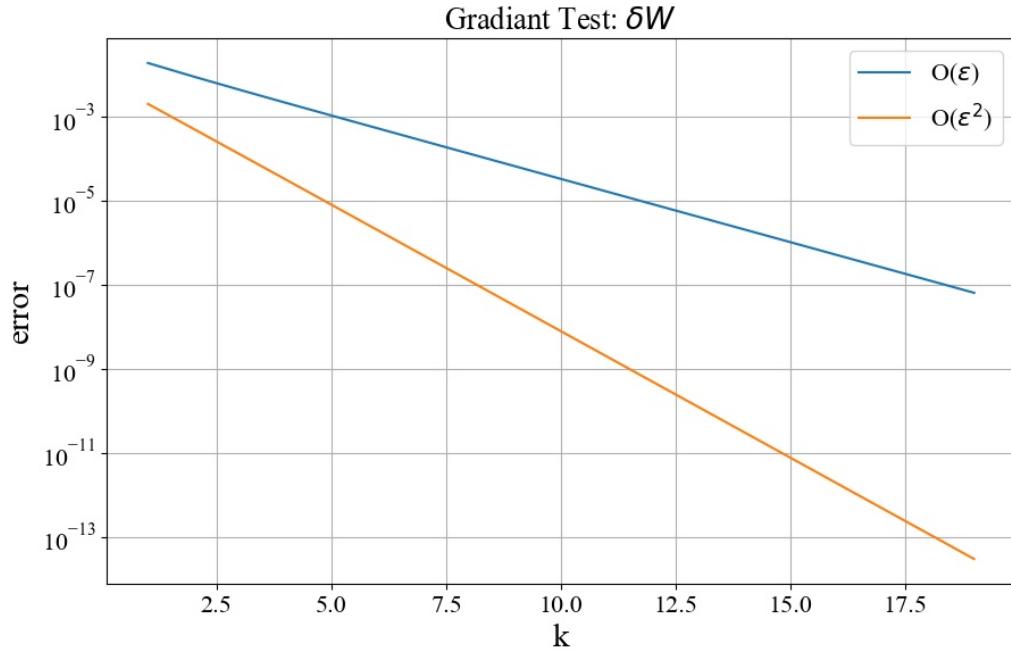


Figure 1.1: Gradient Test with respect to W . Iterating over different k 's which represent the power of ϵ . The initial value is $\epsilon = 0.5$.

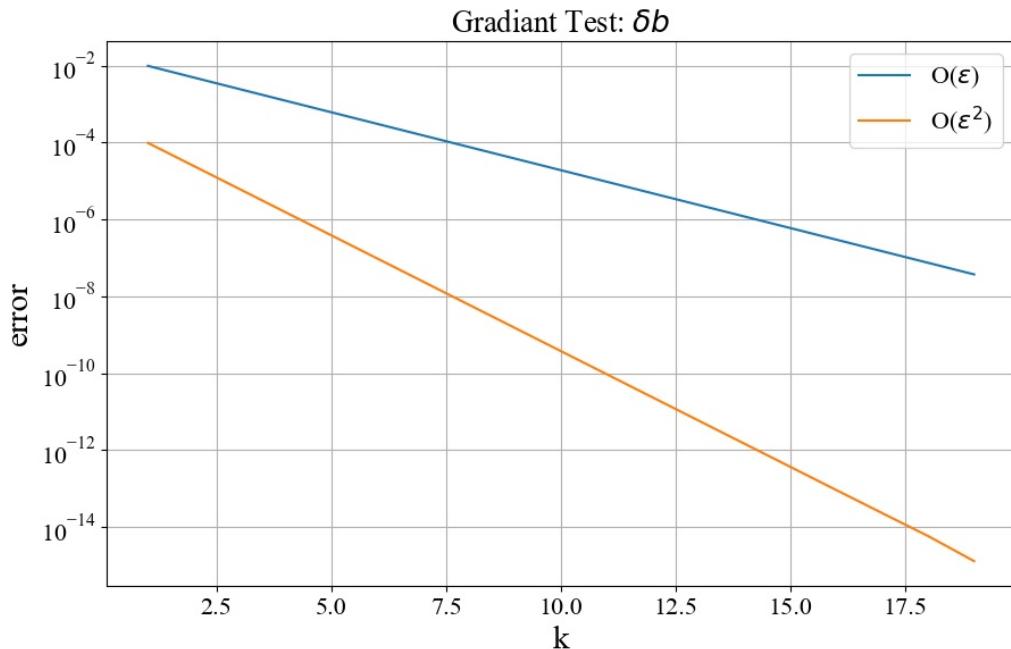


Figure 1.2: Gradient Test with respect to b (the biases). Iterating over different k 's which represent the power of ϵ . The initial value is $\epsilon = 0.5$.

The figures ensure that we are computing the gradients correctly.

1.2 SGD for Least Squares

In this section we show the Stochastic Gradient Descent (SGD) we implemented for the least squares loss function.

Note: This implementation is **specific** for the Least Squares function (and so is the whole example). We later generalize it for the network to come.

SGD for least squares function's code:

```
1 import numpy as np
2
3 xs = np.arange(-1, 1, 0.1)
4 ys = []
5
6
7 def f(x):
8     return 5 * (x ** 2)
9
10
11 data = [np.transpose([x ** 2, x, 1]) for x in xs]
12 expectations = [f(x) for x in xs]
13
14
15 def loss(f_, data_, weights):
16     data_list = list(data_)
17     return sum([f_(x_i, weights, y_i) for (x_i, y_i) in data_list]) / len(data_list)
18
19
20 def lls_func(x, w_, y):
21     return (1 / 2) * (np.transpose(x) @ w_ - y) ** 2
22
23
24 # The gradient.
25 def LLS_grad(x, w_, y):
26     return (np.transpose(x) @ w_ - y) * x
27
28
29 def SGD(mini_loss_func, f_grad, data, expectations, mb_size,
30         max_epochs, lr):
31     loss_hist = []
32     weights = [1, 1, 1] # initial weights.
33     for k in range(max_epochs):
34         # Devide data to mini-batches. TBD
35         num_of_mbs = int(len(data)) / mb_size
36         for j in range(num_of_mbs):
37             x_j = data[j] # Array of 3 - x^2, X, 1.
38             y_j = expectations[j] # Expectation value.
39             grad = f_grad(x_j, weights, y_j)
40             weights = weights - lr * grad
41             loss_hist += [loss(mini_loss_func, zip(data, expectations),
42                               weights)]
43     return weights, loss_hist
```

```
44 epochs = 1000
45 w, l = SGD(lls_func, LLS_grad, data, expectations, 1, epochs, 0.01)
```

The result of the optimization is shown bellow:

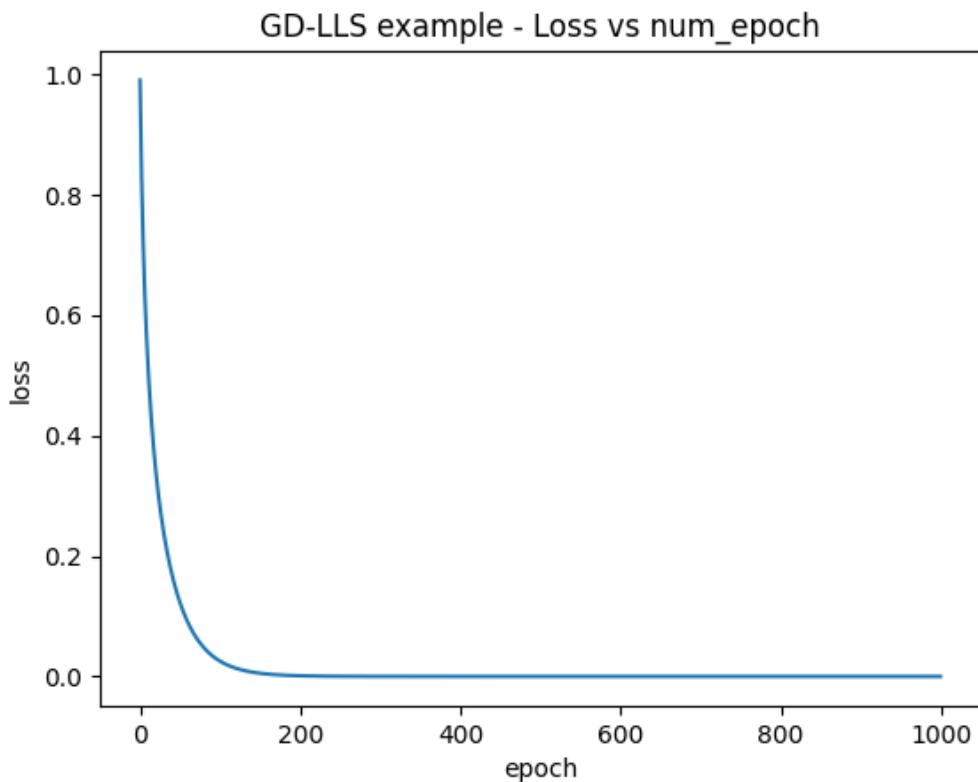


Figure 1.3: SGD for least squares. Iterating over 1000 epochs.

As we can see, the minimization of the Least-Squares function using the SGD algorithm indeed works.

1.3 SGD for softmax

We implemented the SGD for the softmax regression function as a method in the *util* module. We measured the loss and success percentage for different mini-batch sizes and different learning rates, on the Swiss-Roll data-set. We wrote the following code for the Softmax function:

Softmax regression code:

```
1 def soft_max_regression(X: np.array, W: np.array, C: np.array, b:  
2     np.array):  
3     """  
4         Computing the loss function 'Soft-Max regression'.  
5         :param X: The data input as a matrix of size nXm  
6         :param W: The weights, size of lXn (where l is the amount of  
7             labels)  
8         :param C: Indicators matrix. size of mXl.  
9         :return the loss function, and the gradients with respect to X,  
10            W.  
11        """  
12    expr = (W @ X + b).T # m X l  
13    arg = expr - etta(expr) # m X l  
14    prob = np.exp(arg) / np.sum(np.exp(arg), axis=1).reshape(-1, 1)  
15    m = len(X.T)  
16    F = - (1 / m) * np.sum(C * np.log(prob))  
17    grad_W = (1 / m) * (X @ (prob - C)).T  
18    grad_X = (1 / m) * (W.T @ (prob - C).T)  
19    grad_b = (1 / m) * np.sum((prob - C).T, axis=1).reshape(-1, 1)  
20    return F, grad_W, grad_X, grad_b
```

SGD for softmax code:

```
1 def SGD_for_Softmax(X, W, b, C, mb_size, max_epochs, lr):  
2     """  
3         :param loss_func: loss function to be evaluated.  
4         :param loss_func_grad: gradient of loss function.  
5         :param X: X matrix.  
6         :param W: W matrix.  
7         :param b: biases.  
8         :param C: Indicators matrix.  
9         :param mb_size: batches size.  
10        :param max_epochs: Number of epochs.  
11        :param lr: learning rate.  
12        :return: the value of W after the GD, and the loss for each  
13            epoch.  
14        """  
15    loss = []  
16    prob = []  
17    for k in range(max_epochs):  
18        bchs = util.generate_batches(X.T, C, mb_size)  
19        # Partition the data to random mini-batches of size mb_size  
20        .  
21        for curr_Mb, curr_Ind in bchs:  
22            # curr_Mb is a matrix of size n X mb_size.
```

```

21      # curr_Ind is a matrix of size mb_size X 1.
22      _, grad, _, _ = util.soft_max_regression(curr_Mb, W,
23      curr_Ind, b)
24      W -= lr * grad
25      l, p = util.calc_probs(X, W, C, b)
26      loss += [l]
27      prob += [util.get_accuracy(p, C)]
28  return W, loss, prob

```

** Auxiliary codes (such as 'generate_batches') are attached in section 4

The results are shown below:

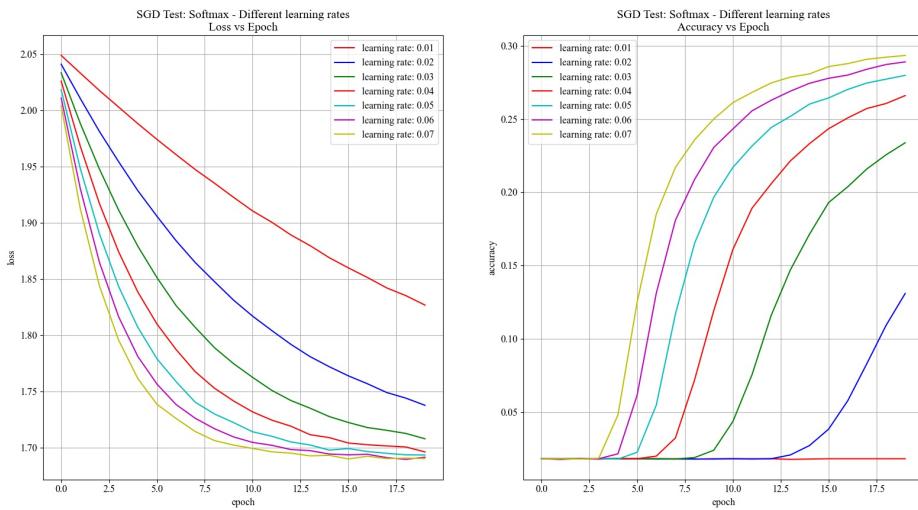


Figure 1.4: SGD for Softmax, on the training data. We iterated over 20 epochs for different learning rates, with mini-batch size of 30.

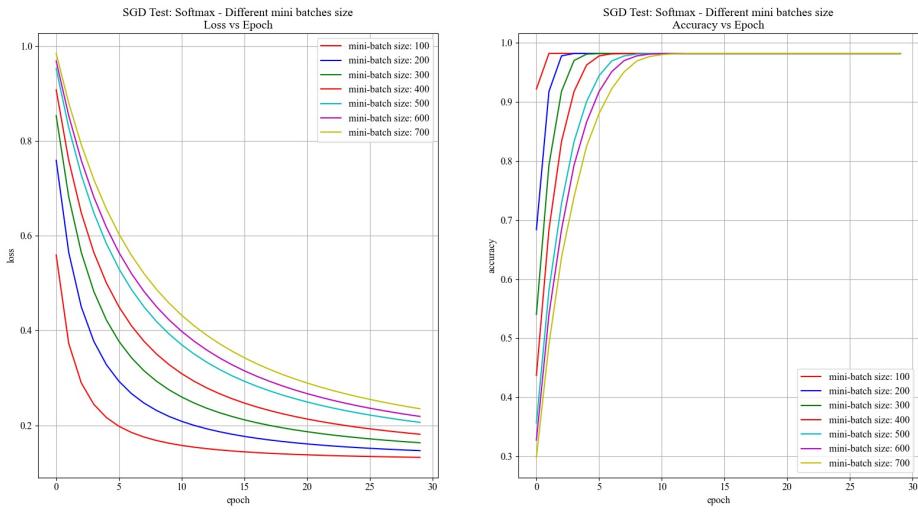


Figure 1.5: SGD for Softmax, on the training data. We iterated over 20 epochs for different mini-batches size, with a learning rate of 0.05.

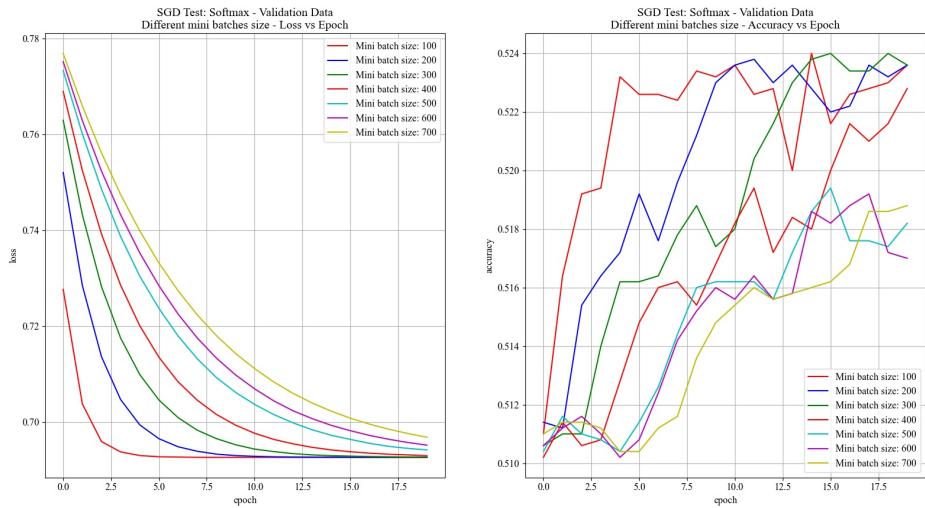


Figure 1.6: SGD for Softmax, on the validation data. We iterated over 20 epochs for different mini-batches size, with a learning rate of 0.05.

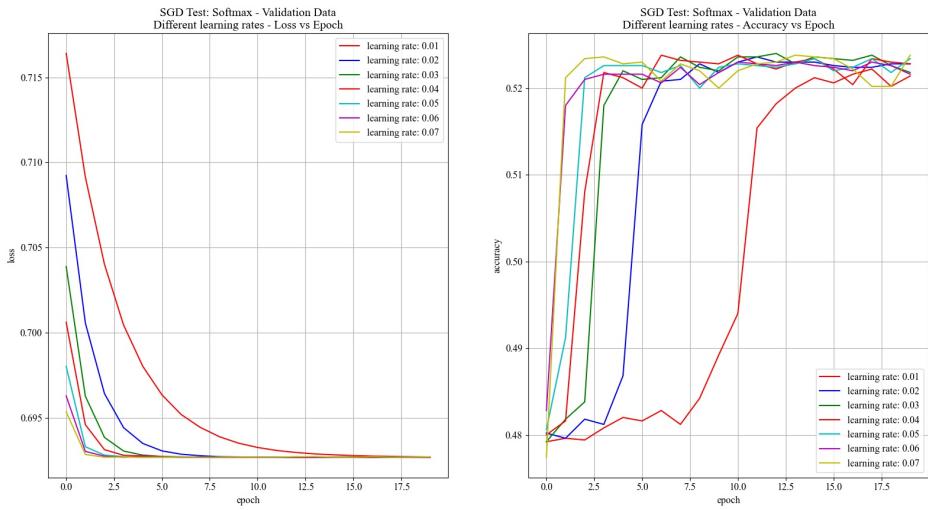


Figure 1.7: SGD for Softmax, on the validation data. We iterated over 20 epochs for different learning rates, with a mini-batches size of 30.

As we can see from the graphs above, the smaller mini-batches sizes and higher learning rates lead for higher accuracy. The latter is partly due to the range of learning rates we tested, as it is clear the for very high learning rates - SGD works badly (the zigzag effect).

2 Part II: The neural network

2.1 Jacobian test - regular layer

We now implement the whole neural-network (initialization, forward-pass, backward-pass, optimization loop).

We present here the Jacobian test we did for W and b for the regular (feed-forward) layer. As mentioned in the assignment we used the Tanh activation function, as it is differentiable.

The code for the test is shown bellow:

Jacobian Test code:

```
1 def JacTMV(X: np.array, W: np.array, b: np.array, V: np.array,
2             active_func: ActiveFunc, policy):
3     temp = W @ X + b
4     arg = active_func.deriv(temp) * V
5     dict_jac = {"W": arg @ X.T, "X": W.T @ arg, "b": np.sum(arg,
6                 axis=1).reshape(-1, 1)}
7     return dict_jac[policy]
8
9
10 def jacobian_test(X: np.array, W: np.array, b: np.array,
11                     active_func: ActiveFunc, policy):
12     """
13         Gradient test with respect for W.
14         :param X matrix.
15         :param W matrix.
16         :param b matrix.
17         :param active_func. The activation function (can be Relu or
18             Tanh).
19         :param policy. Indicates to which parameter we are doing the
20             test.
21         :return matplotlib graph which shows the jacobian test.
22     """
23     dict_param = {"W": W, "X": X, "b": b}
24     dict_name = {"W": "$\delta W$", "X": "$\delta X$",
25                  "b": "$\delta b$"}
26     d = np.random.randn(*dict_param[policy].shape)
27     V = np.random.randn(W.shape[0], X.shape[1])
28
29     d_vector = d.reshape(-1, 1)
30     err_1 = []
31     err_2 = []
32     ks = []
33     g_x = sum([arr[i] for i, arr in enumerate(active_func.activ(W @
34             X + b) @ V.T)])
35     dict_args = {"W": lambda: (W + epsilon * d) @ X + b, "X":
36                 lambda: W @ (X + epsilon * d) + b,
37                 "b": lambda: W @ X + (b + epsilon * d)}
38     for k in range(1, 20):
39         epsilon = 0.5 ** k
40         jac_v = JacTMV(X, W, b, V, active_func, policy).reshape(-1,
41             1)
```

```

33     g_x_d = sum([arr[i] for i, arr in enumerate(active_func.
34         activ(dict_args[policy]()) @ V.T)])
35     err_1.append(abs(g_x_d - g_x))
36     err_2.append(abs(g_x_d - g_x - (epsilon * d_vector.T @
37         jac_v)[0][0]))
38     ks.append(k)
39     print_test(ks, err_1, err_2, "Jacobian Test: " + dict_name[
40         policy])

```

The results we receive:

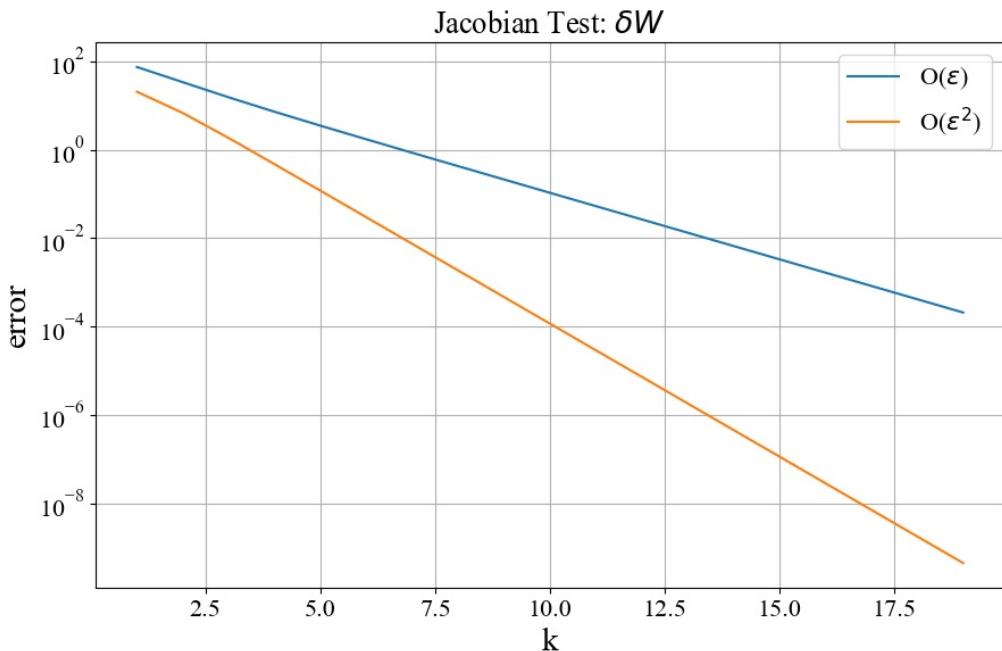


Figure 2.1: Jacobian test with respect to W . Iterating over different k 's which represent the power of ϵ . The initial value is $\epsilon = 0.5$.

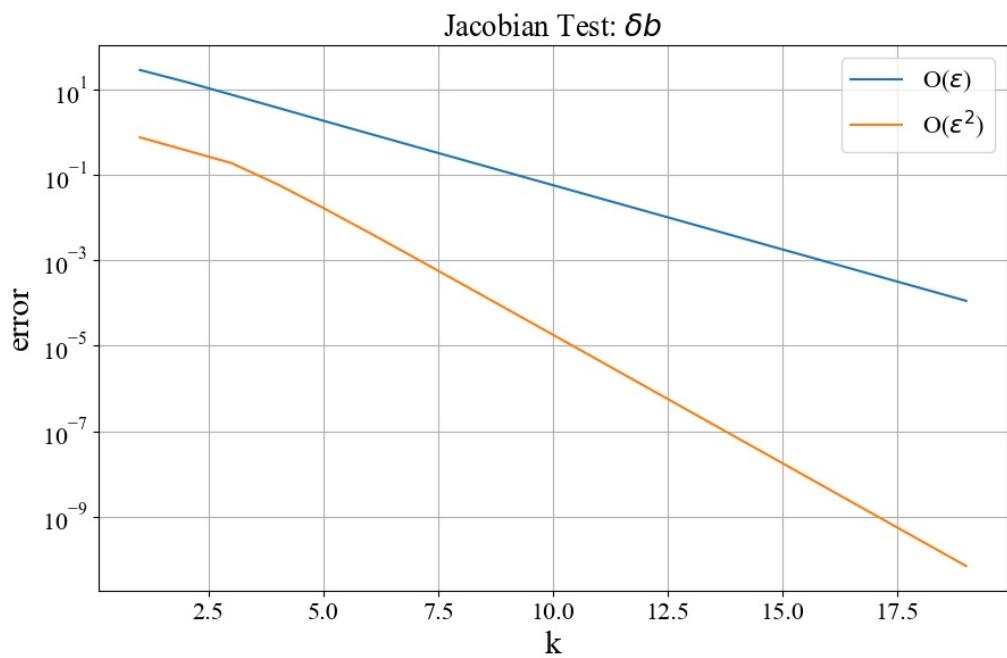


Figure 2.2: Jacobian test with respect to b . Iterating over different k 's which represent the power of ε . The initial value is $\varepsilon = 0.5$.

The results ensure that our derivatives are indeed exact.

2.2 ResNet Jacobian test

We present here the Jacobian test we did for the **ResNet** layer with respect to W and b .

The code for the test is shown bellow:

ResNet Jacobian Test code:

```

1 def resnetJacTMV(X: np.array, W1: np.array, W2: np.array, b: np.
2     array, V: np.array, active_func: ActiveFunc, policy):
3     arg = active_func.der(W1 @ X + b) * (W2.T @ V)
4     dict_jac = {"W1": arg @ X.T, "W2": V @ (active_func.act(W1 @ X
5     + b)).T, "X": V + W1.T @ arg,
6         "b": np.sum(arg, axis=1).reshape(-1, 1)}
7     return dict_jac[policy]
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36

```

```

    def resnet_jacobian_test(X: np.array, W1: np.array, W2: np.array, b
        : np.array, active_func: ActiveFunc, policy):
    """
    :param X: data matrix.
    :param W1: first weights matrix.
    :param W2: second weights matrix.
    :param b: bias matrix.
    :param active_func: activation function.
    :param policy: Indicates to which parameter we are doing the
    test.
    :return: matplotlib graph which shows the jacobian test.
    """
    dict_param = {"W1": W1, "W2": W2, "X": X, "b": b}
    dict_name = {"W1": "$\delta W_1$", "W2": "$\delta W_2$", "X": "$\delta X$",
        "b": "$\delta b$"}
    d = np.random.randn(*dict_param[policy].shape)
    V = np.random.randn(W2.shape[0], X.shape[1])
    d_vector = d.reshape(-1, 1)
    err_1 = []
    err_2 = []
    ks = []
    g_x = sum([arr[i] for i, arr in enumerate((X + W2 @ active_func.
        activ(W1 @ X + b)) @ V.T)])
    dict_args = {"W1": lambda: X + W2 @ active_func.activ((W1 +
        epsilon * d) @ X + b),
        "W2": lambda: X + (W2 + epsilon * d) @ active_func.
        activ(W1 @ X + b),
        "X": lambda: (X + epsilon * d) + W2 @ active_func.
        activ(W1 @ (X + epsilon * d) + b),
        "b": lambda: X + W2 @ active_func.activ(W1 @ X + (
            b + epsilon * d))}
    for k in range(1, 20):
        epsilon = 0.5 ** k
        jac_v = resnetJacTMV(X, W1, W2, b, V, active_func, policy).
        reshape(-1, 1)
        g_x_d = sum([arr[i] for i, arr in enumerate(dict_args[
            policy]() @ V.T)])

```

```

37     err_1.append(abs(g_x_d - g_x))
38     err_2.append(abs(g_x_d - g_x - (epsilon * d_vector.T @
39         jac_v)[0][0]))
40     ks.append(k)
41     print_test(ks, err_1, err_2, "Resnet Jacobian Test: " +
42     dict_name[policy])

```

The results are shown bellow:

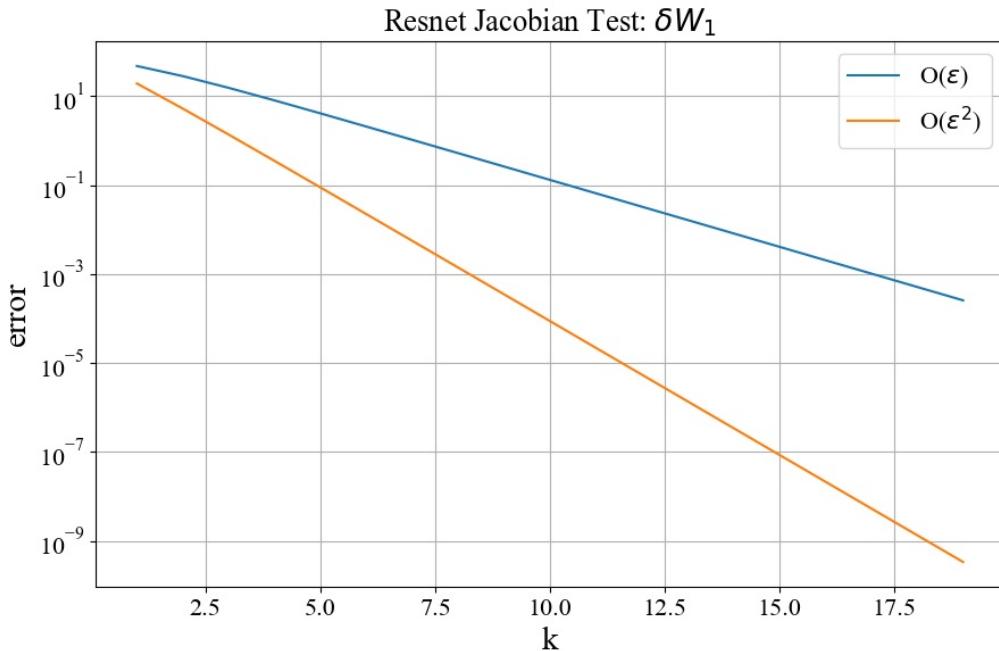


Figure 2.3: ResNet Jacobian test with respect to W_1 . Iterating over different ks which represent the power of ε . The initial value is $\varepsilon = 0.5$.

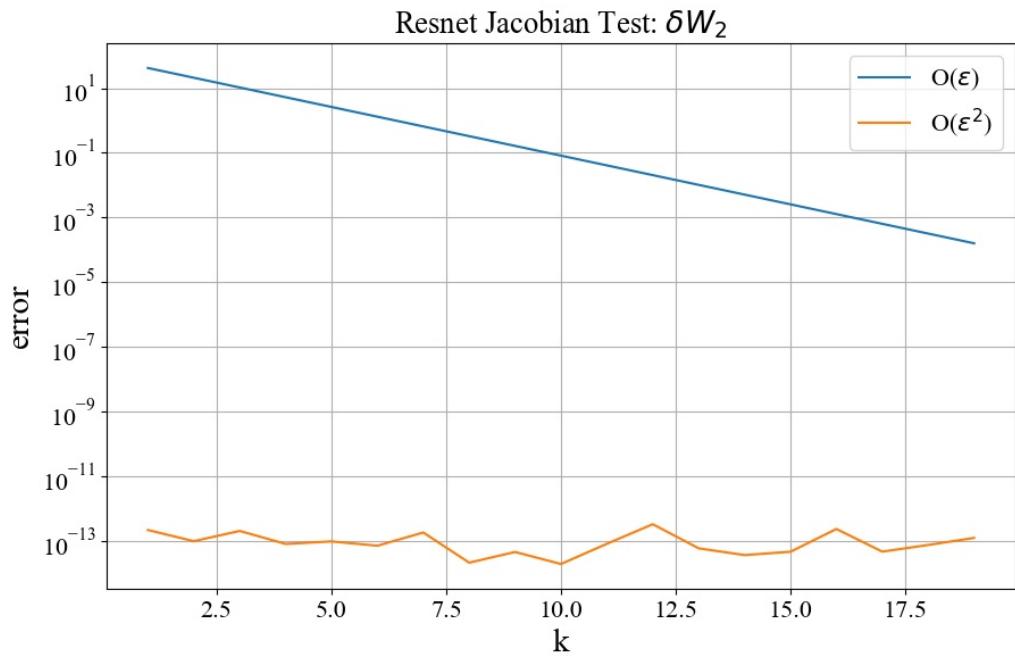


Figure 2.4: ResNet Jacobian test with respect to W_2 . Iterating over different ks which represent the power of ϵ . The initial value is $\epsilon = 0.5$.

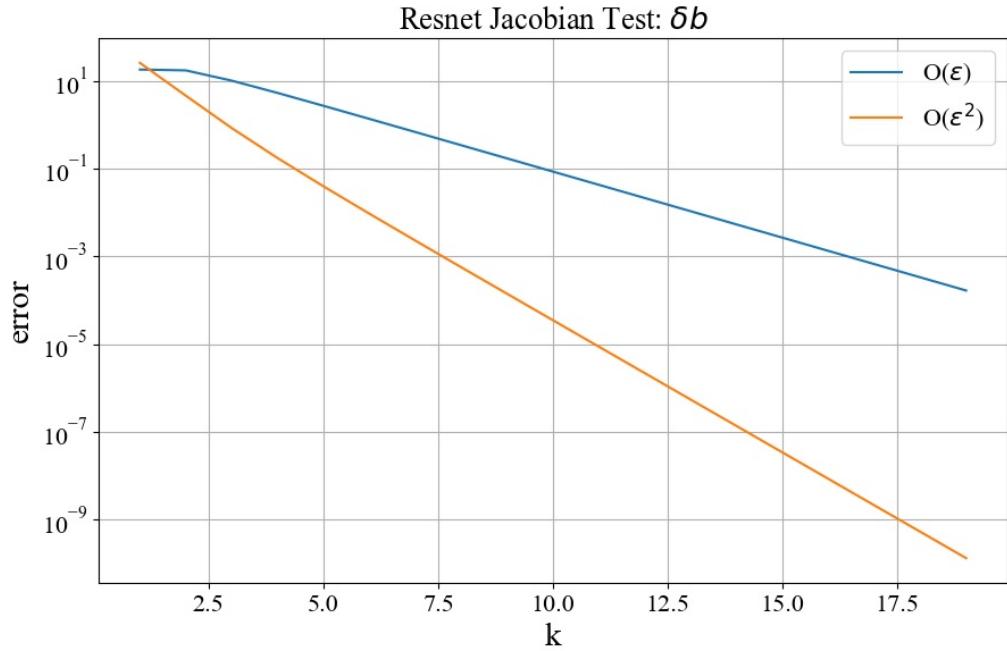


Figure 2.5: ResNet Jacobian test with respect to b . Iterating over different ks which represent the power of ϵ . The initial value is $\epsilon = 0.5$.

It seems that the Jacobian test for W_2 in Fig. 2.4 has an anomaly, but in fact it is the value we expected to see for W_2 because its second derivative is zero (and in our program, zero represented as the smallest value we can reach with 64 bit

Floating-Point representation). For the rest of the tests, as before, we conclude that our derivatives are correct.

2.3 Neural network Gradient test

We present here the Gradient test we did for the whole network with respect to W and b combined. The code for the test is shown bellow:

Network Gradient Test code:

```

1 def nn_gradient_test(nn, X, C):
2     """
3         :param nn: The Neural Network.
4         :param X: Data matrix.
5         :param C: Indicators matrix.
6         :return: matplotlib graph which shows the gradiant test.
7         """
8     Ws = []
9     dWs = []
10    bs = []
11    dbs = []
12
13    nn.forward(X)
14    nn.backward(C)
15    for layer in nn.layers:
16        Ws.append(layer.W.copy())
17        dWs.append(layer.grad_W.copy())
18        bs.append(layer.b.copy())
19        dbs.append(layer.grad_b.copy())
20    f_x, _ = nn.calc_loss_probs()
21
22    ds_W = [np.random.randn(*w.shape) for w in Ws]
23    ds_b = [np.random.randn(*b.shape) for b in bs]
24
25    err_1 = []
26    err_2 = []
27    ks = []
28
29    sum_grad = 0
30    for i, layer in enumerate(nn.layers):
31        sum_grad += (ds_W[i].reshape(1, -1) @ layer.grad_W.reshape
32        (-1, 1))[0][0]
33    for i, layer in enumerate(nn.layers):
34        sum_grad += (ds_b[i].reshape(1, -1) @ layer.grad_b.reshape
35        (-1, 1))[0][0]
36
37    for k in range(1, 10):
38        epsilon = 0.5 ** k
39        for i, layer in enumerate(nn.layers):
40            layer.set_W(Ws[i] + epsilon * ds_W[i])
41        for i, layer in enumerate(nn.layers):
42            layer.set_b(bs[i] + epsilon * ds_b[i])
43        nn.forward(X)
44        f_x_d, _ = nn.calc_loss_probs()
```

```

44     err_1.append(abs(f_x_d - f_x))
45     err_2.append(abs(f_x_d - f_x - epsilon * sum_grad))
46     ks.append(k)
47     print_test(ks, err_1, err_2, "Network Gradiant Test")

```

The results are shown bellow:

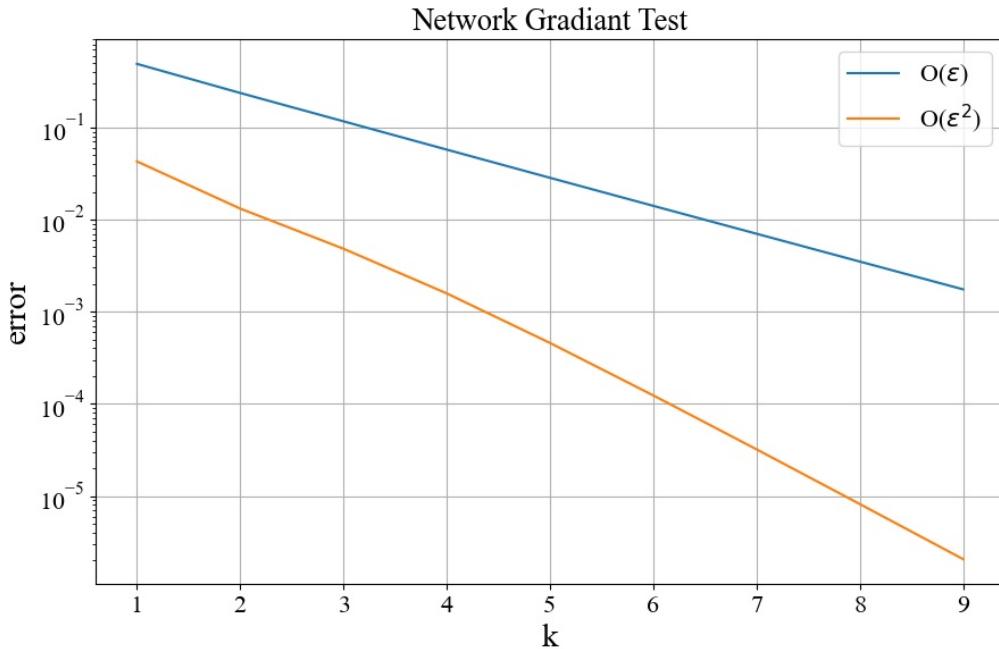


Figure 2.6: ResNet Jacobian test with respect to W_1 . Iterating over different ks which represent the power of ϵ . The initial value is $\epsilon = 0.5$.

The gradient test for the whole network ensures that we are computing the gradient correctly.

2.4 Learning of the whole neural network

2.4.1 Training

We show here the results we received for different hyper-parameters on the **training data** of the **Swiss roll** data-set.

We define the following architectures:

- Regular net: 8 layers with sizes: $[n, 3 \cdot n, 3 \cdot n, 5 \cdot n, 5 \cdot n, 3 \cdot n, 3 \cdot n, n]$.
Note: Regular net 2 is in the figures below has the exact same architecture as the above net. The difference is the mini-batch size we use.
- ResNet: A net of 15 layers, of the sizes: $[n, 2 \cdot n, 3 \cdot n, 4 \cdot n, 3 \cdot n, 2 \cdot n, n, n, n, 2 \cdot n, 3 \cdot n, 4 \cdot n, 3 \cdot n, 2 \cdot n, n]$ where the middle layer (of size n) is a Residual layer.
- Big Net: 8 layers with sizes: $[n, 12 \cdot n, 12 \cdot n, 144 \cdot n, 144 \cdot n, 12 \cdot n, 12 \cdot n, n]$.
- Full ResNet: 8 layers, all of the same size, all layers except the first and last are Residual layers.

Note: n refers to the input data dimension.

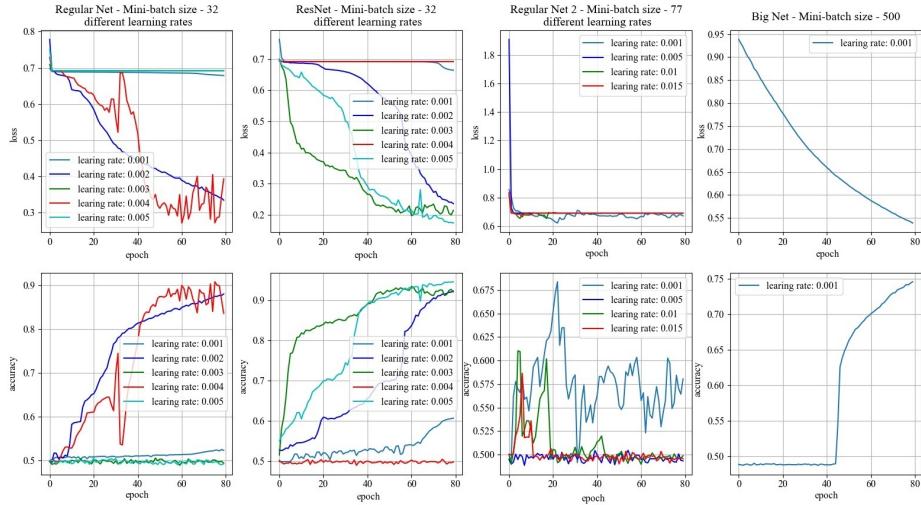


Figure 2.7: Loss and accuracy of the whole network, for the Swiss-Roll training data-set. Layer's size are fixed, varying the learning rate, mini-batch size and adding Resnet layers.

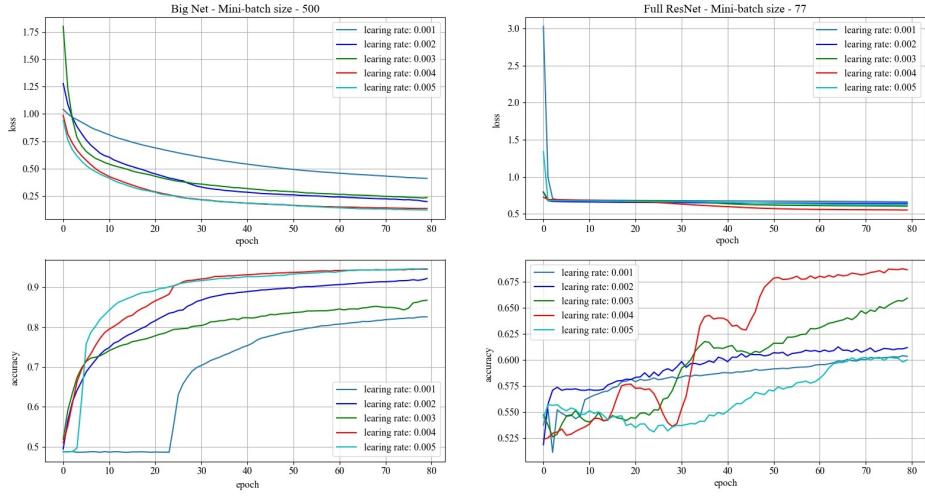


Figure 2.8: Training data of Swiss roll. Loss and accuracy for the whole network. The left side focus's on mini-batch size of 500 with large layers size (as described in the beginning of this section). In the right side we have 8 layers, six of which are ResNet layers (all but first and last).

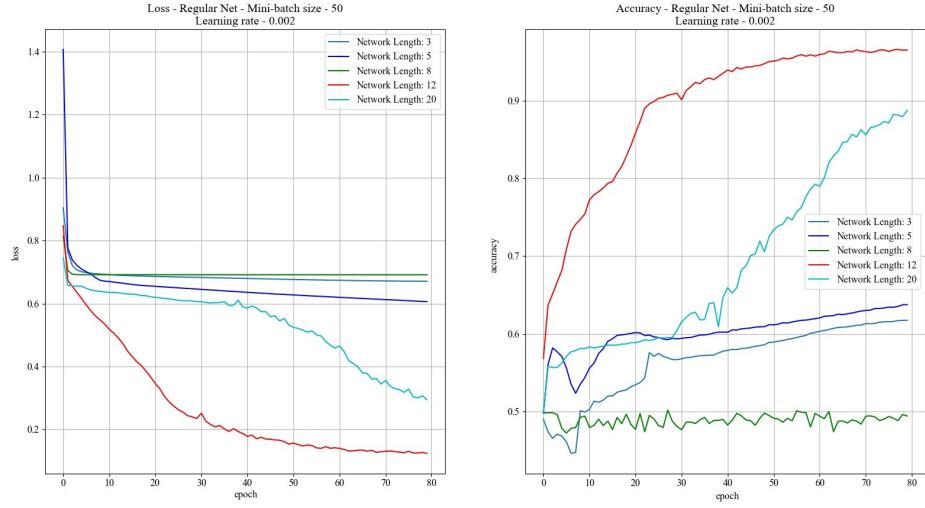


Figure 2.9: Loss and accuracy for the whole network. Different layers length. For all network lengths, we monotonically increase the layer size as we get to the middle of the net, and then we decrease it back as we get the end (as done for the Regular Net described above).

From the above research, we conclude that the best network hyper-parameters and length are:

- Network length: 12.
- Learning rate: 0.03.
- Mini-batch size: 50.

We shall use these parameter values for the learning to come on the data-sets provided.

2.4.2 Swiss Roll data

After finding the best hyper-parameters - we now **train with them** and check the resulting accuracy on the validation data.

As mentioned above, the hyper-parameters we use are network length of 12 (layers), a learning rate of 0.003 and mini-batch size of 50.

We emphasize that although the accuracy and loss are calculated for the validation data as well as for the training data, learning is performed only on the training data, i.e., the weights are not changed due to the values from the validation-data.

The results are shown bellow:

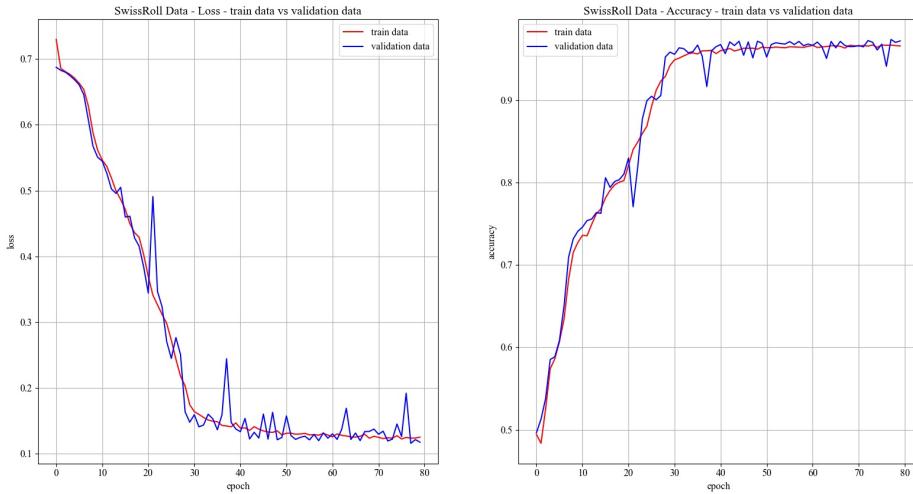


Figure 2.10: Swiss Roll validation data vs training data. After iterating over the whole data-set (in each epoch), we evaluated the loss and probability matrix of the validation data without learning from them.

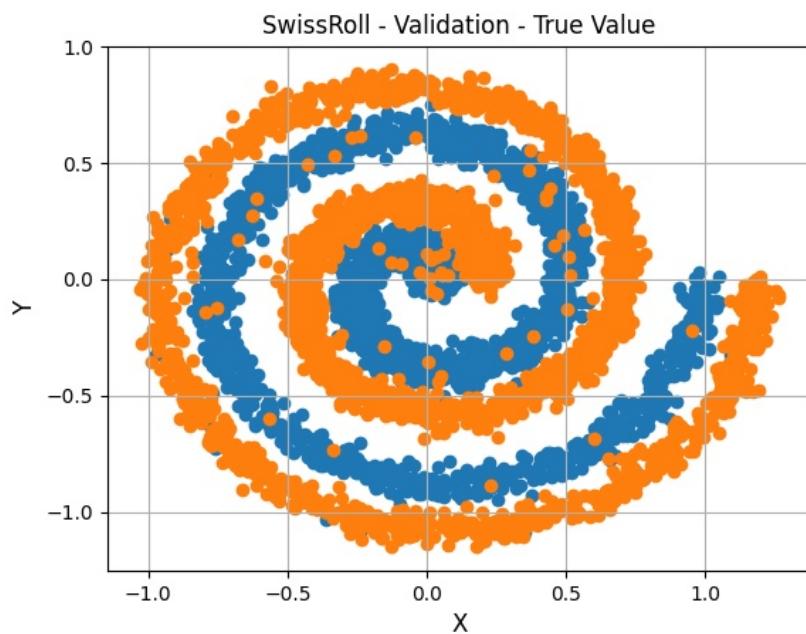


Figure 2.11: Swiss Roll raw validation data.

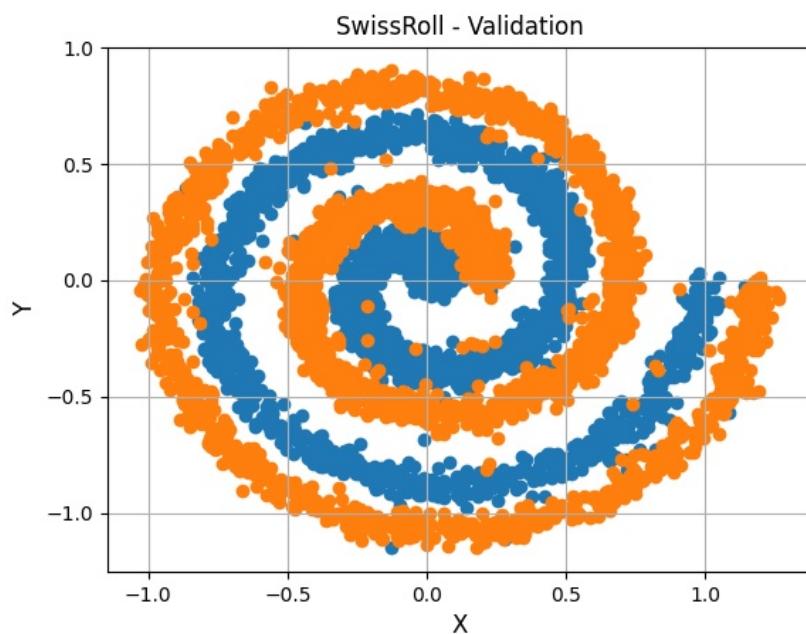


Figure 2.12: Our network's classification of the Swiss Roll validation data. This image represents the values of the points our network learned (classified). The network was trained on 12 layers with mini-batch size of 50 and learning rate of 0.002.

2.4.3 Peaks Data

We show here the result we achieved for the same network architecture in the Swiss Roll case (12 layers, learning rate 0.003 and 50 mini-batch size). Same as in Sec. 2.4.2, we trained our network only on the train data-set and evaluated the performance of the validation data-set compared to the train data-set (see Fig. 2.13).

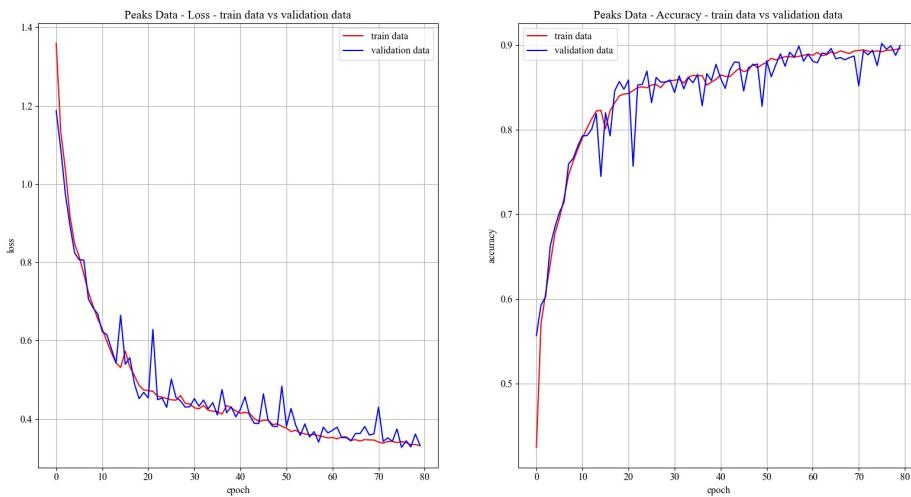


Figure 2.13: Peaks validation data vs training data. We evaluated in each epoch the loss and probability matrix of the validation data and plot it in the graph.

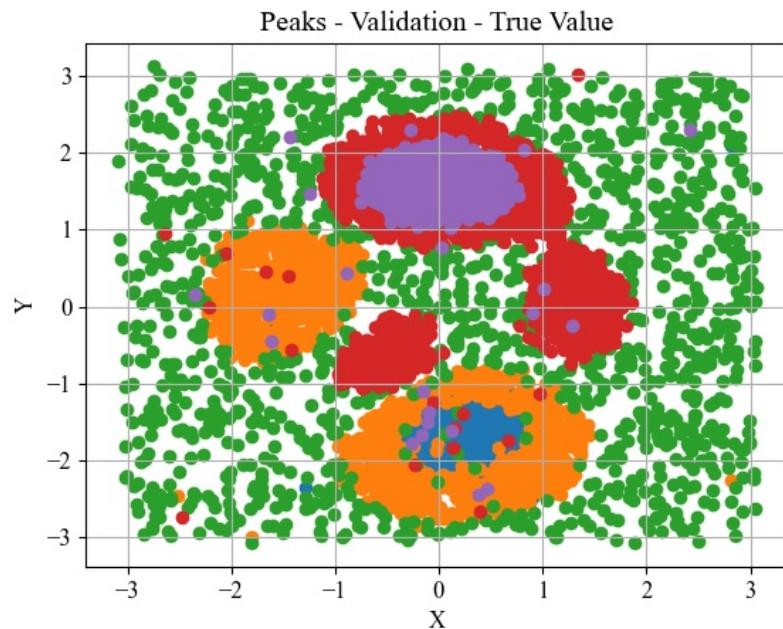


Figure 2.14: Peaks validation raw data.

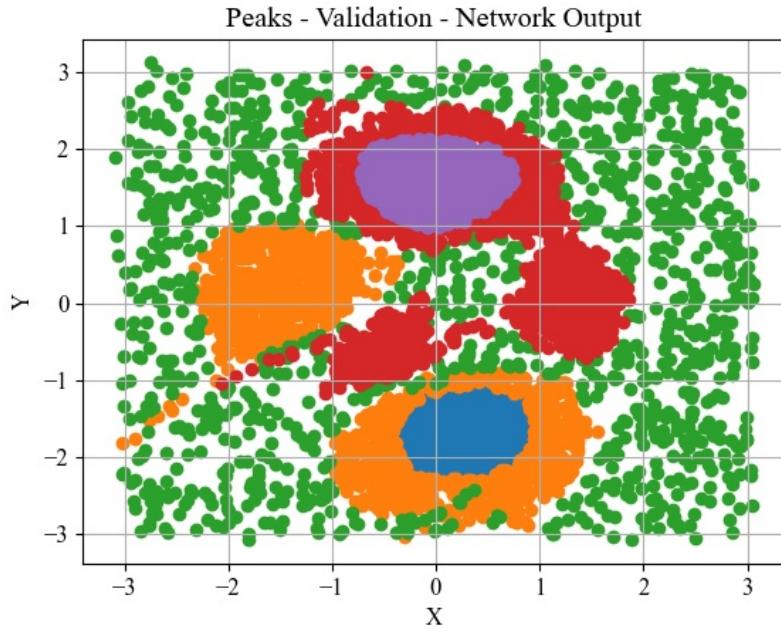


Figure 2.15: Peaks validation data. This graph represent the values of the points (labels) our network learned. The network was trained on 12 layers network with mini-batch size of 50 and learning rate of 0.003.

2.4.4 GMM Data

We show here the results we achieved for the same configuration as the Swiss Roll and the Peaks data (with 12 layers, 0.003 learning rate and 50 mini-batch size). Further explanation in Sec. 2.10.

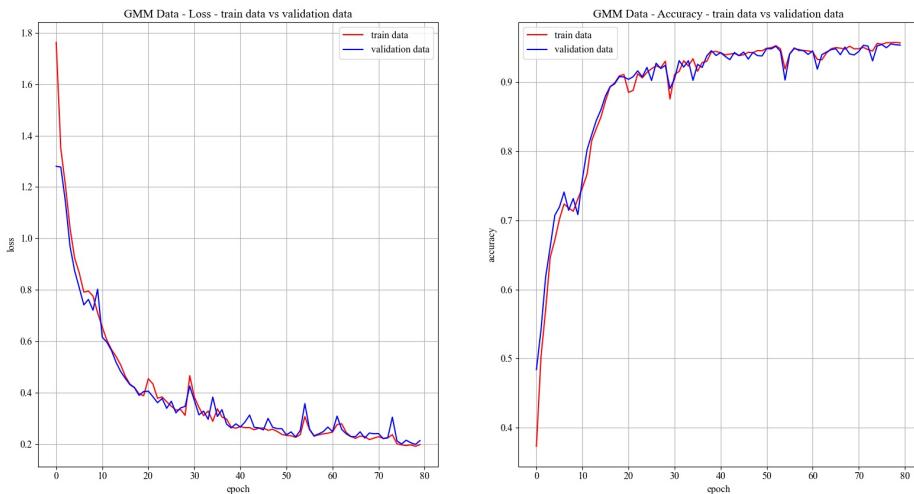


Figure 2.16: GMM validation data vs training data. We evaluated in each epoch the loss and probability matrix of the validation data.

2.4.5 Learning with 200 data points

We repeated the same procedure as in Sec. 2.4.2, 2.4.3, 2.4.4, with a training data set of 200 points. The results are shown bellow (Fig. 2.17, 2.18 and 2.19):

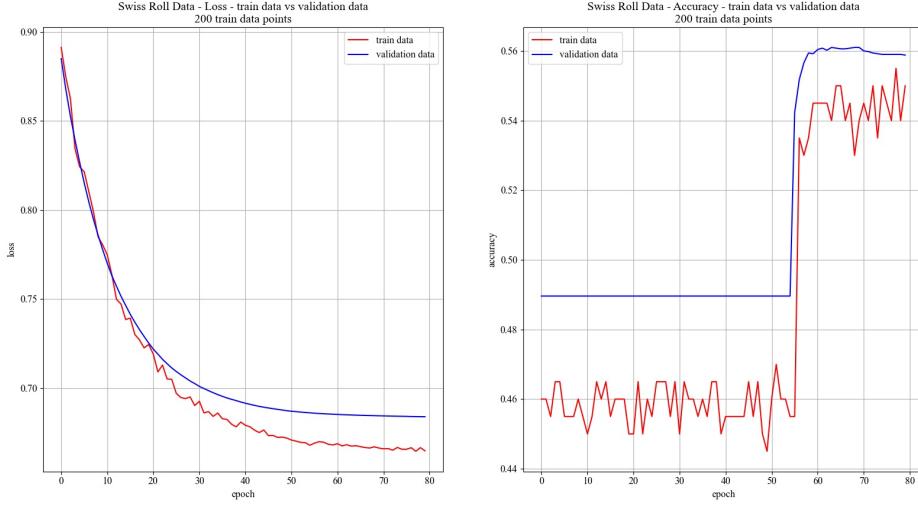


Figure 2.17: Swiss Roll validation data vs training data of 200 data points. We evaluated in each epoch the loss and probability matrix of the validation data.

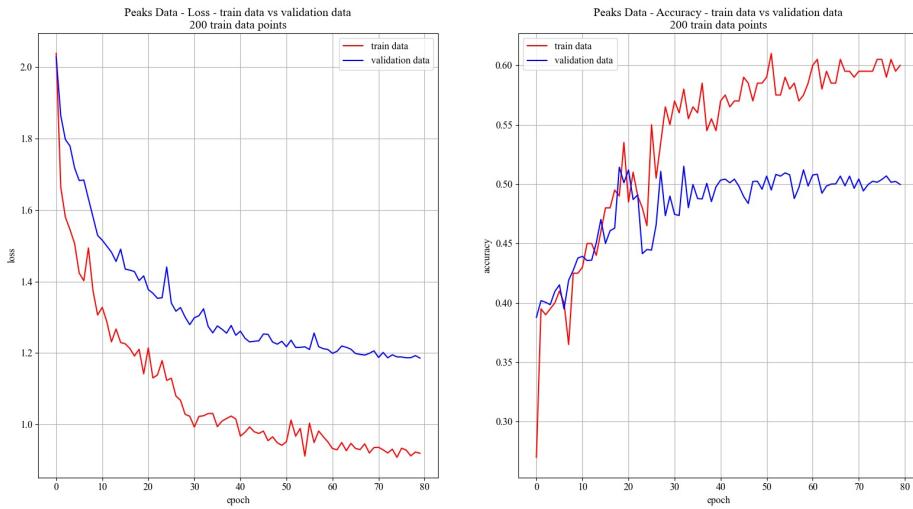


Figure 2.18: Peaks validation data vs training data of 200 data points. We evaluated in each epoch the loss and probability matrix of the validation data.

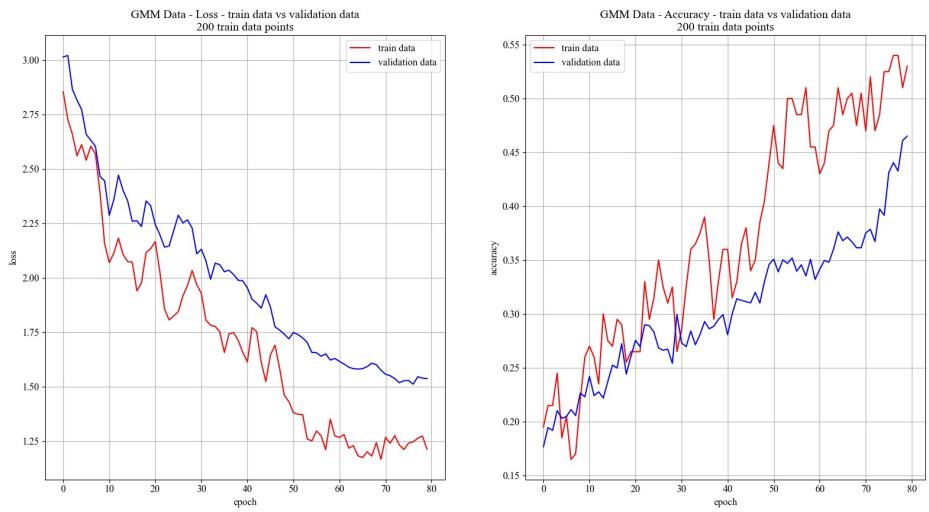


Figure 2.19: GMM validation data vs training data of 200 data points. We evaluated in each epoch the loss and probability matrix of the validation data.

3 Conclusions

1. We see that when we make sure of all our calculations using the gradient and Jacobian test (see Sec. 2.3), we indeed get a learning network, as wanted (emphasizes the importance of checking our derivatives).
2. As we can see in Fig. 2.7, as we add layers and residual layers to our network, we get smoother accuracy curves (then the regular network). This strengthens the claim that when we add more layers and neurons to our network, we can calculate more complex functions (such as the classifier for the Swiss-Roll data-set).
3. As we tested our network for different hyper-parameters for different networks, we saw that their effect on the loss and accuracy varies as we change the network architecture. Their inconsistent effect makes us understand that we need to analyze their effect in an architecture-specific manner (for example, fixing the randomness seed).
4. The results we receive are highly inconsistent from execution to execution, even with the same hyper-parameters (from very 'good' learning to very 'bad'). This emphasizes the effect of the randomness of our initial weights, and our randomly picked seed.
5. From Fig. 2.9 we can see that the best results achieved with 12 layers network. We chose this network rather than the big sized network which also preformed well (see Fig. 2.8) due to the shorter time complexity of this network (big matrices cost longer running time).
6. As we can see from Fig. 2.12 and 2.15, the network seems to 'fix' the raw data, i.e., its fluctuations vanish in the network output.
7. ResNet seems to have good loss reduction but it's accuracy is fluctuated and therefore we did not use it.
8. For all data samples, the comparison between the train network output and validation network output seems consistent (see Fig. 2.10, 2.13 and 2.16), with slightly small fluctuations of the validation graph.
9. We noticed that the correlation between the loss and the accuracy (see Fig. 2.7, 2.8) is not as strong as we initially thought. Smooth and exponentially decaying loss curve does not guarantee exponentially arising accuracy curve. Nevertheless, consistent behaviour of those two parameters increase the confidence and reliability of the architecture that we are choosing.
10. For all 3 data-sets we used (Swiss-Roll, Peaks, GMM) our network was able to achieve accuracy of at least 90%!
11. Upon learning on 200 data-points only (see Sec. 2.4.5), we deduce that the extensive data-set (i.e., a lot of data-points) is a crucial ingredient for sufficient learning.

Furthermore, we witness larger fluctuations in the values of the loss and accuracy, as well as smaller correlations between them, relative to the case of a large (extensive) data-set.

4 Auxiliary code

Generate mini-batch code:

```
1 def generate_batches(X, C, mb_size):
2     """
3     :param X: data matrix.
4     :param C: indicators matrix.
5     :param mb_size: batches size.
6     :return: array with all the mini batches and their
7             correspondent indicators.
8     """
9     data = []
10    mb = []
11    mbs = []
12    # Generate 'data' - An array containing 2-component arrays of [
13    # data, indicator].
13    for i in range(len(X)):
14        data.append([X[i], C[i]]) # C[i] is the i'th row,
15        corresponding to the i'th data-sample (it's indicator).
16        indices = list(range(len(data)))
17        random.shuffle(indices)
18        while len(indices) > mb_size:
19            for i in range(mb_size):
20                mb.append(data[indices.pop()])
21        """
22        Mb: Array of size nXmb_size.
23        Indicator: Matrix of size lXmb_size
24        """
25        Mb, Indicator = functools.reduce(lambda acc, curr: [acc[0]
26            + [curr[0]], acc[1] + [curr[1]]], mb, [[], []])
27        mb = []
28        mbs += [(np.array(Mb).T, np.array(Indicator))]
29
30    return mbs
```

Accuracy calculation code:

```
1 def get_accuracy(Prob: np.array, Indicator: np.array):
2     """
3     Get the probabilities matrix and indicator matrix and calculate
4     the accuracy
5     :param Prob:
6     :type Prob:
7     :param C:
8     :type C:
9     :return:
10    :rtype:
11    """
12    size = len(Prob)
13    probs = np.argmax(Prob, axis=1)
14    indicators = np.argmax(Indicator, axis=1)
15    counter = sum(probs == indicators)
16    return counter / size
```

Probability and Softmax code:

```
1 def soft_max_regression(X: np.array, W: np.array, C: np.array, b: np.array):
2     """
3         Computing the loss function 'Soft-Max regression'.
4         :param X. The data input as a matrix of size nXm
5         :param W. The weights, size of lXn (where l is the amount of
6             labels)
7         :param C. Indicators matrix. size of mXl.
8         :return the loss function, and the gradients with respect to X,
9             W.
10    """
11    expr = (W @ X + b).T # m X l
12    arg = expr - etta(expr) # m X l
13    prob = np.exp(arg) / np.sum(np.exp(arg), axis=1).reshape(-1, 1)
14    m = len(X.T)
15    F = - (1 / m) * np.sum(C * np.log(prob))
16    grad_W = (1 / m) * (X @ (prob - C)).T
17    grad_X = (1 / m) * (W.T @ (prob - C).T)
18    grad_b = (1 / m) * np.sum((prob - C).T, axis=1).reshape(-1, 1)
19    return F, grad_W, grad_X, grad_b
20
21
22 def calc_probs(X, W, C, b):
23     expr = (W @ X + b).T
24     arg = expr - etta(expr)
25     prob = np.exp(arg) / np.sum(np.exp(arg), axis=1).reshape(-1, 1)
26     m = len(X.T)
27     F = - (1 / m) * np.sum(C * np.log(prob))
28     return F, prob
29
30 def etta(A: np.array):
31     """
32         This method calculate the etta vector that required to reduce
33         from A in order to prevent numerical overflow.
34         :return etta vector. this vector is the column with the maximal
35         norm from A.
36    """
37    etta = np.array([])
38    for a in A:
39        etta = np.append(etta, max(a))
40    return etta.reshape(-1, 1)
```