

in3050_in4050_2022_assignment_2

March 9, 2022

0.1 IN3050/IN4050 Mandatory Assignment 2, 2022: Supervised Learning

0.1.1 Rules

Before you begin the exercise, review the rules at this website: <https://www.uio.no/english/studies/examinations/compulsory-activities/mn-ifi-mandatory.html>, in particular the paragraph on cooperation. This is an individual assignment. You are not allowed to deliver together or copy/share source-code/answers with others. By submitting this assignment, you confirm that you are familiar with the rules and the consequences of breaking them.

0.1.2 Delivery

Deadline: Friday, March 25, 2022, 23:59

Your submission should be delivered in Devilry. You may redeliver in Devilry before the deadline, but include all files in the last delivery, as only the last delivery will be read. You are recommended to upload preliminary versions hours (or days) before the final deadline.

0.1.3 What to deliver?

You are recommended to solve the exercise in a Jupyter notebook, but you might solve it in a Python program if you prefer.

If you choose Jupyter, you should deliver the notebook. You should answer all questions and explain what you are doing in Markdown. Still, the code should be properly commented. The notebook should contain results of your runs. In addition, you should make a pdf of your solution which shows the results of the runs. (If you can't export: notebook -> latex -> pdf on your own machine, you may do this on the IFI linux machines.)

If you prefer not to use notebooks, you should deliver the code, your run results, and a pdf-report where you answer all the questions and explain your work.

Your report/notebook should contain your name and username.

Deliver one single zipped folder (.zip, .tgz or .tar.gz) which contains your complete solution.

Important: if you weren't able to finish the assignment, use the PDF report/Markdown to elaborate on what you've tried and what problems you encountered. Students who have made an effort and attempted all parts of the assignment will get a second chance even if they fail initially. This exercise will be graded PASS/FAIL.

0.1.4 Goals of the assignment

The goal of this assignment is to get a better understanding of supervised learning with gradient descent. It will, in particular, consider the similarities and differences between linear classifiers and multi-layer feed forward networks (multi-layer perceptron, MLP) and the differences and similarities between binary and multi-class classification. A main part will be dedicated to implementing and understanding the backpropagation algorithm.

0.1.5 Tools

The aim of the exercises is to give you a look inside the learning algorithms. You may freely use code from the weekly exercises and the published solutions. You should not use ML libraries like scikit-learn or tensorflow.

You may use tools like NumPy and Pandas, which are not specific ML-tools.

The given precode uses NumPy. You are recommended to use NumPy since it results in more compact code, but feel free to use pure python if you prefer.

0.1.6 Beware

There might occur typos or ambiguities. This is a revised assignment compared to earlier years, and there might be new typos. If anything is unclear, do not hesitate to ask. Also, if you think some assumptions are missing, make your own and explain them!

0.1.7 Initialization

```
[1]: import numpy as np
import matplotlib.pyplot as plt
import sklearn #for datasets
```

1 Part 1: Linear classifiers

1.1 Datasets

We start by making a synthetic dataset of 2000 datapoints and five classes, with 400 individuals in each class. (See https://scikit-learn.org/stable/modules/generated/sklearn.datasets.make_blobs.html regarding how the data are generated.) We choose to use a synthetic dataset—and not a set of natural occurring data—because we are mostly interested in properties of the various learning algorithms, in particular the differences between linear classifiers and multi-layer neural networks together with the difference between binary and multi-class data.

When we are doing experiments in supervised learning, and the data are not already split into training and test sets, we should start by splitting the data. Sometimes there are natural ways to split the data, say training on data from one year and testing on data from a later year, but if that is not the case, we should shuffle the data randomly before splitting. (OK, that is not necessary with this particular synthetic data set, since it is already shuffled by default by scikit, but that will not be the case with real-world data.) We should split the data so that we keep the alignment between X and t , which may be achieved by shuffling the indices. We split into 50% for training,

25% for validation, and 25% for final testing. The set for final testing *must not be used* till the end of the assignment in part 3.

We fix the seed both for data set generation and for shuffling, so that we work on the same datasets when we rerun the experiments. This is done by the `random_state` argument and the `rng = np.random.RandomState(2022)`.

```
[2]: from sklearn.datasets import make_blobs
X, t = make_blobs(n_samples=[400,400,400, 400, 400],
                 centers=[[0,1],[4,1],[8,1],[2,0],[6,0]],
                 n_features=2, random_state=2019, cluster_std=1.0)
```

```
[3]: indices = np.arange(X.shape[0])
rng = np.random.RandomState(2022)
rng.shuffle(indices)
indices[:10]
```

```
[3]: array([1018, 1295,  643, 1842, 1669,   86,  164, 1653, 1174,  747])
```

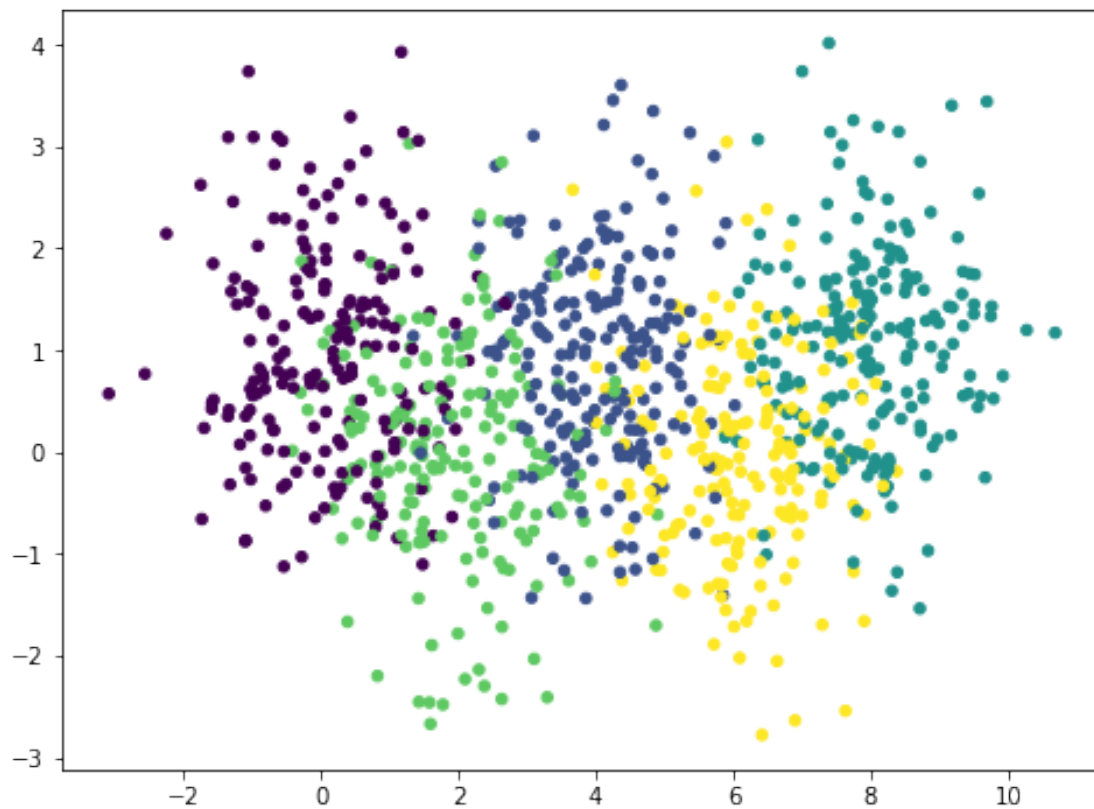
```
[4]: X_train = X[indices[:1000],:]
X_val = X[indices[1000:1500],:]
X_test = X[indices[1500:],:]
t_train = t[indices[:1000]]
t_val = t[indices[1000:1500]]
t_test = t[indices[1500:]]
```

Next, we will make a second dataset by merging the two smaller classes in (X,t) and call the new set (X, t2). This will be a binary set.

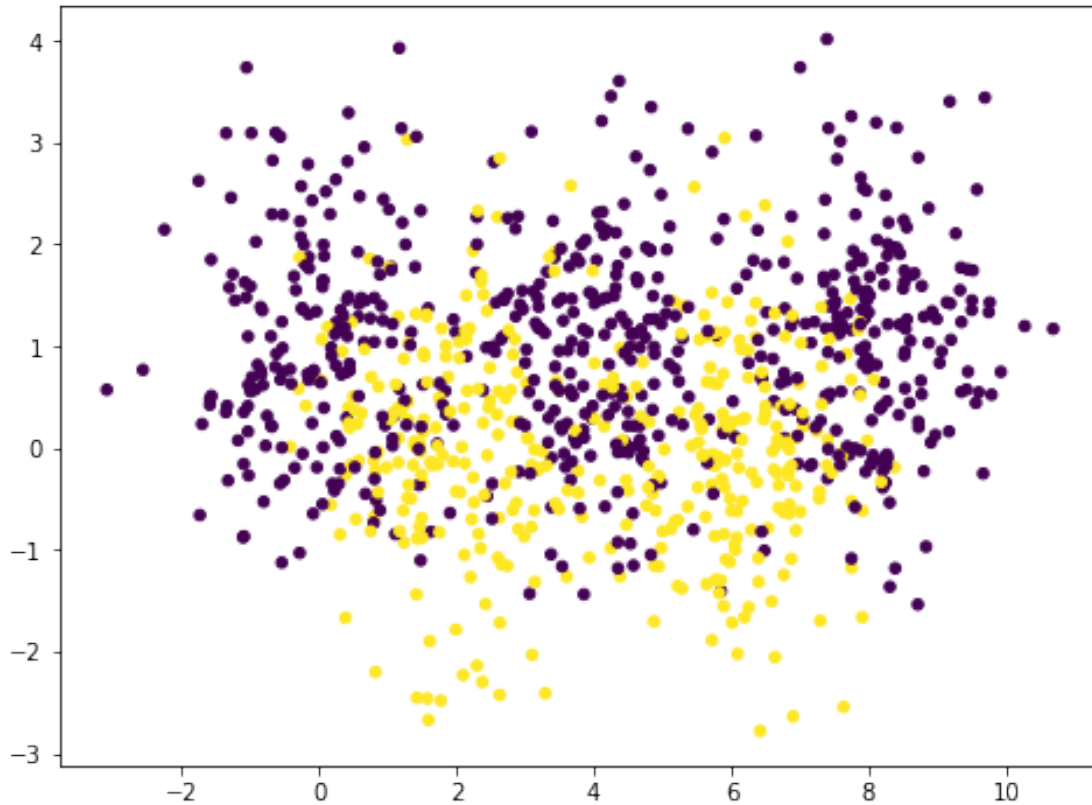
```
[5]: t2_train = t_train >= 3
t2_train = t2_train.astype('int')
t2_val = (t_val >= 3).astype('int')
t2_test = (t_test >= 3).astype('int')
```

We can plot the two training sets.

```
[6]: plt.figure(figsize=(8,6)) # You may adjust the size
plt.scatter(X_train[:, 0], X_train[:, 1], c=t2_train, s=20.0)
plt.show()
```



```
[7]: plt.figure(figsize=(8,6))  
plt.scatter(X_train[:, 0], X_train[:, 1], c=t2_train, s=20.0)  
plt.show()
```



1.2 Binary classifiers

1.2.1 Linear regression

We see that that set (X, t_2) is far from linearly separable, and we will explore how various classifiers are able to handle this. We start with linear regression. You may make your own implementation from scratch or start with the solution to the weekly exercise set 7, which we include here.

```
[8]: def add_bias(X):
    # Put bias in position 0
    sh = X.shape
    if len(sh) == 1:
        #X is a vector
        return np.concatenate([np.array([1]), X])
    else:
        # X is a matrix
        m = sh[0]
        bias = np.ones((m,1)) # Makes a m*1 matrix of 1-s
        return np.concatenate([bias, X], axis = 1)
```

```
[9]: class NumpyClassifier():
    """Common methods to all numpy classifiers --- if any"""

    def accuracy(self, X_test, y_test, **kwargs):
        pred = self.predict(X_test, **kwargs)
        if len(pred.shape) > 1:
            pred = pred[:,0]
        return np.sum(pred==y_test)/len(pred)
```

```
[10]: class NumpyLinRegClass(NumpyClassifier):

    def fit(self, X_train, t_train, eta = 0.1, epochs=10):
        """X_train is a Nxm matrix, N data points, m features
        t_train are the targets values for training data"""

        (k, m) = X_train.shape
        X_train = add_bias(X_train)

        self.weights = weights = np.zeros(m+1)

        for e in range(epochs):
            weights -= eta / k * X_train.T @ (X_train @ weights - t_train)

    def predict(self, x, threshold=0.5):
        z = add_bias(x)
        score = z @ self.weights
        return score > threshold
```

We can train and test a first classifier.

```
[11]: cl = NumpyLinRegClass()
cl.fit(X_train, t2_train)
cl.accuracy(X_val, t2_val)
```

```
[11]: 0.49
```

The result is far from impressive. Experiment with various settings for the hyper-parameters, eta and epochs. Report how the accuracy vary with the hyper-parameter settings. When you are satisfied with the result, you may plot the decision boundaries, as below.

Feel free to improve the colors and the rest av of the graphics. We have chosen a simple set-up which can be applied to more than two classes without substantial modifications.

```
[12]: def plot_decision_regions(X, t, clf=[], size=(8,6)):
    # Plot the decision boundary. For that, we will assign a color to each
    # point in the mesh [x_min, x_max]x[y_min, y_max].
    x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
    y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
```

```

h = 0.02 # step size in the mesh
xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])

plt.figure(figsize=size) # You may adjust this

# Put the result into a color plot
Z = Z.reshape(xx.shape)

plt.contourf(xx, yy, Z, alpha=0.2, cmap = 'Paired')

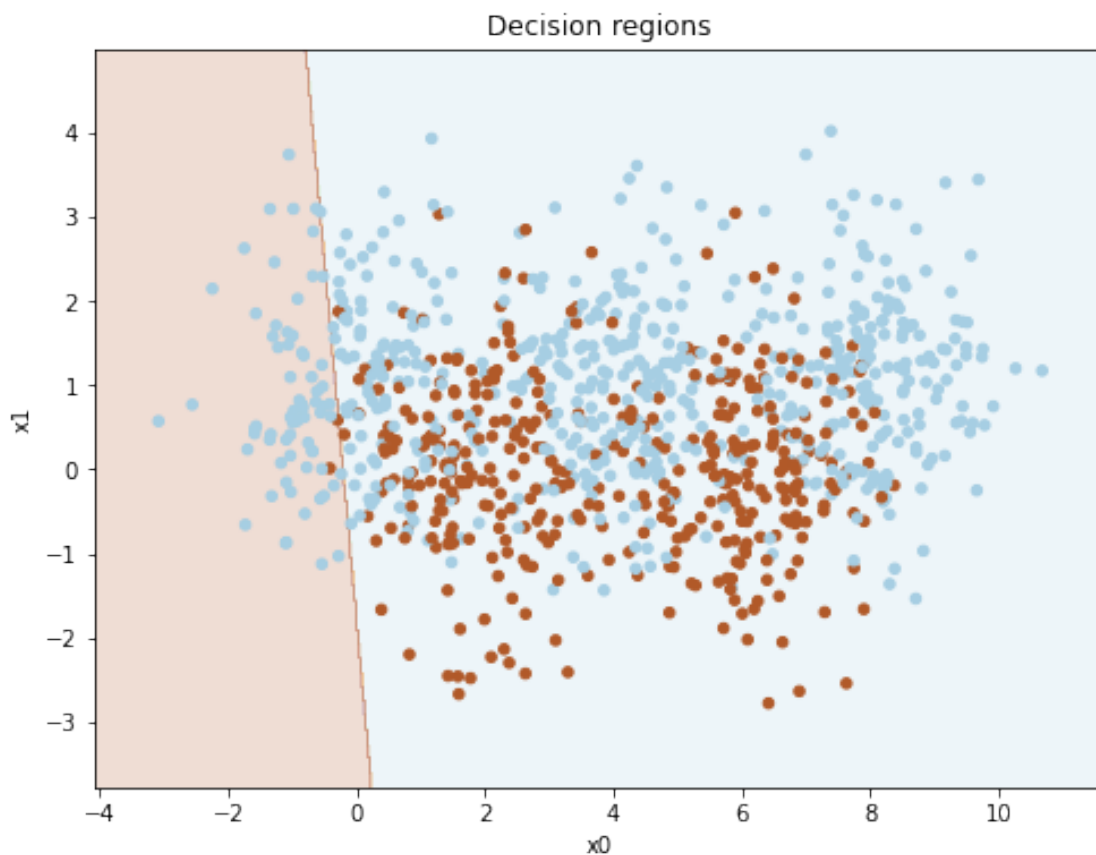
plt.scatter(X[:,0], X[:,1], c=t, s=20.0, cmap='Paired')

plt.xlim(xx.min(), xx.max())
plt.ylim(yy.min(), yy.max())
plt.title("Decision regions")
plt.xlabel("x0")
plt.ylabel("x1")

# plt.show()

```

```
[13]: plot_decision_regions(X_train, t2_train, cl)
```



1.2.2 Loss

The linear regression classifier is trained with mean squared error loss. So far, we have not calculated the loss explicitly in the code. Extend the code to calculate the loss on the training set for each epoch and to store the losses such that the losses can be inspected after training.

Train a classifier with your best settings from last point. After training, plot the loss as a function of the number of epochs.

1.2.3 Control training

The training runs for a number of epochs. We cannot know beforehand for how many epochs it is reasonable to run the training. One possibility is to run the training until the learning does not improve much. Extend the fit-method with a keyword argument, `loss_diff`, and stop training when the loss has not improved with more than `loss_diff`. Also add an attribute to the classifier which tells us after fitting how many epochs were ran.

In addition, extend the fit-method with optional arguments for a validation set (`X_val`, `t_val`). If a validation set is included in the call to fit, calculate the loss for the validation set, and the accuracy for both the training set and the validation set for each epoch.

Train classifiers with the best value for learning rate so far, and with varying values for `loss_diff`. For each run report, `loss_diff`, accuracy and number of epochs ran.

After a succesful training, plot both training loss snd vsvalidation loss as functions of the number of epochs in one figure, and both accuracies as functions of the number of epochs in another figure. Comment on what you see.

1.2.4 Logistic regression

You should now do similarly for a logistic regression classifier. Calculate loss and accuracy for training set and, when provided, also for validation set.

Remember that logistic regression is trained with cross-entropy loss. Hence the loss function is calculated differently than for linear regression.

After a succesful training, plot both losses as functions of the number of epochs in one figure, and both accuracies as functions of the number of epochs in another figure.

Comment on what you see. Do you see any differences between the linear regression classifier and the logistic regression classifier on this dataset?

Starting point: Code from weekly 7

```
[30]: def logistic(x):  
      return 1/(1+np.exp(-x))
```

```
[31]: class NumpyLogReg(NumpyClassifier):  
  
      def fit(self, X_train, t_train, eta = 0.1, epochs=10):
```



```

"""X_train is a Nxm matrix, N data points, m features
t_train are the targets values for training data"""

(k, m) = X_train.shape
X_train = add_bias(X_train)

self.weights = weights = np.zeros(m+1)

for e in range(epochs):
    weights -= eta / k * X_train.T @ (self.forward(X_train) - t_train)

def forward(self, X):
    return logistic(X @ self.weights)

def score(self, x):
    z = add_bias(x)
    score = self.forward(z)
    return score

def predict(self, x, threshold=0.5):
    z = add_bias(x)
    score = self.forward(z)
    return (score>threshold).astype('int')

```

1.3 Multi-class classifiers

We turn to the task of classifying when there are more than two classes, and the task is to ascribe one class to each input. We will now use the set (X, t) .

1.3.1 “One-vs-rest” with logistic regression

We saw in the lecture how a logistic regression classifier can be turned into a multi-class classifier using the one-vs-rest approach. We train one logistic regression classifier for each class. To predict the class of an item, we run all the binary classifiers and collect the probability score from each of them. We assign the class which ascribes the highest probability.

Build such a classifier. Train the resulting classifier on $(X_{\text{train}}, t_{\text{train}})$, test it on $(X_{\text{val}}, t_{\text{val}})$, tune the hyper-parameters and report the accuracy.

Also plot the decision boundaries for your best classifier similarly to the plots for the binary case.

1.3.2 For in4050-students: Multi-nominal logistic regression

The following part is only mandatory for in4050-students. In3050-students are also welcome to make it a try. Everybody has to return for the part 2 on multi-layer neural networks.

In the lecture, we contrasted the one-vs-rest approach with the multinomial logistic regression, also called softmax classifier. Implement also this classifier, tune the parameters, and compare the results to the one-vs-rest classifier. (Don’t expect a large difference on a simple task like this.)

Remember that this classifier uses exponentiation followed by softmax in the forward phase. For loss, it uses cross-entropy loss. The loss has a somewhat simpler form than in the binary case. To calculate the gradient is a little more complicated. The actual gradient and update rule is simple, however, as long as you have calculated the forward values correctly.

2 Part II

2.1 Multi-layer neural networks

We will implement the Multi-layer feed forward network (MLP, Marsland sec. 4.2.1), where we use mean squared loss together with logistic activation in both the hidden and the last layer.

Since this part is more complex, we will do it in two rounds. In the first round, we will go stepwise through the algorithm with the dataset (X, t) . We will initialize the network and run a first round of training, i.e. one pass through the algorithm at p. 78 in Marsland.

In the second round, we will turn this code into a more general classifier. We can train and test this on (X, t) and (X, t_2) , but also on other datasets.

2.2 Round 1: One epoch of training

2.2.1 Scaling

First we have to scale our data. Make a standard scaler (normalizer) and scale the data. Remember, not to follow Marsland on this point. The scaler should be constructed from the training data only, but be applied both to training data and later on to validation and test data.

```
[52]: # Your code
```

2.2.2 Initialization

We will only use one hidden layer. The number of nodes in the hidden layer will be a hyper-parameter provided by the user; let's call it *dim_hidden*. (*dim_hidden* is called *M* by Marsland.) Initially, we will set it to 3. This is a hyper-parameter where other values may give better results, and the hyper-parameter could be tuned.

Another hyper-parameter set by the user, is the learning rate. We set the initial value to 0.01, but also this may need tuning.

```
[55]: eta = 0.01 #Learning rate
      dim_hidden = 3
```

We assume that the input X_{train} (after scaling) is a matrix of dimension $P \times dim_in$, where P is the number of training instances, and dim_in is the number of features in the training instances (L in Marsland). Hence we can read dim_in off from X_{train} .

The target values have to be converted from simple numbers, 0, 2,.. to “one-hot-encoded” vectors similarly to the multi-class task. After the conversion, we can read dim_out off from t_{train} .

```
[56]: # convert t_train
      dim_in = 0 # Calculate the correct value from the input data
```

```
dim_out = 0 # Calculate the correct value from the input data
```

We need two sets of weights: `weights1` between the input and the hidden layer, and `weights2`, between the hidden layer and the output. Make sure that you take the bias terms into consideration and get the correct dimensions. The weight matrices should be initialized to small random numbers, not to zeros. It is important that they are initialized randomly, both to ensure that different neurons start with different initial values and to generate different results when you rerun the classifier. In this introductory part, we have chosen to fix the random state to make it easier for you to control your calculations. But this should not be part of your final classifier.

```
[57]: # Your code
```

```
[59]: rng = np.random.RandomState(2022)
weights1 = (rng.rand(dim_in + 1, dim_hidden) * 2 - 1)/np.sqrt(dim_in)
weights2 = (rng.rand(dim_hidden+1, dim_out) * 2 - 1)/np.sqrt(dim_hidden)
```

```
[60]: weights1
```

```
[60]: array([[ -0.6938717 , -0.00133246, -0.54675803],
          [-0.63643285,  0.26220593, -0.01840165],
          [ 0.56237224,  0.20852872,  0.56139063]])
```

2.2.3 Forwards phase

We will run the first step in the training, and start with the forward phase. Calculate the activations after the hidden layer and after the output layer. We will follow Marsland and use the logistic (sigmoid) activation function in both layers. Inspect whether the results seem reasonable with respect to format and values.

```
[61]: # Your code
# hidden_activations =
```

```
[62]: # Your code
# output_activations =
```

To control that you are on the right track, you may compare your first output value with our result. We have put the bias term -1 in position 0 in both layers. If you have done anything differently from us, you will not get the same numbers. But you may still be on the right track!

```
[68]: outputs[0, :]
```

```
[68]: array([0.28969058, 0.44120276, 0.41012141, 0.38135763, 0.44130415])
```

2.2.4 Backwards phase

Calculate the delta terms at the output. We assume, like Marsland, that we use sum of squared errors. (This amounts to the same as using the mean square error).

```
[69]: # Your code
```

Calculate the delta terms in the hidden layer.

```
[73]: # Your code
```

Update the weights in both layers.. See whether the weights have changed.

As an aid, you may compare your new weights with our results. But again, you may have done everything correctly even though you get a different result. For example, there are several ways to introduce the mean squared error. They may give different results after one epoch. But if you run sufficiently many epochs, you will get about the same classifier.

```
[78]: print("New weights:")  
      print(weights1)
```

New weights:

```
[[-0.64918987  0.0049323 -0.57494453]  
 [-0.63388739  0.32480283  0.08132063]  
 [ 0.51939233  0.17555818  0.58007288]]
```

2.3 Step 2: A Multi-layer neural network classifier

2.3.1 Make the classifier

You want to train and test a classifier on (X, t) . You could have put some parts of the code in the last step into a loop and run it through some iterations. But instead of copying code for every network we want to train, we will build a general Multi-layer neural network classifier as a class. This class will have some of the same structure as the classifiers we made for linear and logistic regression. The task consists mainly in copying in parts from what you did in step 1 into the template below. Remember to add the *self*- prefix where needed, and be careful in your use of variable names. And don't fix the random numbers within the classifier.

```
[79]: class MNClassifier():  
      """A multi-layer neural network with one hidden layer"""  
  
      def __init__(self, eta = 0.001, dim_hidden = 6):  
          """Initialize the hyperparameters"""  
          self.eta = eta  
          self.dim_hidden = dim_hidden  
  
          # Should you put additional code here?  
  
      def fit(self, X_train, t_train, epochs = 100):  
          """Initialize the weights. Train *epochs* many epochs."""  
  
          # Initilaization  
          # Fill in code for initalization
```

```

for e in range(epochs):
    # Run one epoch of forward-backward
    #Fill in the code
    pass

def forward(self, X):
    """Perform one forward step.
    Return a pair consisting of the outputs of the hidden_layer
    and the outputs on the final layer"""
    #Fill in the code

def accuracy(self, X_test, t_test):
    """Calculate the accuracy of the classifier for the pair (X_test,
    ↪t_test)
    Return the accuracy"""
    #Fill in the code

```

2.3.2 Multi-class

Train the network on (X_train, t_train) (after scaling), and test on (X_val, t_val). Tune the hyperparameters to get the best result: - number of epochs - learning rate - number of hidden nodes.

When you are content with the hyperparameters, you should run the same experiment 10 times, collect the accuracies and report the mean value and standard deviation of the accuracies across the experiments. This is common practise when you apply neural networks as the result may vary slightly between the runs. You may plot the decision boundaries for one of the runs.

Discuss shortly how the results and decision boundaries compare to the “one-vs-rest” classifier.

2.3.3 Binary class

Let us see whether a multilayer neural network can learn a non-linear classifier. Train a classifier on (X_train, t2_train) and test it on (X_val, t2_val). Tune the hyper-parameters for the best result. Run ten times with the best setting and report mean and standard deviation. Plot the decision boundaries.

2.4 For in4050-students: Early stopping

The following part is only mandatory for in4050-students. In3050-students are also welcome to make it a try. Everybody has to return for the part 2 on multi-layer neural networks.

There is a danger of overfitting if we run too many epochs of training. One way to control that is to use early stopping. We can use (X_val, t_val) as valuation set when training on (X_train, t_train).

Let $e=50$ or $e=10$ (You may try both or choose some other number) After e number of epochs, calculate the loss for both the training set ($X_{\text{train}}, t_{\text{train}}$) and the validation set ($X_{\text{val}}, t_{\text{val}}$), and store them.

Train a classifier for many epochs. Plot the losses for both the training set and the validation set in the same figure and see whether you get the same effect as in figure 4.11 in Marsland.

Modify the code so that the training stops if the loss on the validation set is not reduced by more than t after e many epochs, where t is a threshold you provide as a parameter.

Run the classifier with various values for t and report the accuracy and the number of epochs ran.

3 Part III: Final testing

We can now perform a final testing on the held-out test set.

3.1 Binary task (X, t_2)

Consider the linear regression classifier, the logistic regression classifier and the multi-layer network with the best settings you found. Train each of them on the training set and evaluate on the held-out test set, but also on the validation set and the training set. Report in a 3 by 3 table.

Comment on what you see. How do the three different algorithms compare? Also, compare the result between the different data sets. In cases like these, one might expect slightly inferior results on the held-out test data compared to the validation data. Is so the case?

Also report precision and recall for class 1.

3.2 Multi-class task (X, t)

For IN3050 students compare the one-vs-rest classifier to the multi-layer perceptron. Evaluate on test, validation and training set as above. In4050-students should also include results from the multi-nomial logistic regression.

Comment on the results.