

FYS-3012 Home Exam

CANDIDATE 78

UiT - Norges Arktiske Universitet

October 30, 2023

Task 1

1a)

To find J_{min} we look at the general l dimensional case where all the classes are identical. We have:

$$J_3 = \text{trace}\{S_w^{-1} S_m\}$$

We also have that:

$$S_m = S_w + S_b$$

We can first look at S_b , which is defined as:

$$S_b = \sum_{i=1}^M P_i (\mu_i - \mu_0)(\mu_i - \mu_0)^T$$

Since all classes are identical we know that all the mean vectors, as well as the global mean vector are equal. Meaning S_b is a $l \times l$ zero matrix, and thus $S_m = S_w$

Which gives:

$$J_{min} = \text{trace}\{S_w^{-1} S_w\} = \text{trace}\{I\}$$

Where I is the $l \times l$ identity matrix. Since the trace operation is defined as the sum of diagonal elements we get that $\text{trace}\{I\}$ is equal to l

$$J_{min} = l$$

But we can still simplify our modified J_3 score

$$\begin{aligned} J_3 &= \text{trace}\{S_w^{-1} S_m\} - J_{min} \\ &= \text{trace}\{S_w^{-1} (S_w + S_b)\} - J_{min} \\ &= \text{trace}\{S_w^{-1} S_w + S_w^{-1} S_b\} - J_{min} \\ &= \text{trace}\{S_w^{-1} S_b\} + \text{trace}\{S_w^{-1} S_w\} - J_{min} \\ &= \text{trace}\{S_w^{-1} S_b\} + J_{min} - J_{min} \\ &= \text{trace}\{S_w^{-1} S_b\} \end{aligned}$$

So for our modified J_3 score we have

$$J_3 = \text{trace}\{S_w^{-1} S_m\} - J_{min} = \text{trace}\{S_w^{-1} S_b\}$$

This modified J_3 score is very useful for feature selection, often in a classification task one will have an abundance of features, and for simplification and reducing the computational load be required to only keep the ones with most information. With this modified J_3 score we will now get a value of 0 if the classes are identical in a given feature space, for any number of dimensions. This means that we can check all or most of the combinations of features and check their J_3 score, and if one is zero or close to zero we know that the class separability is low and we don't get much information from that feature(s). If we didn't use the modified J_3 score this process would be more cumbersome, since we would have to remember how many dimensions we have for a given feature space and take that into account for the feature selection process.

1b)

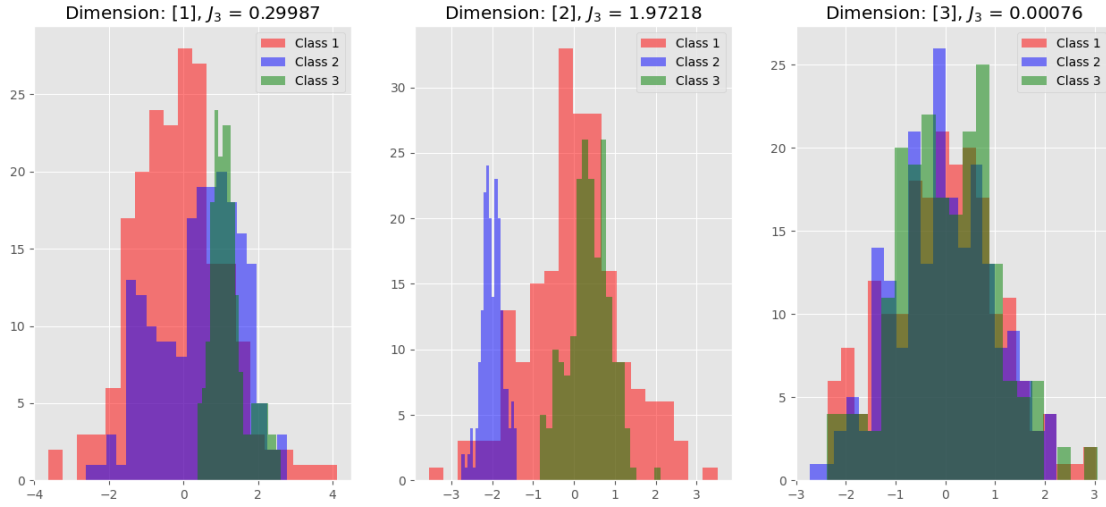


Figure 1: Histograms of the training data in all three dimensions

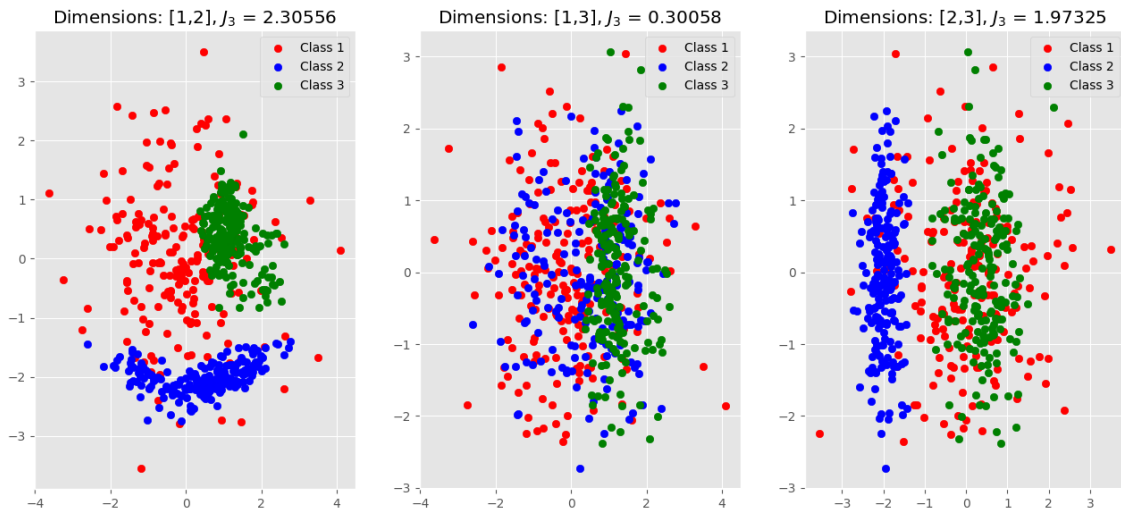


Figure 2: Scatter plots of all 2d combinations of the training data

1c)

In figure 1 we can see that the J_3 score correlates with the separability of the classes. The histogram plot of dimension 3 is a good example of this, we can see that the classes are more or less completely on top of each other, meaning the mean vectors and variance in this dimension would be close to identical for all classes. Which means we can't really separate the classes in this dimension, and the J_3 score is basically 0. J_3 is not exactly zero, but this won't really happen in the real world, but this is close enough that we can say that it is zero. On the other hand we can look at dimension 2, where we can see three distinct 'peaks', just by looking at the plot we can clearly see that the separability is higher, and that separating the classes is easier than in dimension 3. The J_3 score also backs this up, as it is much higher. Dimension 1 is somewhere in the middle of these two.

Figure 2 also shows that the J_3 score correlates with class separability. We can easily see that the most separable combination of 2 dimensions has the highest J_3 score. That being dimensions one and two. This makes sense considering our conclusion from figure 1. Since dimension 3 has the lowest separability, the combination of dimensions which don't include dimension 3 will be the one with the highest separability.

The 3-dimensional data set has a J_3 score of 2.30658, which is more or less the same as dimensions 1 and 2. Again this matches our conclusion earlier. Since the J_3 score for the 3 dimensions is more or less the same as for dimensions 1 and 2 we can get rid of the third dimension without any difference in our classification accuracy. For this reason the following tasks will be done in the 2-dimensional feature space, omitting dimension 3.

Task 2

2a)

This task was done by using the built in numpy functions for calculating the mean vectors and covariance matrices, for each class. These values from the training set were then used to calculate the Gaussian probability density function for each class, given by:

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{l/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

This Gaussian PDF was calculated for each class, and the classification was done by assigning each feature vector to the class whose PDF has the largest value at that specific point in the feature space.

Gaussian classifier		
Data set	Accuracy	Misclassifications
Training set	89.5%	63
Test set	86.8%	79

Table 1: Classification results for the Gaussian classifier

Table 1 shows the results of the Gaussian maximum likelihood classifier, as we can see the results

are quite good, with an 86.5% accuracy for the test set. As we saw in task 1 the class separability is not that high, so this result is quite good, especially for such a simple classifier.

2b)

The basic principle for the Parzen window density estimator is that we create a probability density function (pdf) that is the sum of N Gaussians, each centered on a feature vector in the training set. Where N is the number of training vectors. Other kernels than the Gaussian can be used, but for this task we are only concerned with the Gaussian. This type of pdf can estimate more or less any kind of distribution, given that the training set is large enough. The kernel itself is a symmetric Gaussian distribution, but since we take the sum of many different Gaussians, the resulting pdf can take up any shape, Gaussian or not. The expression for the Parzen pdf is:

$$\hat{p}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{(2\pi)^{\frac{l}{2}} h^l} \exp\left(-\frac{(\mathbf{x} - \mathbf{x}_i)^T (\mathbf{x} - \mathbf{x}_i)}{2h^2}\right)$$

The only free parameter here is h, which is a measure of the size of each Gaussian kernel. A low value for h will give a narrow and tall kernel. While a large h will give more flat and spread out kernel. To get a good approximation of the distribution of the training set a sufficiently small value for h should be chosen, but there is a problem with choosing h to be too small. In the limit where h goes to zero, each Gaussian will approximate the delta function, becoming infinitely tall and infinitely narrow. Meaning that we will get a perfect 100% correct estimation of the training set, but the classifier will not be able to classify anything else correctly, essentially making the classifier useless. This phenomenon where the classifier becomes too well adapted to the training set to the point where it can't classify anything else correctly is called overfitting. This means that we have to choose h carefully to get an accurate classifier, without overfitting.

2c)

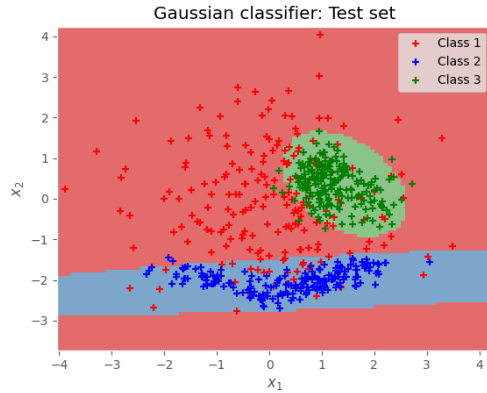
Implementing this classifier in Python is very straightforward. The PDF for each class is created by implementing the equation listed above into the code. Then the feature vectors are classified by evaluating the PDFs for each class at that feature vector, and choosing the one with the highest value

Parzen classifier		
Data set	Accuracy	Misclassifications
Training set	94.3%	34
Test set	87.3%	76

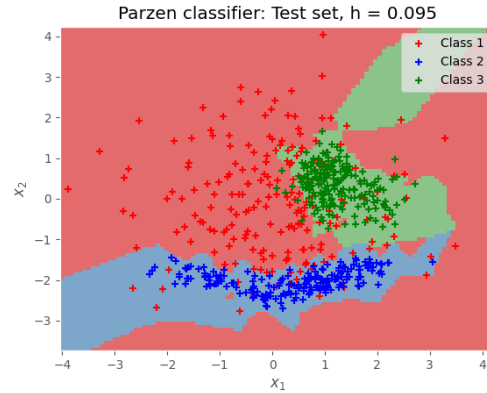
Table 2: Classification results for the Parzen window density estimator

The results of the Parzen classifier can be seen in table 2, for this result h was set to 0.095, which is the most optimal with regard to the classification accuracy for the test set. This value for h was found through trial and error. We can get a higher accuracy for the training set if we choose a low value for h, but as mentioned earlier this will cause overfitting, thus decreasing the accuracy for the test set. As we can see the results for the Parzen classifier are a bit better than the Gaussian classifier, but not by much. The test set has 3 fewer misclassifications than for the Gaussian classifier, which is good, but it is of course not a huge difference.

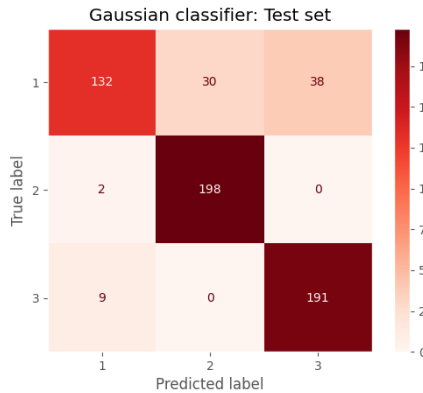
2d)



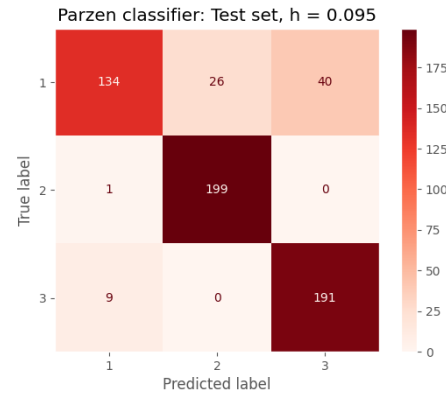
(a) Decision boundary of the Gaussian ML classifier



(b) Decision boundary of the Parzen classifier



(a) Confusion matrix for the Gaussian ML classifier



(b) Confusion matrix for the Parzen classifier

2e)

As we can see in figures 3a and 3b, the decision boundary for the Gaussian classifier is a bit simpler than the Parzen window classifier. The Gaussian classifier creates three Gaussians, one for each class, and the resulting decision boundary gives a visual understanding of the shapes of these Gaussians. The third class is very 'clumped' together, which gives a narrow and tall peak, which we can see gives a small region for class 3. The feature vectors in class 2 are placed more linearly, and narrow. This is also reflected in the decision boundary, which is a narrow band. The first class has the highest variance, meaning it's most spread out, which means that the pdf for this class is very spread out, we can see this in the decision boundary as the 'region' for class 1 is very large.

The Parzen window classifier has a rather different decision boundary, here we don't only have contiguous regions as for the Gaussian classifier. Here the decision boundary is shaped more to the specific 'shape' of the classes in the training set. The specific peculiarities of the training set are more clear here, for example outliers have more of an effect on the decision boundary. This is apparent in the top left corner, where one outlier in class 3 'stretch' the decision boundary to include the points in the top left in class 3. This doesn't help the classifier, as we can see that this region contains feature vectors from class 1. This could likely be 'fixed' by modifying h , but this would lead to misclassifications elsewhere.

Figures 4a and 4b show the confusion matrices of the two classifiers, these figures show us how many feature vectors are classified correctly, and how many are misclassified. We can also see what class the misclassified feature vectors are classified as. As we can see both classifiers are very good at classifying classes 2 and 3. Feature vectors in class 1 however often get misclassified into the other two classes.

As mentioned earlier we can cause overfitting in the Parzen classifier by choosing a small value for h . To visualize this we can look at an example of what happens when we choose h to be small. With $h = 0.001$ the decision boundary for the Parzen classifier becomes this:

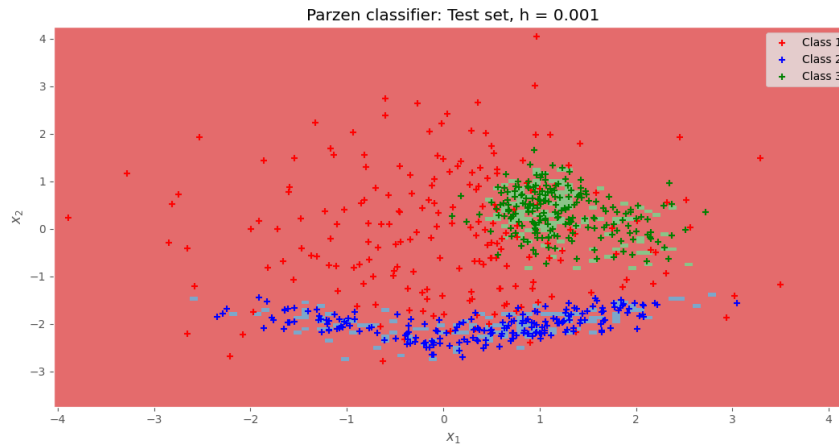


Figure 5: Decision boundary for Parzen classifier with $h=0.001$

Here we can see a good example of overfitting, the decision boundary is only focused at the feature vectors in the training set. This causes a drop in accuracy for the test set. For this classifier the accuracy for the test set is 50.2%, a huge drop from earlier. The accuracy for the training set however is at 100%. This shows the importance of choosing a good value for h .

Task 3

3a)

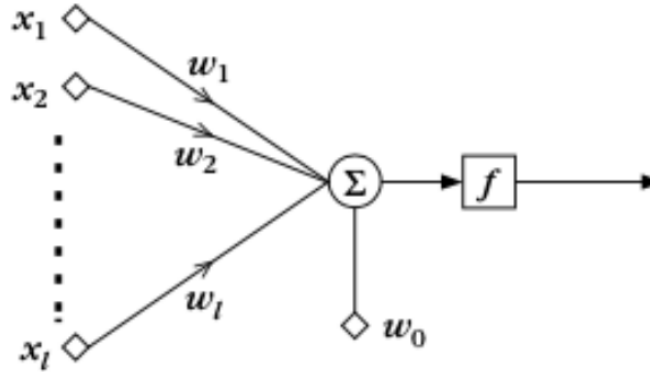


Figure 6: Perceptron node (image borrowed from exam sheet)

Figure 6 shows a Perceptron node.

First the input $\mathbf{x}=[x_1, x_2, \dots, x_l]$ gets sent to the Perceptron, for a single layer Perceptron network, or for the input layer Perceptrons of a multi-layer Perceptron, this is the feature vector we want to classify. For the hidden layers in a multi-layer network, this input will be the output of the Perceptrons in the previous layer. Then the dot product of the input with the weight vector is taken. The weight vector $\mathbf{w}=[w_1, w_2, \dots, w_l]$ will give a weighting to each feature, or input to the Perceptron. At this step we also add the threshold w_0 . The weight vector along with the threshold will define the decision hyperplane used for class separation. This step, denoted by Σ in figure 6 is the discriminant function $g_j(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$. The value from the discriminant function is then 'sent' to the activation function f . The unit step function is often used for the activation function, that is the output is 1 if the input is positive, and 0 if the input is negative. More commonly used is a sigmoid function, this is due to the discontinuity of the unit step function. In the backpropagation algorithm we take the derivative of the activation function, which is not possible for the unit step function. The sigmoid functions however are differentiable, so they are more commonly used, and what will be used later in the task (more specifically the logistic function will be used). Thus the output of the Perceptron will be a value in the range $[0,1]$, where the higher the value is, the higher likelihood that the feature vector belongs to a given class. More accurately, the output value is a measure of the distance between the feature vector and the decision hyperplane defined by $g_j(\mathbf{x})$.

3b)

This single layer Perceptron network can be coded by creating a class for the Perceptron, where the class attributes are the weights, prediction(output) and the true labels of the set of feature vectors. The methods for this class will be to run forward, which will be to calculate the output of the Perceptron. The other method will be to update the weights of the Perceptron. The expression

for updating the weights is:

$$\Delta \mathbf{w}_j = -\rho \frac{\partial J}{\partial \mathbf{w}_j}$$

Where $\Delta \mathbf{w}_j$ is the change in w_j in consecutive epochs. j denotes the weight vector for the j th Perceptron. Using the squared error cost function gives us:

$$\frac{\partial J}{\partial \mathbf{w}_j} = \frac{\partial}{\partial \mathbf{w}_j} \frac{1}{2} \sum_i^N (y(i) - \hat{y}(i))^2$$

$y(i)$ is not dependent on w_j so we can treat it as a constant during derivation, and using the chain rule we get:

$$\frac{\partial J}{\partial \mathbf{w}_j} = - \sum_i^N (y(i) - \hat{y}(i)) \frac{\partial}{\partial \mathbf{w}_j} \hat{y}(i)$$

Again, we can use the chain rule to get:

$$\begin{aligned} \frac{\partial}{\partial w_j} \hat{y}(i) &= \frac{\partial}{\partial w_j} f(w_j^T \mathbf{x}) \\ &= \frac{\partial f(w_j^T \mathbf{x})}{\partial w_j^T \mathbf{x}} \frac{\partial w_j^T \mathbf{x}}{\partial w_j} \\ &= f'(w_j^T \mathbf{x}) \mathbf{x} \end{aligned}$$

Giving us the final expression for updating the weights:

$$\Delta \mathbf{w}_j = \sum_i^N \rho (y(i) - \hat{y}(i)) f'(w_j^T \mathbf{x}) \mathbf{x}$$

This expression is very straight forward to implement in the code.

To set up the network itself we can create a class called `neural_network` which creates instances of the Perceptron class, trains them by calling on the methods in the Perceptron class T times, where T is the number of epochs. The Perceptrons will take a one vs. all approach, which means it compares the likelihood that a feature vector belongs to one class against all the others. So to interpret the output of the three perceptrons in the network we can just choose the perceptron with the highest valued output, if that corresponds to class 1, then we classify the feature vector to class 1 etc.. To get the Perceptrons to do this one vs. all approach we have to modify the training labels. We will have three Perceptrons, one for each class, each will output a measure of the likelihood that the given feature vector belongs to their class. To give each Perceptron a class we have to modify the training labels so that for each Perceptron the label for the class they are to check is 1 if the feature vector belongs to that class, and 0 else. So for example for the Perceptron comparing class 1 against the other two classes we want to modify the training labels such that the label is 1 for feature vectors belonging to class 1, and 0 for feature vectors belonging to classes 2 and 3. These training labels are $y(i)$ in the equation above.

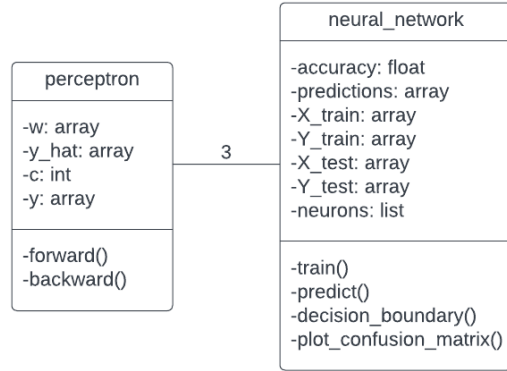
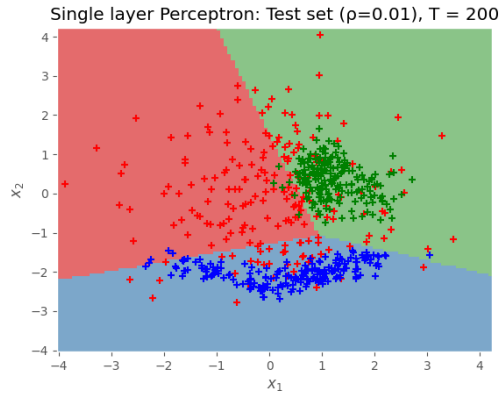


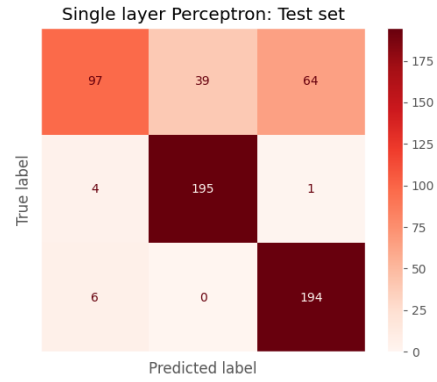
Figure 7: Simplified UML-diagram of single layer perceptron code

Figure 7 shows a simple schematic of the code used in the implementation of the single layer Perceptron network.

3c)



(a) Decision boundary for single layer Perceptron



(b) Confusion matrix for single layer Perceptron

Single layer Perceptron		
Data set	Accuracy	Misclassifications
Training set	86.7%	80
Test set	81.0%	115

Table 3: Classification results for the single layer Perceptron

Figure 8a shows the decision boundary for the single layer Perceptron network. The decision boundary looks more or less exactly as we would expect. The weight vector for each Perceptron describes a decision hyperplane (a line in this 2d case), and we have 3 Perceptrons, so we would expect to see three straight lines separating the feature vectors. Of course this data set is not linearly separable, nor is it very close to being linearly separable, so the accuracy of the network is not amazing. As table 3 shows we get an 81% accuracy with this classifier, which is pretty good. Figure 8b shows the confusion matrix for the network, and we can see that it is quite similar as for the previous classifiers in that it does a good job classifying classes 2 and 3, but struggles more on class 1. The parameters we have to adjust the network are the number of epochs T , which is the amount of times we update our weights, and we have the learning rate ρ . The learning rates decides how 'big' of a step we take when updating the weights. There is not necessarily an optimal value for the learning rate, but we have to choose it carefully depending on the classification task. If we have a big value for ρ we will approach the optimal solution faster, but we might overshoot it. Meaning that as we descend down the gradient of the cost function, we never reach the optimal solution. Choosing a small ρ we are sure to reach the optimal solution, but it will take a long time. With a small learning rate we also increase the chance of getting trapped in a local minima. Meaning we will get stuck at a suboptimal solution. For this task ρ was set to 0.01, which was the optimal value found by trial and error.

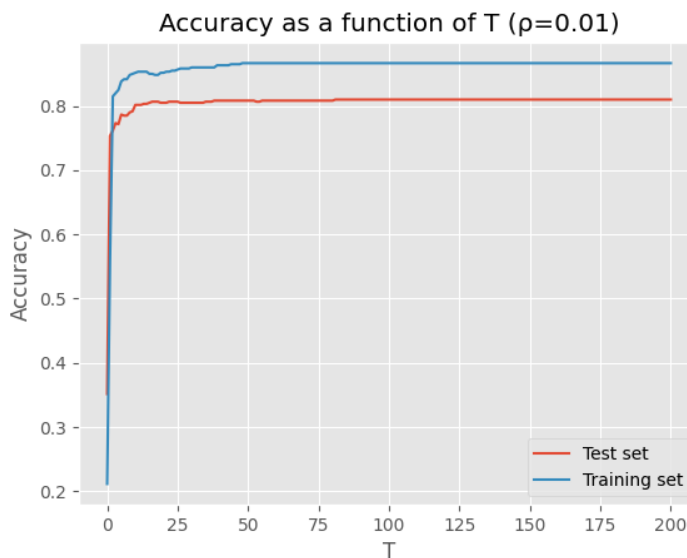


Figure 9: Accuracy as a function of the number of epochs

Figure 9 shows the accuracy of the Perceptron network as a function of the number of epochs. As we can see the accuracy begins at around 33%, which is expected since at $T=0$ the classifier is just guessing, so at $T=0$ we expect to get approximately 33% accuracy. Then it rapidly reaches around 80% for the test set, and after around 80 epochs it doesn't improve. For the training set it has reached its max value at around 50 epochs, and is constant after that. This means that we likely

have the optimal solution, since our guess doesn't improve any more.

Compared to the Bayesian classifiers from the previous task this classifier obviously doesn't perform as well. This stems from the fact that this is a linear classifier, and our data set is not linearly separable. For other data sets however the Perceptron network might perform much better than the Bayesian classifiers, but that is very dependant on the distribution of the data set. One thing that is better with this classifier than the others, especially the Gaussian classifier, is that we don't need to know anything about the distribution of the feature vectors before we classify them. For the Gaussian classifier we need knowledge about the data set to know if the Gaussian classifier will do a good job. In terms of computational efficiency the Perceptron network outperforms the Parzen window density estimator and the Gaussian classifier, given that the number of epochs is not too large. So what the Perceptron network wins out on is that it is quite general, in that it will do a good job on many types of distributions, and it is computationally efficient.

Task 4

4a)

The learning rate is, as mentioned previously, a parameter that decides the size of the step we take along the gradient of the cost function when updating the weights. A high learning rate means that we take a longer step along the gradient, and vice versa. If we choose a learning rate which is too low, then we take too long to reach the optimal solution, and we have a higher chance of getting stuck in a local minima. If we set the learning rate to be too high we risk never getting the optimal solution, and we can end up 'climbing up' the minima of the cost function. The learning rate that is just right is when we steadily approach the minima of the cost function in a reasonable amount of epochs.

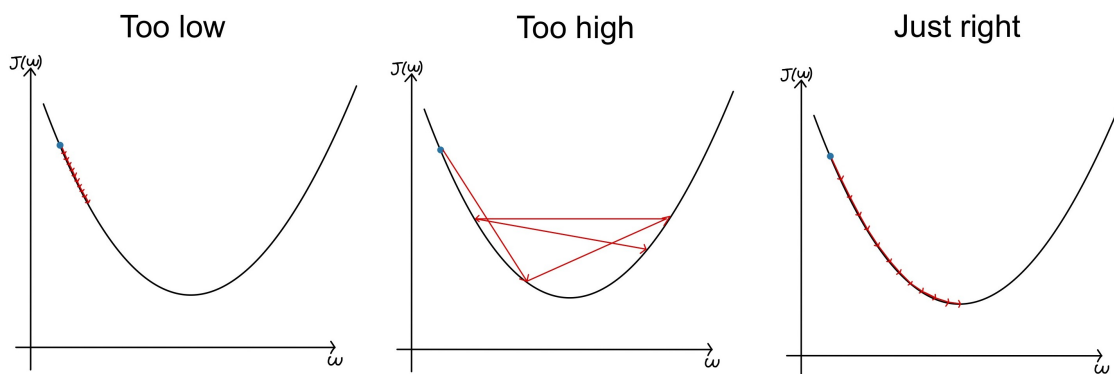


Figure 10: Gradient descent with different learning rates

Figure 10 is an attempt to illustrate the importance of choosing a correct learning rate. One could also choose an adaptive learning rate, which varies according to some predefined rule. Typically one will compare the cost function before and after backpropagation, and if it has increased we will decrease the learning rate and vice versa, but this adaptive learning rate has not been used here,

only the fixed learning rate. Note that figure 10 would be a simplified version of the cost function for a one dimensional feature space, in our two dimensional case the cost function can be thought of as a surface in 3d space, as opposed to the 2d function shown in the figure.

Momentum is often used in calculating the new weights, this is done by adding a new term to the updated weights which is dependant on the previous change in weights. What this does is that it smoothes out the oscillations we get on approach to the minima of the cost function, these oscillations are seen in figure 10 with high learning rate. So we get a smoothing of the oscillations, and faster convergence to the minimum. As T increases, the momentum term will decrease, which makes the smoothing effect of the momentum apparent. Since the momentum term is dependant on the change in weights in the previous iteration, momentum will also help get the network out of small local minima. The new expression for the updated weights for the jth Perceptron in layer r is given by:

$$\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) + \alpha \Delta \mathbf{w}_j^r(\text{old}) - \rho \sum_i^N \delta_j^r(i) \mathbf{y}^{r-1}(i)$$

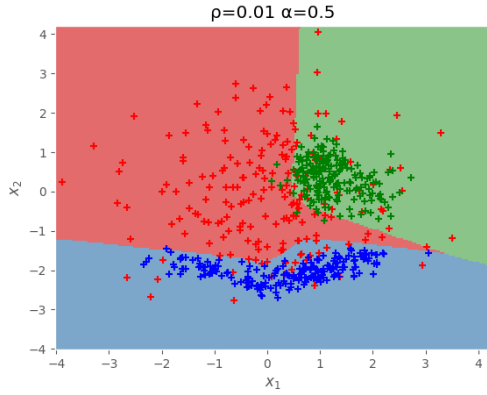
$\Delta \mathbf{w}_j^r(\text{old})$ here refers to the last two terms in the equation, for the previous iteration.

The final term is the 'gradient descent term', which is the gradient of the cost function multiplied by the learning rate ρ . The second term is the momentum term, which is the previous update in the weights, multiplied by α . α is the momentum factor, it is just a number between 0 and 1 which defines how much 'impact' the momentum factor has on the training of the network.

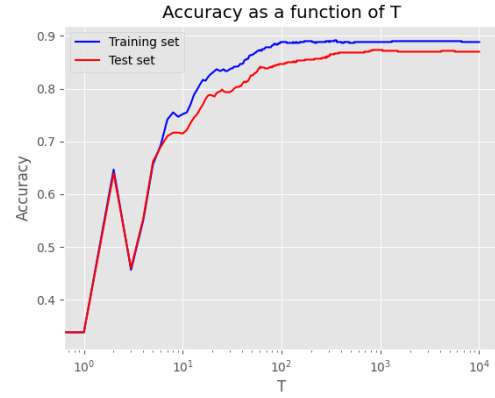
4b)

(Note: A quick explanation of the code used for task 4 is given at the end of this section.)

To see the effect of network complexity we can create different architectures and analyze the resulting decision boundaries. These following figures show the results of classifying the data sets with learning rate set to 0.01, momentum factor set to 0.5. For simpler notation the architecture of the neural network is written as a list where each number in the list represents the number of nodes in the respective layer, e.g. [4,3] means two hidden layers with 4 nodes in the first layer and 3 nodes in the second layer. The output layer with 3 output nodes comes in addition to these (As well as the input layer), the list only refers to the hidden layers.

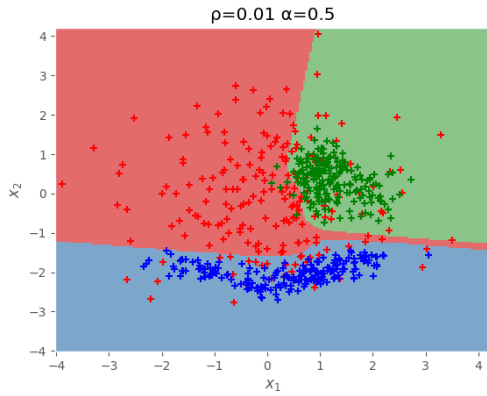


(a) Decision boundary for network with one hidden layer with 3 nodes

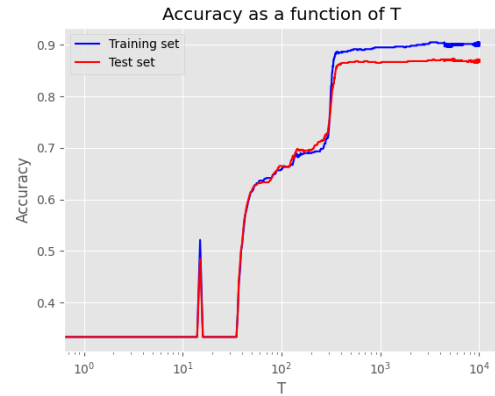


(b) Accuracy vs Epoch for network with one hidden layer with 3 nodes

Figure 11: [3]

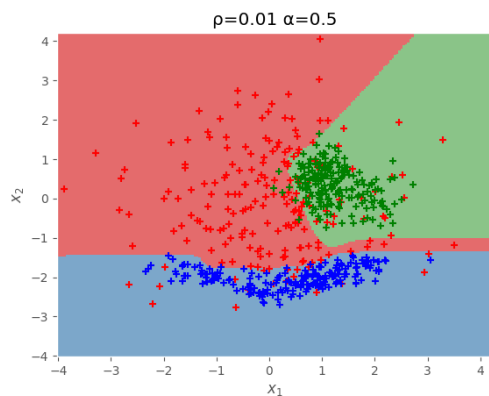


(a) Decision boundary for network with three hidden layers with 3 nodes each

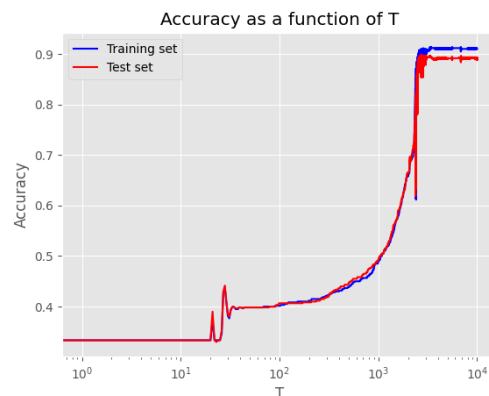


(b) Accuracy vs Epoch for network with three hidden layers with 3 nodes each

Figure 12: [3, 3, 3]

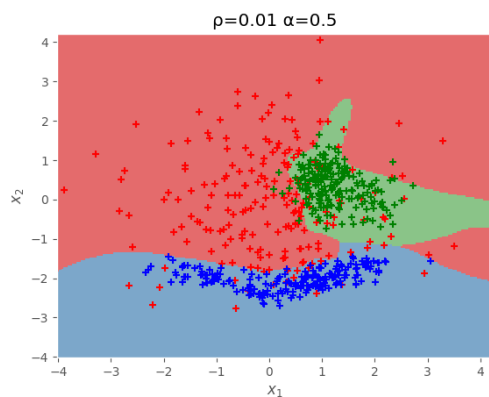


(a) Decision boundary for network with five hidden layers with 3 nodes each

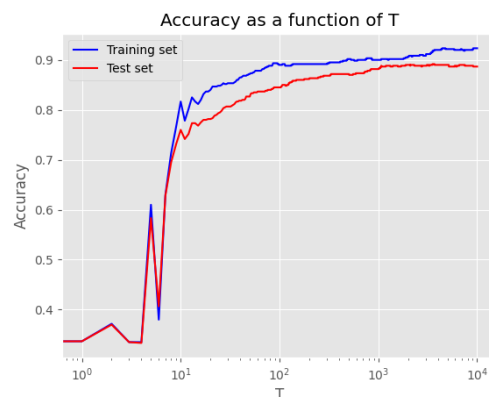


(b) Accuracy vs Epoch for network with five hidden layers with 3 nodes each

Figure 13: [3, 3, 3, 3, 3]



(a) Decision boundary for network with one hidden layer with 15 nodes



(b) Accuracy vs Epoch for network with one hidden layers with 15 nodes

Figure 14: [15]

Multi-layer Perceptron network			
Architecture	Data set	Accuracy	Misclassifications
[3]	Training	88.8%	67
	Test	87.0%	78
[3, 3, 3]	Training	90.0%	60
	Test	86.7%	80
[3, 3, 3, 3, 3]	Training	91.0%	54
	Test	88.8%	67
[15]	Training	92.3%	46
	Test	88.7%	68

Table 4: Classification results for different network complexities

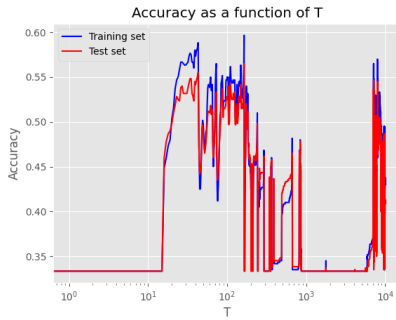
We can learn a lot by viewing these figures. The most obvious conclusion we can draw comes from the complexity of the decision boundaries for the different networks. From figures 11a, 12a and 13a we can see that as the number of hidden nodes increases we get a more complex decision boundary. This makes sense if we think back to the single layer Perceptron network, where each discriminant function describes a line (hyperplane in the general l dimensional case) that separates the feature space. More nodes should, by the same logic, give more lines to divide the feature space giving us regions that are more complex, and not just separated by single lines. Meaning we have a network that can separate classes that are not linearly separable. By looking at the accuracies of the network as a function of T we can see the accuracy steadily increases, and after a while reaches what appears to be an asymptotic value. This tells us two things, first that the network has found a minima and is 'stuck' there, thus the accuracy won't increase much with more epoch. The other thing we can learn from this plot is that our learning rate is quite good, if it were too low we would see a shallower rise in accuracy, and a higher learning rate would give us more erratic behaviour. One very interesting thing we can see in many of the figures, e.g. fig 12b is that we can clearly see that the network ended up in a small local minima, and then climb out of that local minima. It is very interesting to actually see that this happens in our code, fig 11b is especially interesting, since here we can see that the accuracy jumped to 66% only to drop down again right after. What this means is that the network found a solution where two classes (probably 2 and 3) are classified accurately, but no vectors in class 1 are correctly classified. Then the network is able to get out of this local minima and find a much better solution.

Now we know that the decision boundary gets more complex with more hidden nodes, but what about how they are distributed in the hidden layers? To answer this we can look at figures 13a and 14a. These networks have the same amount of nodes, but distributed differently, and the resulting decision boundaries are very different. It is clear that the decision boundary in figure 14a is more complex than figure 13a, even though the number of hidden nodes is the same. This points to having many nodes in one layer causes further complexity in the decision boundary.

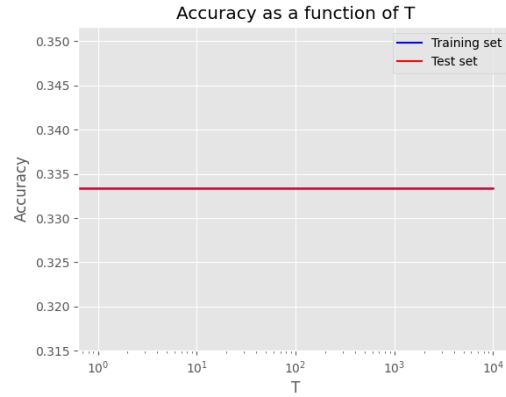
The fact that the complexity of the decision boundary increases as we increase the number of nodes means that we run the risk of overfitting if we have too many nodes and many epochs. So we need to be aware of how the neural network is designed.

4c)

For this task the architecture of the network will be set to have only one hidden layer with 15 nodes, the momentum factor has been set to 0.5 for all figures. Since the accuracy vs epoch plot is most interesting for comparing the learning rates, this plot will be the only shown in this section to save space and time, but the accompanying plots, as well as for more values of ρ will be listed in the appendix.



(a) Accuracy vs epochs for $\rho = 1$



(b) Accuracy vs epochs for $\rho = 10^{-6}$

Figure 15a is an excellent example of what happens when the learning rate is set too high. Here we can see very erratic behaviour as the classifier is 'bouncing around' the cost function. There is no convergence taking place here, only jumping around erratically. The accuracy on the test set for these parameters was only 41.0%, which is extremely low compared to what we got with a good learning rate.

Figure 15b shows the complete other end of the spectrum. Here the learning rate is so low that the change in the weights being made each iteration is too small to make any progress towards a solution. That is why we see no improvement over 10 000 epochs for this classifier, it is just stuck in the case where all feature vectors are classified as the same class.

These two cases are of course very extreme, for more 'moderate' values of ρ we would expect to see a steeper gradient in the accuracy for a higher learning rate, until we reach the point where we see the erratic behaviour shown in figure 15a. Both of these figures match what we would expect to see, and they both match figure 10.

We can also take a look at how the momentum factor affects the accuracy as a function of epochs. As mentioned previously the momentum dampens the oscillations that happen near the minima of the cost function, thus speeding up convergence. So we should expect to see this if we vary the momentum factor while keeping all other parameters constant. (Only the accuracy vs epoch will be plotted here, since the decision boundary does not contain much information in the sense of the effect of the momentum.)

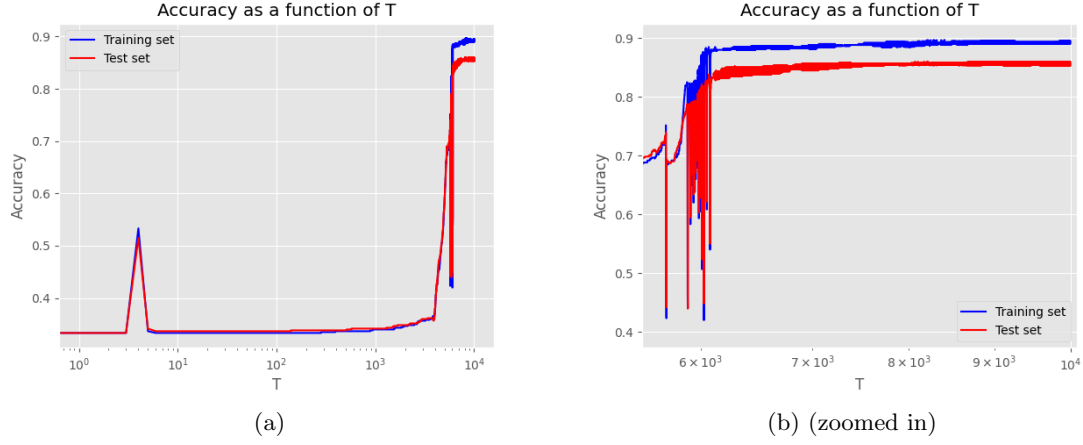


Figure 16: Accuracy for $\alpha = 0$

Figure 16 shows us how much of an effect the momentum has on the convergence of the network. We can very clearly see a lot of oscillations here that are not present when we have momentum (such as the plots in task 4b). This shows us how important the momentum is since we don't really get any proper convergence without momentum, meaning we don't get as good a result for our classifier.

4d)

The best performing Neural Network for this task is $[10, 8, 6]$ with $\rho = 0.008$ and $\alpha = 0.65$. This design was found by trial and error. Ideally a technique for optimizing the design should have been used, like e.g. pruning, but this was not implemented due to time. The result for this design is:

Multi-layer Perceptron		
Data set	Accuracy	Misclassifications
Training	92.5%	45
Test	90.2%	59

Table 5: Result for optimal multi-layer Perceptron

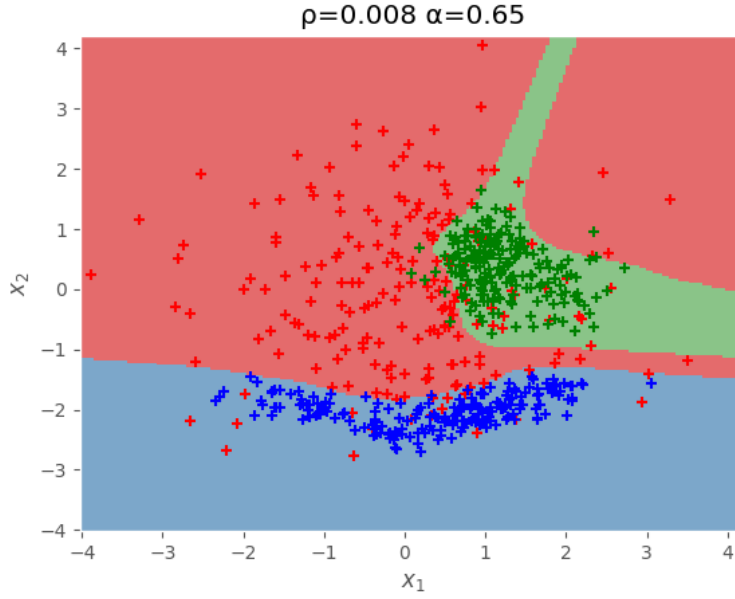
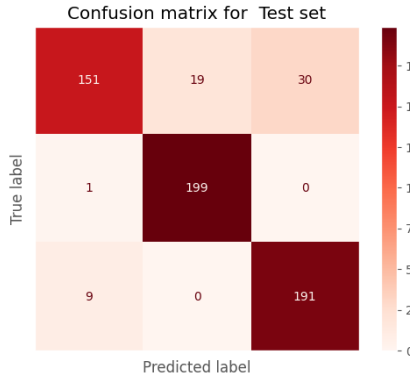
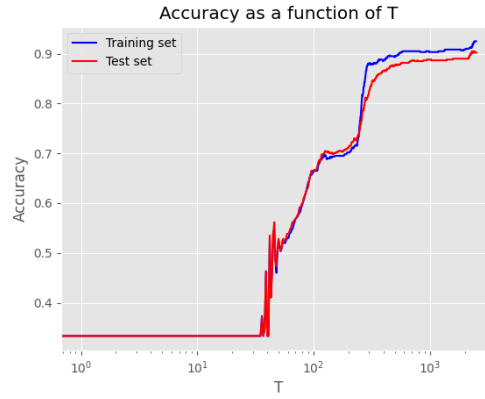


Figure 17: Decision boundary for optimal neural network



(a) Confusion matrix for optimal neural network



(b) Accuracy plot for optimal neural network

As we can see this network performs extremely well, getting an accuracy of 90.2% on the test set after 2500 epochs. This design performs so well since it has many hidden nodes, so it can create more complex decision regions, but not too much complexity as to cause overfitting. Also the parameters (ρ , α) are tuned for this network to perform optimally.

There are probably more optimal designs for this network, but in the short time span of this exam no better design was found. One limitation in finding a better design is the computation time for this network, especially with many nodes this network takes quite a long time to run. The most

optimal solution possible with such a network is probably a few percent better than this design (speculation).

This result is the best of all the classifiers used on this data set, by a decent margin. The training set had a higher accuracy with the Parzen classifier, but obviously we don't care as much about this since this was when we got overfitting, causing the classification of the test set to drop drastically. To properly compare all classifiers we can compare their decision boundaries, as well as their performance:

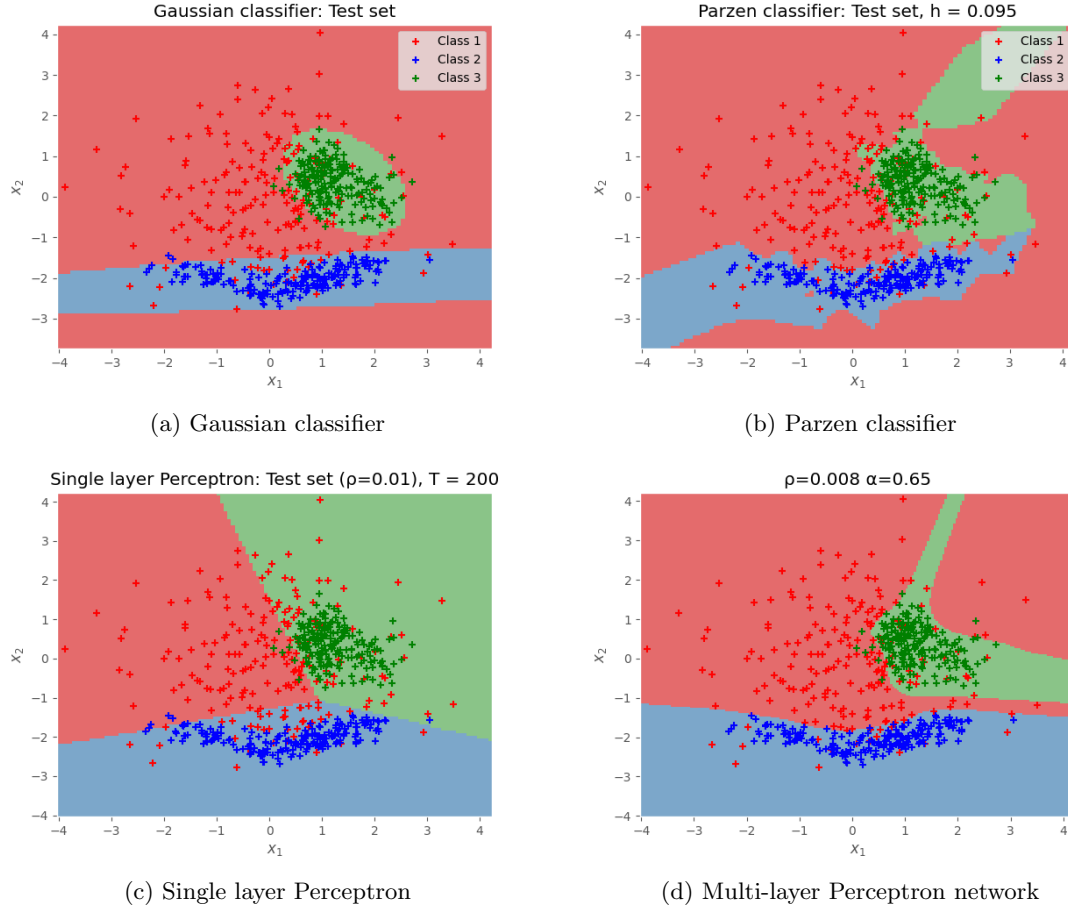


Figure 19: Decision boundaries for all classifiers

Comparison of classifiers				
Classifier	Data set	Accuracy	Misclassifications	Time (s)
Gaussian	Training	89.5%	63	$1.42 * 10^{-2}$
	Test	86.8%	79	
Parzen	Training	94.3%	34	$1.936 * 10^0$
	Test	87.3%	76	
Perceptron	Training	86.7%	80	$3.344 * 10^{-2}$
	Test	81.0%	115	
MLP-network	Training	92.5%	45	$2.016 * 10^1$
	Test	90.2%	59	

Table 6: Comparison of all classifiers used on the data set

Here we can clearly see performance of all the classifiers and compare them properly against each other. The most accurate classifier is as mentioned the Multi layer neural network, with an accuracy of 90.2%. The decision boundary for this classifier is a bit complicated, especially when compared to the Gaussian classifier and the single layer Perceptron. The second best classifier is the Parzen classifier, followed closely by the Gaussian, and then the single layer Perceptron trails a bit behind. The difference in number of misclassified feature vectors for the three best classifiers is not that really that big at only 20 feature vectors. However the data set is pretty small, so we should look at the percentage difference, and a 2.9% difference (between MLP and Parzen) is quite a bit. Another very interesting point of comparison is the computation time for each classifier. We can see in table 6 that the single layer Perceptron and Gaussian classifier are by far the quickest, both taking on the order of 10ms to classify the test set. The Parzen classifier takes quite a bit more time, 2 orders of magnitude more, taking roughly 2 seconds to classify the data set. Then comes the MLP-network, using approximately 20 seconds to classify the test data. This is a huge difference in time for the classifiers, the MLP-network is over 1000 times slower than the Gaussian classifier. For this relatively small data set the difference doesn't matter that much, but if the data set was much larger, and the network more complicated, then this difference in magnitude is huge. 3 orders of magnitude doesn't necessarily sound like much, but it is the difference of running the code for 1 day and running it for around 3 years. So computation time is really important, and especially for large data sets one has to take computation time into account when designing the classifier.

The easiest classifier to code was by far the Gaussian ML classifier, this classifier is about as easy as it gets and is very straight forward to implement. The Parzen classifier was also relatively easy, but trickier than the Gaussian classifier. The multi-layer network was by far the hardest to create, mainly because the backpropagation algorithm is very tricky to properly implement due to its complex nature, especially since there are so many indices to keep track of. The multi-layer network is also a bit more tricky to use since there are more free parameters, and if one doesn't know how they work then it can be tricky to get good results in the beginning. My personal favorite is the multi-layer neural network. The reason being that it is very general, and works on many types of classification problems. Also it was extremely satisfying to get it working. The single layer network is in a close second place, mainly since the concept of finding the ideal line (hyperplane) to separate the feature vectors is so simple, yet very effective. Also the concept that these two networks 'learn' is very fascinating and super interesting.

The code for the MLP-network

The multi-layer Perceptron network was coded in a similar way to the single layer network in task 3. First a class for the Perceptrons was created, these Perceptrons have all the necessary attributes and methods for both forward, and backward propagation. A new class for the output Perceptrons, which inherits from the Perceptron class was also made which is almost the same, but contains different methods for calculating the error and delta. This was done since the expressions for calculating these values are different for the output layer Perceptrons, than for the hidden layers. Then a class for the network itself is created, which contains the Perceptrons, stored in a nested list to keep track of which Perceptrons are in which layer. The forward- and backpropagation are done as methods in this class which implements the respective algorithms as they are explained in the textbook/ lecture notes. The most tricky part here is to keep track of which node is which, and not messing up the indices (I won't go into more detail on the backpropagation since it is a straightforward implementation of the steps laid out in the textbook). The training, prediction etc. is done in the same way as in task 3. A simplified UML diagram of the code used is given here:

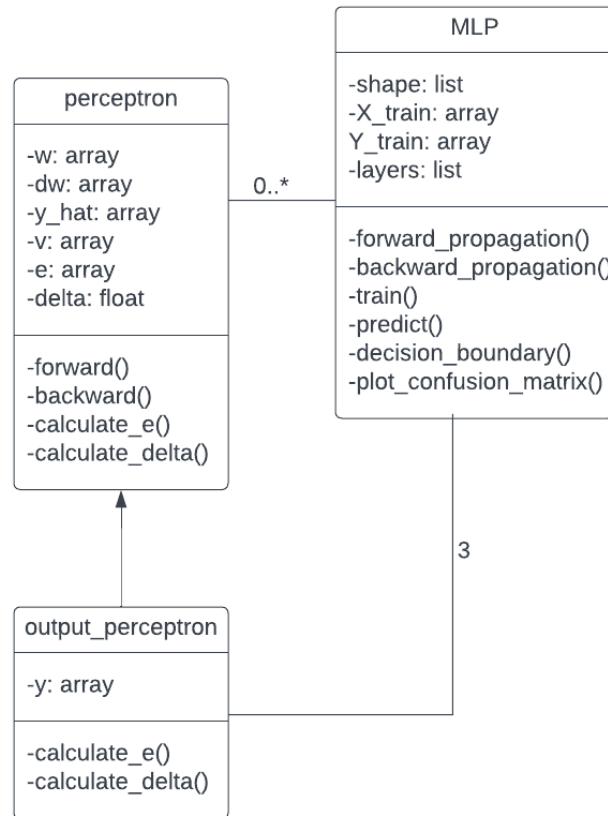


Figure 20: UML diagram of code used in task 4

Appendix

Additional plots for task 4c

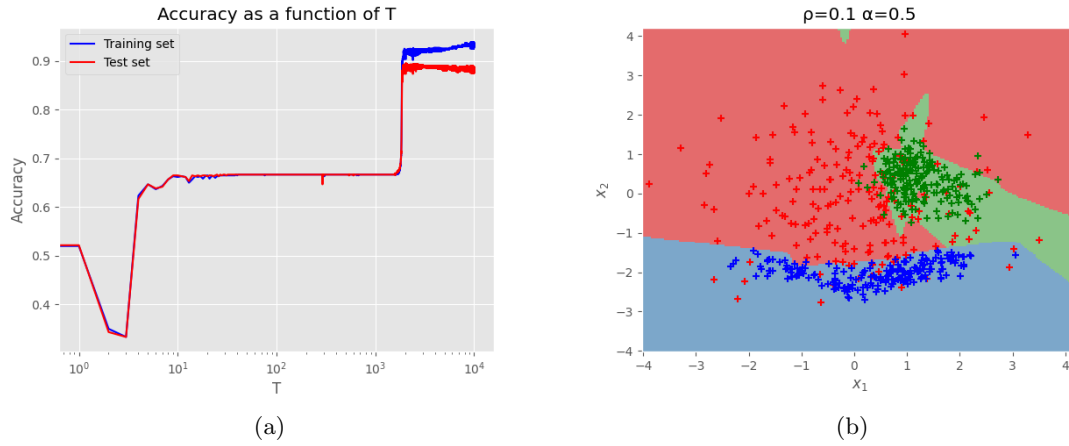


Figure 21: Accuracy and decision boundary for $\rho = 0.1$

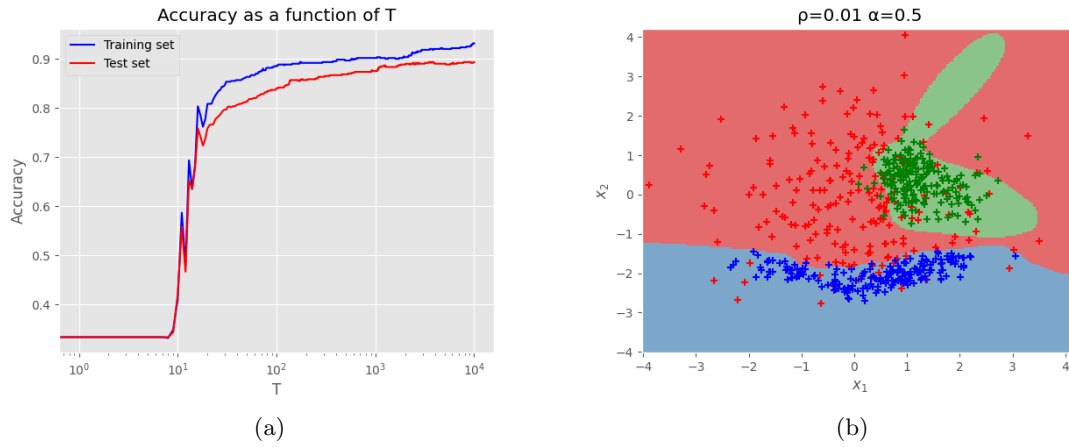
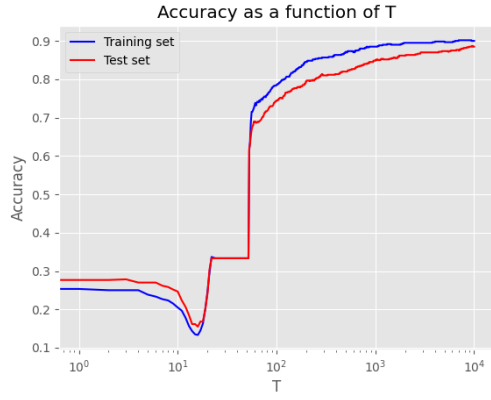
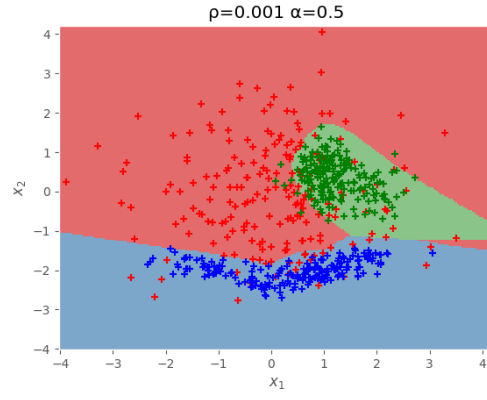


Figure 22: Accuracy and decision boundary for $\rho = 0.01$

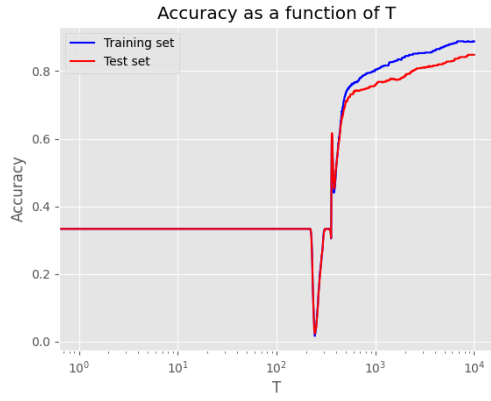


(a)

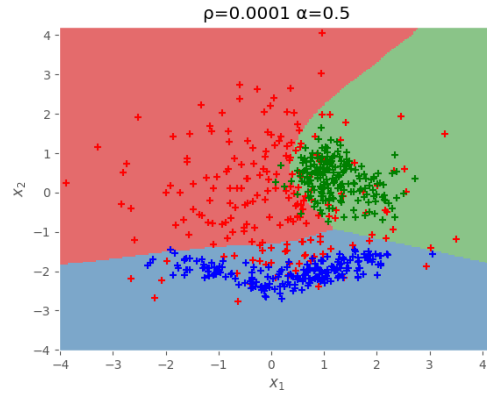


(b)

Figure 23: Accuracy and decision boundary for $\rho = 0.001$



(a)



(b)

Figure 24: Accuracy and decision boundary for $\rho = 0.0001$

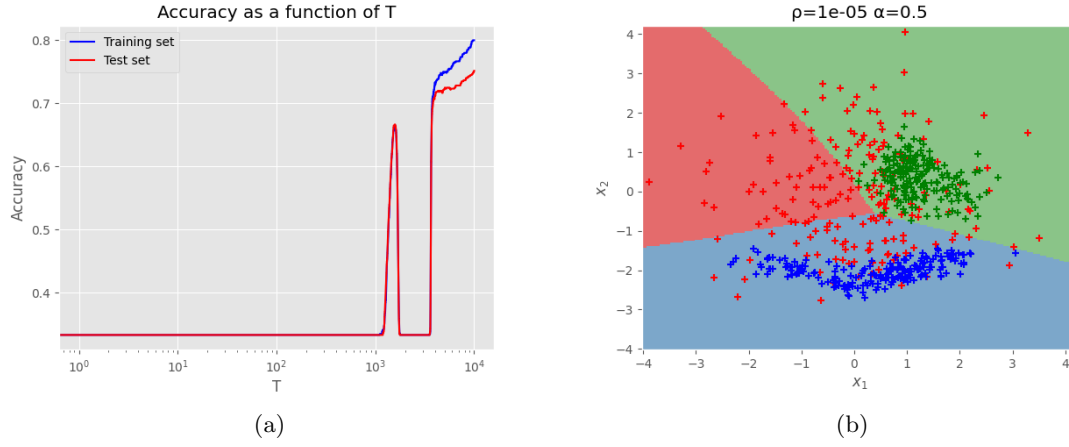


Figure 25: Accuracy and decision boundary for $\rho = 0.00001$

Multi-layer network ([15]) ($\alpha = 0.5$)			
Learning rate	Data set	Accuracy	Misclassifications
10^{-6}	Training	33.3%	400
	Test	33.3%	400
10^{-5}	Training	80.0%	120
	Test	75.2%	150
10^{-4}	Training	88.8%	67
	Test	84.8%	91
10^{-3}	Training	90.0%	60
	Test	88.5%	69
10^{-2}	Training	93.2%	41
	Test	89.3%	64
10^{-1}	Training	93.3%	40
	Test	88,7%	68
0.5	Training	48.3%	310
	Test	48.8%	307
1	Training	43.0%	342
	Test	41.0%	354

Table 7: Result of neural network with varying learning rate

Code

(All the code used for this exam is also available in a zipped folder delivered along with this PDF)

Tasks 1 & 2

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 plt.style.use('ggplot')
4 import scipy as sc
5 from sklearn import metrics
6 from prettytable import PrettyTable
7
8 # Importing the data
9 dataset = sc.io.loadmat('Home-exam\ExamData3D.mat')
10
11 # Splitting the data into training and test sets
12 X_train = dataset['X_train']
13 X_test = dataset['X_test']
14 Y_train = dataset['Y_train'][0]
15 Y_test = dataset['Y_test'][0]
16
17 # defining the probability for the classes
18 P = 1/3 # using the a priori knowledge that the classes are equiprobable (to do
19         properly it should be calculated from the data set)
20
21 def mean_vector(x, y, class_nr):
22     """
23     Calculates the mean vector for a given class
24
25     Parameters
26     -----
27     - x : The set of feature vectors
28     - y: Known labels for the feature vectors
29     - class_nr: The class for which the mean vector is calculated
30
31     Returns
32     -----
33     - mu: The mean vector for the given class
34     """
35     # special case for 1 dimension (only difference is the indexing of the array)
36     if x.ndim == 1:
37         mu = np.mean(x[y == class_nr]) # calculating the mean vector for the given
38         class
39     else:
40         mu = np.mean(x[:,y == class_nr], axis = 1) # calculating the mean vector for
41         the given class
42
43     return np.matrix(mu).T # returning the mean vector as a column matrix
44
45 def J3_score(x, y):
46     """
47     Calculates the modified J3 score for a given data set
48
49     Parameters
50     -----
```

```

49     - x: The set of feature vectors
50     - y: Known labels for the feature vectors
51
52     Returns
53     -----
54     - J3: The modified J3 score for the given data set
55     """
56
57     # calculating the mean vector for each class
58     mu1 = mean_vector(x, y, 1)
59     mu2 = mean_vector(x, y, 2)
60     mu3 = mean_vector(x, y, 3)
61
62     # calculating the global mean vector
63     mu0 = np.matrix((mu1 + mu2 + mu3)*P)
64
65     # calculating the covariance matrix for each class
66     if x.ndim == 1:
67         cov1 = np.cov(np.array((x[y == 1])))
68         cov3 = np.cov(np.array((x[y == 3])))
69         cov2 = np.cov(np.array((x[y == 2])))
70     else:
71         cov1 = np.cov(np.array((x[:,y == 1])))
72         cov3 = np.cov(np.array((x[:,y == 3])))
73         cov2 = np.cov(np.array((x[:,y == 2])))
74
75     # calculating the within class scatter matrix
76     Sw = (cov1 + cov2 + cov3)*P
77
78     # calculating the between class scatter matrix
79     Sb = ((mu1-mu0)@((mu1-mu0).T) + (mu2-mu0)@((mu2-mu0).T) + (mu3-mu0)@((mu3-mu0).T))*P
80
81     # calculating the modified J3 score
82     if x.ndim == 1:
83         J3 = Sb/Sw
84     else:
85         J3 = np.trace(np.linalg.inv(Sw)@Sb)
86     return J3
87
88 def print_J3():
89     """
90     Prints the J3 score for each dimension and combination of dimensions
91     """
92     print("-----Task 1-----")
93     table = PrettyTable()
94     table.field_names = ["Dimension", "J3 score"]
95     table.add_row(["[1]", np.round(float(J3_score(X_train[0], Y_train)), 3)])
96     table.add_row(["[2]", np.round(float(J3_score(X_train[1], Y_train)), 3)])
97     table.add_row(["[3]", np.round(float(J3_score(X_train[2], Y_train)), 3)])
98     table.add_row(["[1,2]", np.round(float(J3_score(np.delete(X_train, 2, 0), Y_train)), 3)])
99     table.add_row(["[1,3]", np.round(float(J3_score(np.delete(X_train, 1, 0), Y_train)), 3)])
100    table.add_row(["[2,3]", np.round(float(J3_score(np.delete(X_train, 0, 0), Y_train)), 3)])
101    table.add_row(["[1,2,3]", np.round(float(J3_score(X_train, Y_train)), 5)])
102    print(table)

```

```

103
104 def plot_data():
105     """
106     Plots all combinations of dimensions of the training set
107     """
108     plt.subplot(1,3,1)
109     plt.hist(X_train[0][Y_train == 1], bins = 20, alpha = 0.5, label = 'Class 1',
110             color = 'red')
111     plt.hist(X_train[0][Y_train == 2], bins = 20, alpha = 0.5, label = 'Class 2',
112             color = 'blue')
113     plt.hist(X_train[0][Y_train == 3], bins = 20, alpha = 0.5, label = 'Class 3',
114             color = 'green')
115     plt.legend()
116     plt.title('Dimension: [1], $J_3$ = ' + str(np.round(float(J3_score(X_train[0],
117     Y_train)), 5)))
118
119     plt.subplot(1,3,2)
120     plt.hist(X_train[1][Y_train == 1], bins = 20, alpha = 0.5, label = 'Class 1',
121             color = 'red')
122     plt.hist(X_train[1][Y_train == 2], bins = 20, alpha = 0.5, label = 'Class 2',
123             color = 'blue')
124     plt.hist(X_train[1][Y_train == 3], bins = 20, alpha = 0.5, label = 'Class 3',
125             color = 'green')
126     plt.legend()
127     plt.title('Dimension: [2], $J_3$ = ' + str(np.round(float(J3_score(X_train[1],
128     Y_train)), 5)))
129
130     plt.subplot(1,3,3)
131     plt.hist(X_train[2][Y_train == 1], bins = 20, alpha = 0.5, label = 'Class 1',
132             color = 'red')
133     plt.hist(X_train[2][Y_train == 2], bins = 20, alpha = 0.5, label = 'Class 2',
134             color = 'blue')
135     plt.hist(X_train[2][Y_train == 3], bins = 20, alpha = 0.5, label = 'Class 3',
136             color = 'green')
137     plt.legend()
138     plt.title('Dimension: [3], $J_3$ = ' + str(np.round(float(J3_score(X_train[2],
139     Y_train)), 5)))
140     plt.show()
141
142     # plotting the 2D data
143     plt.subplot(1,3,1)
144     plt.scatter(X_train[0][Y_train == 1], X_train[1][Y_train == 1], label = 'Class 1
145     ', color = 'red')
146     plt.scatter(X_train[0][Y_train == 2], X_train[1][Y_train == 2], label = 'Class 2
147     ', color = 'blue')
148     plt.scatter(X_train[0][Y_train == 3], X_train[1][Y_train == 3], label = 'Class 3
149     ', color = 'green')
150     plt.legend()
151     plt.title('Dimensions: [1,2], $J_3$ = ' + str(np.round(float(J3_score(np.delete(
152     X_train, 2, 0), Y_train)), 5)))
153
154     plt.subplot(1,3,2)
155     plt.scatter(X_train[0][Y_train == 1], X_train[2][Y_train == 1], label = 'Class 1
156     ', color = 'red')
157     plt.scatter(X_train[0][Y_train == 2], X_train[2][Y_train == 2], label = 'Class 2
158     ', color = 'blue')
159     plt.scatter(X_train[0][Y_train == 3], X_train[2][Y_train == 3], label = 'Class 3
160     ', color = 'green')

```

```

142 plt.legend()
143 plt.title('Dimensions: [1,3], $J_3$ = ' + str(np.round(float(J3_score(np.delete(
144 X_train, 1, 0), Y_train)), 5)))
145
146 plt.subplot(1,3,3)
147 plt.scatter(X_train[1][Y_train == 1], X_train[2][Y_train == 1], label = 'Class 1
148 ', color = 'red')
149 plt.scatter(X_train[1][Y_train == 2], X_train[2][Y_train == 2], label = 'Class 2
150 ', color = 'blue')
151 plt.scatter(X_train[1][Y_train == 3], X_train[2][Y_train == 3], label = 'Class 3
152 ', color = 'green')
153 plt.legend()
154 plt.title('Dimensions: [2,3], $J_3$ = ' + str(np.round(float(J3_score(np.delete(
155 X_train, 0, 0), Y_train)), 5)))
156 plt.show()
157 plt.clf()
158
159 print("The difference between dimensions 1,2 and 3d is: ",
160       np.round(float(J3_score(X_train, Y_train))-float(J3_score(np.delete(X_train,
161 2, 0), Y_train)), 5),
162       "which is so small that we can get rid of dimensin 3 without losing much
163 information.")
164
165 class Gaussian:
166     def __init__(self, X_train, Y_train):
167         """Gaussian classifier for 3 classes"""
168         self.X_train = X_train
169         self.Y_train = Y_train
170         self.accuracy = 0
171
172     def mean_vector(self):
173         """Returns the mean vector for each class"""
174         mean_1 = mean_vector(X_train, Y_train, 1)
175         mean_2 = mean_vector(X_train, Y_train, 2)
176         mean_3 = mean_vector(X_train, Y_train, 3)
177         return mean_1, mean_2, mean_3
178
179     def covariance_matrix(self):
180         """Returns the covariance matrix for each class"""
181         cov_1 = np.cov(np.array((X_train[:,Y_train == 1])))
182         cov_2 = np.cov(np.array((X_train[:,Y_train == 2])))
183         cov_3 = np.cov(np.array((X_train[:,Y_train == 3])))
184         return cov_1, cov_2, cov_3
185
186     def pdf(self, x):
187         """
188         Calculates the probability density function for each class
189
190         Parameters
191         -----
192         - x: The set of feature vectors
193
194         Returns
195         -----
196         - p1, p2, p3: The probability density function for each class
197         """
198         # calculating the mean vector and covariance matrix for each class

```

```

193     mean_1, mean_2, mean_3 = self.mean_vector()
194     cov_1, cov_2, cov_3 = self.covariance_matrix()
195
196     # calculating the probability for each class (using the multivariate
197     gaussian distribution)
198     denominator = (2*np.pi)**(len(mean_1)/2)*np.linalg.det(cov_1)**(1/2)
199     exponent = -0.5*(x-mean_1).T@np.linalg.inv(cov_1)@(x-mean_1)
200     p1 = np.diag(np.exp(exponent))/denominator # note: we take the diagonal
201     since the other values don't matter for us
202
203     denominator = (2*np.pi)**(len(mean_2)/2)*np.linalg.det(cov_2)**(1/2)
204     exponent = -0.5*(x-mean_2).T@np.linalg.inv(cov_2)@(x-mean_2)
205     p2 = np.diag(np.exp(exponent))/denominator
206
207     denominator = (2*np.pi)**(len(mean_3)/2)*np.linalg.det(cov_3)**(1/2)
208     exponent = -0.5*(x-mean_3).T@np.linalg.inv(cov_3)@(x-mean_3)
209     p3 = np.diag(np.exp(exponent))/denominator
210
211     return p1, p2, p3
212
213 def predict(self, X, Y, accuracy = True):
214     """
215     Predicts the class for each feature vector in X
216
217     Parameters
218     -----
219     - X: The set of feature vectors
220     - Y: Known labels for the feature vectors
221     - accuracy: If True, the accuracy of the prediction is calculated and stored
222     in self.accuracy
223
224     Returns
225     -----
226     - pred: The predicted class for each feature vector in X
227     """
228     p1, p2, p3 = self.pdf(X)
229     pred = np.argmax(np.array([p1, p2, p3]), axis = 0) + 1
230     if accuracy == True:
231         self.accuracy = np.sum(pred == Y)/len(Y)
232     return pred
233
234 def plot(self, X, Y, title, resolution = 100):
235     """
236     Plots the dataset and the decision boundary
237
238     Parameters
239     -----
240     - X: The set of feature vectors
241     - Y: Known labels for the feature vectors
242     - title: The title of the plot (Gaussian classifier for "title")
243     - resolution: The resolution of the plot
244     """
245     x, y = np.linspace(-4.0, 4.2, num=resolution), np.linspace(-3.7, 4.2, num=
246     resolution)
247     XY = np.asarray(np.meshgrid(x,y)).reshape(2, -1)
248
249     pred = self.predict(XY, Y, accuracy = False)

```

```

247     plt.pcolormesh(x, y, pred.reshape(resolution,resolution), cmap = 'Set1',
alpha=0.6, vmin=0, vmax=12)
248     plt.scatter(X[0][Y == 1], X[1][Y == 1],marker='+', label = 'Class 1', color
= 'red')
249     plt.scatter(X[0][Y == 2], X[1][Y == 2],marker='+', label = 'Class 2', color
= 'blue')
250     plt.scatter(X[0][Y == 3], X[1][Y == 3],marker='+', label = 'Class 3', color
= 'green')
251     plt.xlabel('$x_1$')
252     plt.ylabel('$x_2$')
253     plt.legend()
254     plt.title('Gaussian classifier:' + title)
255     plt.show()
256
257     def plot_confusionmatrix(self, X, Y, title):
258         pred = self.predict(X, Y, accuracy = False)
259         confusion_matrix = metrics.confusion_matrix(Y, pred)
260         confusion_matrix_display = metrics.ConfusionMatrixDisplay(confusion_matrix,
display_labels = [1,2,3])
261         confusion_matrix_display.plot(cmap="Reds")
262         plt.title("Gaussian classifier: " + title)
263         plt.grid(False)
264         plt.show()
265         plt.clf()
266
267     class Parzen:
268         def __init__(self, X_train, Y_train):
269             """Parzen window density estimator for 3 classes"""
270             self.X_train = X_train
271             self.Y_train = Y_train
272             self.accuracy = 0
273
274         def pdf(self, x, h):
275             """
276             Calculates the probability density function for each class
277
278             Parameters
279             -----
280             - x: The set of feature vectors
281             - h: The window size
282
283             Returns
284             -----
285             - p1, p2, p3: The probability density function for each class
286             """
287
288             # separating the training data into the respective classes
289             X1 = self.X_train[:,Y_train == 1]
290             X2 = self.X_train[:,Y_train == 2]
291             X3 = self.X_train[:,Y_train == 3]
292
293             N = len(X1[0]) # number of feature vectors in each class
294             l = len(X1) # numbers of dimensions
295
296             denominator = ((2*np.pi)**(l/2))*h**l # denominator of the probability
density function
297             p1, p2, p3 = np.zeros(len(x[0])), np.zeros(len(x[0])), np.zeros(len(x[0])) #
initializing empty arrays

```

```

298         for i in range(N):
299             exponent1 = np.diag(-((x.T-X1[:,i]) @ (x.T-X1[:,i]).T)/(2*h**2)) #
exponent of the probability density function for class 1
300             p1 += np.exp(exponent1)/(denominator*N) # adding the probability density
function for each feature vector
301
302             exponent2 = np.diag(-((x.T-X2[:,i]) @ (x.T-X2[:,i]).T)/(2*h**2))
303             p2 += np.exp(exponent2)/(denominator*N)
304
305             exponent3 = np.diag(-((x.T-X3[:,i]) @ (x.T-X3[:,i]).T)/(2*h**2))
306             p3 += np.exp(exponent3)/(denominator*N)
307
308         return p1, p2, p3
309
310     def predict(self, X, Y, h, accuracy = True):
311         """
312         Predicts the class for each feature vector in X
313
314         Parameters
315         -----
316         - X: The set of feature vectors
317         - Y: Known labels for the feature vectors
318         - h: The window size
319         - accuracy: If True, the accuracy of the prediction is calculated and stored
in self.accuracy
320
321         Returns
322         -----
323         - pred: The predicted class for each feature vector in X
324         """
325
326         p1, p2, p3 = self.pdf(X, h) # getting the pdfs
327         pred = np.argmax(np.array([p1, p2, p3]), axis = 0) + 1 # predicting the
class for each feature vector
328         if accuracy == True:
329             self.accuracy = np.sum(pred == Y)/len(Y) # calculating the accuracy
330         return pred
331
332     def plot(self, X, Y, h, title, resolution = 50):
333         """
334         Plots the dataset and the decision boundary
335
336         Parameters
337         -----
338         - X: The set of feature vectors
339         - Y: Known labels for the feature vectors
340         - h: The window size
341         - title: The title of the plot (Parzen window density estimator for "title",
h = "h")
342         - resolution: The resolution of the plot
343
344         Returns
345         -----
346         - pred: The predicted class for each feature vector in X
347         """
348
349         x, y = np.linspace(-4.0, 4.2, num=resolution), np.linspace(-3.7, 4.2, num=
resolution) # creating the grid

```



```

350     XY = np.asarray(np.meshgrid(x,y)).reshape(2, -1)
351
352     pred = self.predict(XY, Y, h, accuracy = False) # predicting the class for
each point on the grid
353
354     plt.pcolormesh(x, y, pred.reshape(resolution,resolution), cmap = 'Set1',
alpha=0.6, vmin=0, vmax=12)
355     plt.scatter(X[0][Y == 1], X[1][Y == 1], marker = '+', label = 'Class 1',
color = 'red')
356     plt.scatter(X[0][Y == 2], X[1][Y == 2], marker = '+', label = 'Class 2',
color = 'blue')
357     plt.scatter(X[0][Y == 3], X[1][Y == 3], marker = '+', label = 'Class 3',
color = 'green')
358     plt.xlabel('$x_1$')
359     plt.ylabel('$x_2$')
360     plt.legend()
361     plt.title('Parzen classifier:' + title + ', h = ' + str(h))
362     plt.show()
363
364 def plot_confusionmatrix(self, X, Y, h, title):
365     """
366     Plots the confusion matrix
367
368     Parameters
369     -----
370     - X: The set of feature vectors
371     - Y: Known labels for the feature vectors
372     - h: The window size
373     - title: The title of the plot (Confusion matrix for "title", h = "h")
374     """
375     pred = self.predict(X, Y, h, accuracy = False)
376     confusion_matrix = metrics.confusion_matrix(Y, pred)
377     confusion_matrix_display = metrics.ConfusionMatrixDisplay(confusion_matrix,
display_labels = [1,2,3])
378     confusion_matrix_display.plot(cmap="Reds")
379     plt.title("Parzen classifier: " + title + ", h = " + str(h))
380     plt.grid(False)
381     plt.show()
382
383 if __name__ == "__main__":
384     # task 1
385     print_J3() # printing the J3 scores
386     plot_data() # plotting the training data
387
388     # removing dimension 3 from the training and test sets
389     X_train = np.delete(X_train, 2, 0)
390     X_test = np.delete(X_test, 2, 0)
391
392     # task 2
393     print("-----Task 2-----")
394
395     # Gaussian classifier
396     # plotting the results
397     gaussian = Gaussian(X_train, Y_train)
398     gaussian.plot(X_train, Y_train, 'Training set')
399     gaussian.plot(X_test, Y_test, 'Test set')
400     gaussian.plot_confusionmatrix(X_test, Y_test, 'Test set')
401

```

```

402 # printing the results
403 table = PrettyTable()
404 table.title = "Gaussian classifier"
405 table.field_names = ["Data set", "Accuracy", "Misclassifications"]
406 gaussian.predict(X_train, Y_train, accuracy = True)
407 table.add_row(["Training set", np.round(gaussian.accuracy, 3)*100, np.round(len(
Y_train)*(1-gaussian.accuracy), 3)])
408 gaussian.predict(X_test, Y_test, accuracy = True)
409 table.add_row(["Test set", np.round(gaussian.accuracy, 3)*100, np.round(len(
Y_test)*(1-gaussian.accuracy), 3)])
410 print(table)
411
412 # Parzen classifier
413 h = 0.095 # optimal (found by trial and error)
414 #h = 0.001 # overfitting (extreme case)
415
416 # plotting the results
417 parzen = Parzen(X_train, Y_train)
418 parzen.plot(X_test, Y_test, h, 'Test set')
419 parzen.plot_confusionmatrix(X_test, Y_test, h, 'Test set')
420
421 # printing the results
422 table = PrettyTable()
423 table.title = "Parzen classifier (h = " + str(h) + ")"
424 table.field_names = ["Data set", "Accuracy", "Misclassifications"]
425 parzen.predict(X_train, Y_train, h, accuracy = True)
426 table.add_row(["Training set", np.round(parzen.accuracy, 3)*100, np.round(len(
Y_train)*(1-parzen.accuracy), 3)])
427 parzen.predict(X_test, Y_test, h, accuracy = True)
428 table.add_row(["Test set", np.round(parzen.accuracy, 3)*100, np.round(len(Y_test
)*(1-parzen.accuracy), 3)])
429 print(table)

```

Task 3

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 plt.style.use('ggplot')
4 import scipy as sc
5 from sklearn import metrics
6
7 # Importing the data
8 dataset = sc.io.loadmat('Home-exam\ExamData3D.mat')
9
10 # Splitting the data into training and test sets
11 X_train = dataset['X_train']
12 X_test = dataset['X_test']
13 X_train, X_test = np.delete(X_train, 2, 0), np.delete(X_test, 2, 0)
14 Y_train = dataset['Y_train'][0]
15 Y_test = dataset['Y_test'][0]
16
17 # adding the extra 'dimension' to X_train and X_test
18 X_train = np.insert(X_train, 0, 1, axis=0)
19 X_test = np.insert(X_test, 0, 1, axis=0)
20
21 class perceptron:
22     def __init__(self, c, y, N=len(X_train[0])):
23         """

```

```

24     Perceptron
25     -----
26
27     Parameters
28     -----
29     - c: class which the perceptron is trained to recognize
30     - y: training set
31     - N: number of training feature vectors
32     """
33
34     self.w = np.random.rand(3,1) # random initial weights
35     self.y_hat = np.zeros(N)
36     self.c = c
37     self.y = np.zeros(N)
38     self.y[y == self.c] = 1 # y = 1 if the feature vector belongs to class c, 0
39     else
40
41     def forward(self, x):
42         """
43         Forward pass (calculates the output of the perceptron)
44
45         Parameters
46         -----
47         - x: array of feature vectors
48         """
49
50         g = self.w.T @ x # discriminant function
51         f = 1/(1+np.exp(-g)) # activation function
52         self.y_hat = f
53
54     def backward(self, x):
55         """
56         Backward pass (updates the weights of the perceptron)
57
58         Parameters
59         -----
60         - x: array of feature vectors
61         """
62
63         df = 1/(1+np.exp(-self.w.T @ x))*(1-(1/(1+np.exp(-self.w.T @ x)))) #
64         derivative of the activation function
65         dw = (rho*(self.y-self.y_hat)*df @ x.T).T # weight update
66         self.w += dw
67
68     class neural_network:
69         def __init__(self, X_train, Y_train, X_test, Y_test):
70             """
71             Neural network (single layer with 3 nodes)
72
73             Parameters
74             -----
75             - X_train: training set
76             - Y_train: training labels
77             - X_test: test set
78             - Y_test: test labels
79             """
80
81             self.accuracy = 0

```

```

80     self.predictions = np.zeros(len(Y_test))
81     self.X_train = X_train
82     self.Y_train = Y_train
83     self.X_test = X_test
84     self.Y_test = Y_test
85
86     # setting up the perceptrons
87     self.neurons = [perceptron(1, self.Y_train), perceptron(2, self.Y_train),
88                     perceptron(3, self.Y_train)]
89
90 def train(self, T):
91     """
92     Training the neural network (updates the weights T times)
93
94     Parameters
95     ----
96     - T: number of epochs
97     """
98     for t in range(T):
99         for neuron in self.neurons:
100             neuron.forward(self.X_train)
101             neuron.backward(self.X_train)
102
103 def predict(self):
104     """
105     Predicting the labels of the test set
106
107     Returns
108     ----
109     - accuracy: accuracy of the predictions
110     """
111
112     for neuron in self.neurons:
113         neuron.forward(self.X_test) # forward pass to get the predictions
114
115     # taking the most probable class
116     self.predictions = np.argmax([self.neurons[0].y_hat, self.neurons[1].y_hat,
117 self.neurons[2].y_hat], axis=0)+1
118     self.accuracy = np.sum(self.predictions == self.Y_test)/len(self.Y_test)
119
120     return self.accuracy
121
122 def decision_boundary(self, txt, resolution=120):
123     """
124     Plotting the decision boundary
125
126     Parameters
127     ----
128     - txt: title of the plot
129     - resolution: resolution of the decision boundary
130     """
131
132     x, y = np.linspace(-4.0, 4.2, num=resolution), np.linspace(-4.0, 4.2, num=
133 resolution)
134     XY = np.asarray(np.meshgrid(x,y)).reshape(2, -1)
135     XY = np.insert(XY, 0, 1, axis=0)
136     for neuron in self.neurons:

```

```

135         neuron.forward(XY)
136         predictions = np.argmax([self.neurons[0].y_hat, self.neurons[1].y_hat, self.
neurons[2].y_hat], axis=0)+1
137         plt.pcolormesh(x, y, predictions.reshape(resolution, resolution), cmap = '
Set1', alpha=0.6, vmin=0, vmax=12)
138         plt.scatter(self.X_test[1, self.Y_train == 1], self.X_test[2, self.Y_train
== 1], marker = '+', color='red', label='Class 1')
139         plt.scatter(self.X_test[1, self.Y_train == 2], self.X_test[2, self.Y_train
== 2], marker = '+', color='blue', label='Class 2')
140         plt.scatter(self.X_test[1, self.Y_train == 3], self.X_test[2, self.Y_train
== 3], marker = '+', color='green', label='Class 3')
141         plt.title("Single layer Perceptron: " + txt)
142         plt.xlabel("$x_1$")
143         plt.ylabel("$x_2$")
144
145     def plot_confusion_matrix(self, title):
146         """
147         Plotting the confusion matrix
148
149         Parameters
150         ----
151         - title: title of the plot
152         """
153
154         confusion_matrix = metrics.confusion_matrix(self.Y_test, self.predictions
[0])
155         confusion_matrix_display = metrics.ConfusionMatrixDisplay(confusion_matrix)
156         confusion_matrix_display.plot(cmap="Reds")
157         plt.title("Single layer Perceptron: " + title)
158         plt.xlabel("Predicted label")
159         plt.ylabel("True label")
160         confusion_matrix_display.ax_.set_xticks([])
161         confusion_matrix_display.ax_.set_yticks([])
162         plt.show()
163
164     if __name__ == "__main__":
165         """---Testing the neural network---"""
166         # defining parameters
167         rho = 0.01 # learning rate
168         T = np.linspace(0, 200, 201, dtype=int) # epochs
169
170         # initiating the neural network, for both training and test sets
171         NN_test = neural_network(X_train, Y_train, X_test, Y_test)
172         NN_train = neural_network(X_train, Y_train, X_train, Y_train)
173
174         # get the accuracy as a function of T
175         test_accuracies = np.zeros(len(T))
176         train_accuracies = np.zeros(len(T))
177         for t in range(len(T)):
178             #print(t)
179             NN_test.train(T[t])
180             NN_train.train(T[t])
181             test_accuracies[t] = NN_test.predict()
182             train_accuracies[t] = NN_train.predict()
183
184         # plotting the decision boundary
185         plt.subplot(1,2,1)
186         NN_test.decision_boundary("Test set" + " (" + str(chr(961)) + "=" + str(rho)+ ")

```

```

187     " + ", T = " + str(T[-1]))
188     plt.subplot(1,2,2)
189     NN_train.decision_boundary("Training set")
190     plt.show()
191     print("Training set: ", np.round(train_accuracies[-1], 3), ", Test set: ", np.
192         round(test_accuracies[-1], 3))
193
194     # plotting the accuracy as a function of T
195     plt.plot(T, test_accuracies, label="Test set")
196     plt.plot(T, train_accuracies, label="Training set")
197     plt.title("Accuracy as a function of T (" + str(chr(961)) + "=" + str(rho) + ")")
198     plt.xlabel("T")
199     plt.ylabel("Accuracy")
200     plt.legend()
201     plt.show()
202     print("Highest accuracy for test set: ", np.round(np.max(test_accuracies), 3), "
203         , at T = ", np.argmax(test_accuracies)+1)
204
205     # plotting the confusion matrix
206     NN_test.plot_confusion_matrix("Test set")
207     NN_train.plot_confusion_matrix("Training set")

```

Task 4

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3  plt.style.use('ggplot')
4  import scipy as sc
5  from sklearn import metrics
6
7  # Importing the data
8  dataset = sc.io.loadmat('Home-exam\ExamData3D.mat')
9
10 # Splitting the data into training and test sets
11 X_train = dataset['X_train']
12 X_test = dataset['X_test']
13 X_train, X_test = np.delete(X_train, 2, 0), np.delete(X_test, 2, 0)
14 Y_train = dataset['Y_train'][0]
15 Y_test = dataset['Y_test'][0]
16
17 # adding the extra 'dimension' to X_train
18 X_train = np.insert(X_train, 0, 1, axis=0)
19 X_test = np.insert(X_test, 0, 1, axis=0)
20
21 class perceptron:
22     def __init__(self, dimension, N=len(X_train[0])):
23         """
24         Perceptron
25         -----
26         Parameters
27         -----
28         - dimension: dimension of the feature vectors
29         - N: number of training feature vectors
30         """
31         # initializing random weights, and all other attributes to zero
32         self.w = np.random.rand(dimension,1)
33         self.dw = np.zeros(self.w.shape)
34         self.y_hat = np.zeros(N)

```

```

35     self.v = np.zeros(N)
36     self.e = np.zeros(N)
37     self.delta = 0
38
39     def forward(self, x):
40         """
41         Forward pass (calculates the output of the perceptron)
42
43         Parameters
44         -----
45         - x: array of feature vectors
46         """
47         self.v = self.w.T @ x
48         self.y_hat = 1/(1+np.exp(-self.v))
49
50     def calculate_e(self, deltas, weights):
51         """
52         Calculate the error of the perceptron
53
54         Parameters
55         -----
56         - deltas: the deltas of the next layer
57         - weights: the weights of the next layer (that are connected to the
58         perceptron)
59         """
60         self.e = (deltas @ weights.T).T
61
62     def calculate_delta(self):
63         """
64         Calculates the value for delta
65         """
66         f = 1/(1+np.exp(-self.v))
67         df = f*(1-f)
68         self.delta = np.multiply(self.e,df)
69
70     def backward(self, y):
71         """
72         Backward pass (updates the weights of the perceptron)
73
74         Parameters
75         -----
76         - y: array of outputs from the previous layer
77         """
78         self.dw = alpha*self.dw - (rho*y*self.delta.T) # (momentum term + cost
79         function term)
80         self.w += self.dw
81
82 class output_perceptron(perceptron):
83     def __init__(self, label, dimension, y, N=len(X_train[0])):
84         """
85         Output perceptron (only difference from regular perceptron if calculating e
86         and delta)
87
88         Parameters
89         -----
90         - label: the label of the output perceptron
91         - dimension: dimension of the feature vectors
92         - y: training set

```

```

90         - N: number of training feature vectors
91         """
92         super().__init__(dimension, N=len(X_train[0]))
93         self.y = np.zeros(N)
94         self.y[y == label] = 1
95
96     def calculate_e(self):
97         """
98         Calculates the error
99         """
100         self.e = self.y_hat - self.y
101
102     def calculate_delta(self):
103         """
104         Calculates the value for delta
105         """
106         f = 1/(1+np.exp(-self.v))
107         df = f*(1-f)
108         self.delta = self.e*df
109
110 class MLP:
111     def __init__(self, shape, X_train, Y_train):
112         """
113         Multi-layer perceptron
114         -----
115         Parameters
116         -----
117         - shape: the shape of the hidden layers of the network, input layer is 2
118         since the dataset is 2d and output layer is 3 nodes from the task description
119         - X_train: training set
120         - Y_train: training labels
121         """
122         # initializing attributes
123         self.shape = shape
124         self.X_train = X_train
125         self.Y_train = Y_train
126         self.layers = []
127         layer = []
128
129         # creating the layers
130         # first hidden layer
131         for i in range(shape[0]):
132             layer.append(perceptron(3))
133         self.layers.append(layer)
134
135         # other hidden layers
136         for i in range(1, len(shape)):
137             layer = []
138             for j in range(shape[i]):
139                 layer.append(perceptron(shape[i-1]+1))
140             self.layers.append(layer)
141
142         # output layer
143         output_layer = [output_perceptron(1, shape[-1]+1, self.Y_train),
144             output_perceptron(2, shape[-1]+1, self.Y_train), output_perceptron(3, shape
145             [-1]+1, self.Y_train)]
146         self.layers.append(output_layer)

```



```

145
146 def forward_propagation(self, X_train):
147     """
148     Propogates the network forwards (calculates the output of the network)
149
150     Parameters
151     -----
152     - X_train: training set (input to the network)
153     """
154
155     # inpput the training set into the first layer
156     for neuron in self.layers[0]:
157         neuron.forward(X_train)
158
159     # propagate the output of the first layer to the next layer, and so on
160     for i in range(1, len(self.layers)):
161         input = []
162
163         for n in self.layers[i-1]: # for all neurons in the previous layer
164             input.append(n.y_hat.tolist()[0]) # add the output to the input list
165         input = np.insert(input, 0, 1, axis=0) # add the bias term
166
167         for neuron in self.layers[i]: # propogate the input to the next layer
168             neuron.forward(np.asarray(input))
169
170 def backward_propagation(self):
171     """
172     Propogates the network backwards (updates the weights of the network)
173     """
174
175     # calculate e and delta for the last layer
176     for neuron in self.layers[-1]:
177         neuron.calculate_e()
178         neuron.calculate_delta()
179
180     for i in reversed(range(0, len(self.layers)-1)): # iterate backwards through
the layers
181         for j in range(len(self.layers[i])): # iterate through the neurons in
the layer
182             deltas = []
183             weights = []
184             for neuron in self.layers[i+1]: # for all neurons in the next layer
185
186                 deltas.append(neuron.delta) # add the delta to the list
187                 weights.append(neuron.w[j+1]) # add the weights to the list
188             deltas, weights = np.matrix(np.asarray(deltas)).T, np.matrix(np.
asarray(weights)).T # convert to matrices
189
190             self.layers[i][j].calculate_e(deltas, weights) # calculate the new
error for the neuron
191             self.layers[i][j].calculate_delta() # calculate the new delta for
the neuron
192
193     # update the weights for the first layer
194     for neuron in self.layers[0]:
195         neuron.backward(self.X_train)
196
197     # update the weights for the other layers

```

```

198         for i in range(1, len(self.layers)):
199             for j in range(len(self.layers[i])):
200                 y = []
201                 for neuron in self.layers[i-1]:
202                     y.append(neuron.y_hat)
203                 y = np.asarray(y)
204                 y = np.insert(y, 0, 1, axis=0) # bias term
205                 self.layers[i][j].backward(np.matrix(y))
206
207     def train(self, T, X_train, Y_train, X_test, Y_test, predict=False):
208         """
209         Training the network (updates the weights T times)
210
211         Parameters
212         ----
213         - T: number of epochs
214         - X_train: training set
215         - Y_train: training labels
216         - X_test: test set
217         - Y_test: test labels
218         - predict: (bool) if we want to get the accuracy for each epoch
219         """
220
221         # if we want to get the accuracy for each epoch
222         if predict == True:
223             predictions_train = np.zeros(T)
224             predictions_test = np.zeros(T)
225             for t in range(T):
226                 # Just for checking the progress
227                 if t % 500 == 0:
228                     print(t)
229                 self.forward_propagation(self.X_train)
230                 self.backward_propagation()
231                 self.predict(X_train, Y_train)
232                 predictions_train[t] = self.accuracy
233                 self.predict(X_test, Y_test)
234                 predictions_test[t] = self.accuracy
235             return predictions_train, predictions_test
236
237         # if we only want the final accuracy
238         else:
239             for t in range(T):
240                 self.forward_propagation(self.X_train)
241                 self.backward_propagation()
242
243     def predict(self, X_test, Y_test):
244         """
245         Predicting the labels of the test set
246
247         Parameters
248         ----
249         - X_test: test set
250         - Y_test: test labels
251         """
252         # calculate the output of the network
253         self.forward_propagation(X_test)
254
255         # taking the most probable class

```

```

256         self.prediction = np.argmax([self.layers[-1][0].y_hat, self.layers[-1][1].
y_hat, self.layers[-1][2].y_hat], axis=0)+1
257
258         # calculating the accuracy
259         self.accuracy = np.sum(self.prediction == Y_test)/len(Y_test)
260
261     def decision_boundary(self, Y_test, resolution=100):
262         """
263         Plotting the decision boundary
264
265         Parameters
266         -----
267         - Y_test: The true labels of the test set
268         - resolution: The resolution of the decision boundary
269         """
270
271         # creating the meshgrid
272         x, y = np.linspace(-4.0, 4.2, num=resolution), np.linspace(-4.0, 4.2, num=
resolution)
273         XY = np.asarray(np.meshgrid(x,y)).reshape(2, -1)
274         XY = np.insert(XY, 0, 1, axis=0)
275
276         # calculating the output of the network
277         self.forward_propagation(XY)
278
279         # predicting the labels
280         predictions = np.argmax([self.layers[-1][0].y_hat, self.layers[-1][1].y_hat,
self.layers[-1][2].y_hat], axis=0)+1
281
282         # plotting the decision boundary
283         plt.pcolormesh(x, y, predictions.reshape(resolution, resolution), cmap = '
Set1', alpha=0.6, vmin=0, vmax=12)
284         plt.scatter(X_test[1, Y_test == 1], X_test[2, Y_test == 1], marker = '+',
color='red', label='Class 1')
285         plt.scatter(X_test[1, Y_test == 2], X_test[2, Y_test == 2], marker = '+',
color='blue', label='Class 2')
286         plt.scatter(X_test[1, Y_test == 3], X_test[2, Y_test == 3], marker = '+',
color='green', label='Class 3')
287         plt.xlabel("$x_1$")
288         plt.ylabel("$x_2$")
289
290     def plot_confusion_matrix(self, Y_test, title):
291         """
292         Plotting the confusion matrix for the prediction
293
294         Parameters
295         -----
296         - Y_test: The true labels of the test set
297         - title: The title of the plot
298         """
299
300         confusion_matrix = metrics.confusion_matrix(Y_test, self.prediction[0])
301         confusion_matrix_display = metrics.ConfusionMatrixDisplay(confusion_matrix)
302         confusion_matrix_display.plot(cmap="Reds")
303         plt.title("Confusion matrix for " + title)
304         confusion_matrix_display.ax_.set_xticks([])
305         confusion_matrix_display.ax_.set_yticks([])
306         plt.show()

```

```

307
308 if __name__ == "__main__":
309     # defining the parameters
310     rho = 0.008
311     alpha = 0.65
312     shape = [10, 8, 6]
313     T = 2500
314
315     # creating and training the network
316     nn = MLP(shape, X_train, Y_train)
317     accuracies_train, accuracies_test = nn.train(T, X_train, Y_train, X_test, Y_test
, predict=True)
318
319     # plotting the decision boundary
320     print("Test set: ", np.round(accuracies_test[-1], 3), ", Training set: ", np.
round(accuracies_train[-1], 3))
321     plt.title(chr(961) + "=" + str(rho) + " " + chr(945) + "=" + str(alpha))
322     nn.decision_boundary(Y_test, resolution = 200)
323     plt.show()
324
325     # plotting the accuracy as a function of T
326     plt.plot(np.arange(T), accuracies_train, label="Training set", color = 'blue')
327     plt.plot(np.arange(T), accuracies_test, label="Test set", color = 'red')
328     plt.xscale("log")
329     plt.title("Accuracy as a function of T")
330     plt.xlabel("T")
331     plt.ylabel("Accuracy")
332     plt.legend()
333     plt.show()
334
335     # plotting the confusion matrix
336     nn.plot_confusion_matrix(Y_test, " Test set")

```