

INTRODUCTION MACHINE LEARNING

EXERCISE 4

The logo of Bauhaus-Universität Weimar, featuring the university's name in white sans-serif font on a solid red rectangular background.

Bauhaus-
Universität
Weimar

TEACHER:

Johannes Kiesel

GROUP:

Group 16

SUBMITTED BY:

Aaron Perez Herrera
Cesar Fernando Gamba Tiusaba
Chun Ting Lin
Olubunmi Emmanuel Ogunleye

Exercise 1: Gradient Descent (1+1+1=3 Points)

(a) Name one difference between the perceptron training rule and the gradient descent method.

The Perceptron training rule is not based on residuals in the $(p+1)$ dimensional input- output-space) but refers to the input space only, where it simply evaluates the side of the hyperplane as a binary feature (correct side or not). Gradient descent is a regression approach that exploits the residuals provided by a loss function of choice, whose differential is evaluated to guide hyperplane search

(b) What are the main requirements for the hypothesis space and error function for gradient descent to be successfully applied?

In batch gradient descent, the weight vector is determined by calculating the differences in weights for each example in the training data and subsequently computing the weight vector for the entire batch in a single step. Incremental gradient descent involves computing the weight vector immediately after calculating the difference in weights for each example in the training dataset, leading to a more iterative and continuous update process. Compared to batch gradient descent, the example-based weight adaptation of incremental gradient descent can better avoid getting stuck in a local minimum of the loss function.

(c) Name the key difference between the algorithm of batch gradient descent (BGD) and the algorithm of incremental gradient descent (IGD).

The perceptron $y_1()$, with a non-zero bias term w_0 , is more general than $y_0()$, which has a bias term of zero. The bias term acts as an additional flexibility factor, allowing $y_1()$ to represent all the functions $y_0()$'s weights can capture, plus additional functions that arise from the non-zero bias. In essence, $y_1()$ can model a broader range of relationships in the data, making it more versatile and general than $y_0()$. The bias term allows models to fit the data better by providing additional flexibility.

Exercise 2 : Perceptron Learning (2 Points)

Consider two perceptrons, $y_0()$ and $y_1()$, both defined by the function $\text{heaviside}(\sum w_j x_j)$. Both perceptrons have identical weights except for the bias term: $w_0 = 0$ for $y_0()$ and $w_0 = 1$ for $y_1()$. Determine if one of the perceptrons, $y_0()$ or $y_1()$, is more general than the other (as defined in the lecture units on concept learning). If one is more general, specify which one and explain your answer.

We have two perceptrons, $y_0()$ and $y_1()$, which are defined using the Heaviside function (a step function) and both have identical weights except for the bias term. Specifically.

- $w_0 = 0$ for $y_0()$
- $w_0 = 1$ for $y_1()$

And each perceptron is represented as:

$$y_0 = \text{heaviside} \left(\sum_{j=0}^p w_j x_j \right) \text{ with } w_0 = 0$$

$$y_1 = \text{heaviside} \left(\sum_{j=0}^p w_j x_j \right) \text{ with } w_0 = 1$$

The bias between $y_0()$ and $y_1()$ is 0 and 1, respectively. This bias shifts the decision boundary of the perceptron. The general decision rule for a perceptron is:

$$\text{Output} = \text{heaviside} \left(\sum_{j=0}^p w_j x_j + w_0 \right)$$

For perceptron $y_0()$, the output is determined by the sum of the weighted inputs, and since the bias is 0, the decision boundary is where:

$$\sum_{j=0}^p w_j x_j = 0$$

For perceptron $y_1()$, the decision boundary is where:

$$\sum_{j=0}^p w_j x_j + 1 = 0 \text{ or } \sum_{j=0}^p w_j x_j = -1$$

- The perceptron $y_0()$ has the decision boundary at $\sum_{j=0}^p w_j x_j = 0$
- The perceptron $y_1()$ has the decision boundary shifted by -1, i.e., at $\sum_{j=0}^p w_j x_j = -1$

Since the bias in $y_0()$ is 0, it has a decision boundary at the origin. On the other hand, the bias in $y_1()$ shifts the decision boundary. This shift in the decision boundary means that $y_1()$ is more general than $y_0()$ because it allows for a broader range of inputs to be classified as positive, due to the shift in the threshold.

In conclusion $y_1()$ is more general than $y_0()$ because the bias in $y_1()$ shifts the decision boundary, allowing for a greater variety of inputs to be classified as 1. Therefore, $y_1()$ can classify a broader range of inputs, making it more general.

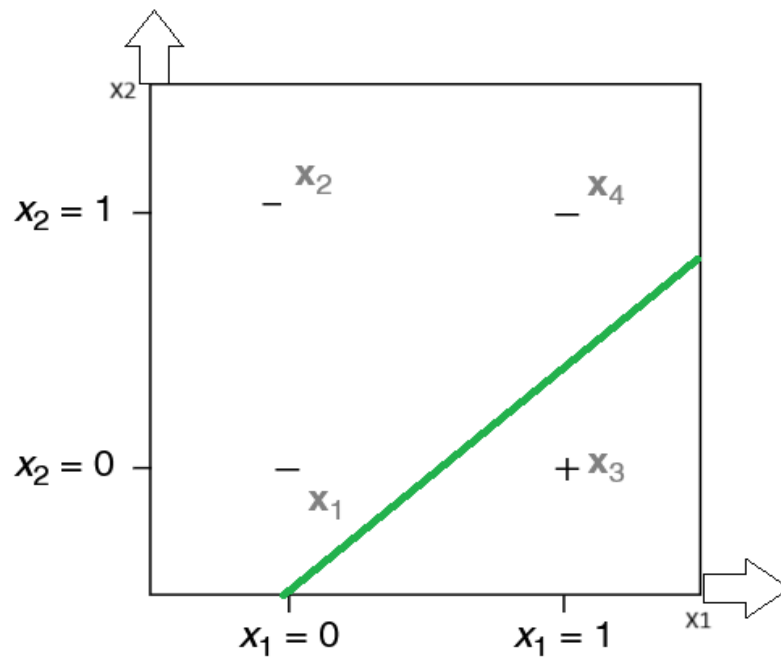
Exercise 3 : Perceptron Learning (1+1+2+2+1=6 Points)

In this exercise, you design a single perceptron with two inputs x_1 and x_2 . This perceptron shall implement the Boolean formula $A \wedge \neg B$ with a suitable function $y(x_1, x_2)$. Use the values 0 for false and 1 for true.

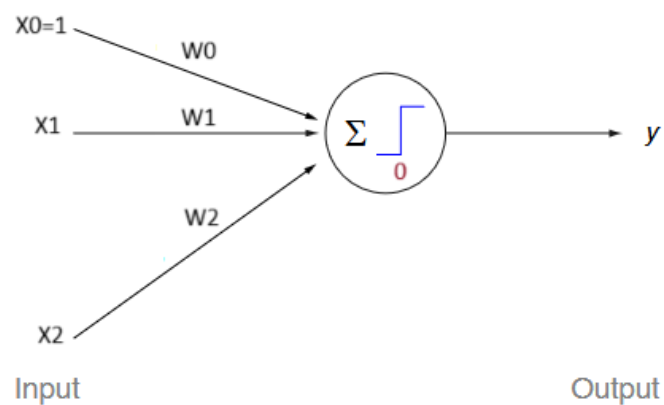
(a) Draw all possible examples and a suitable decision boundary in a coordinate system.

First, we build the table

Example	A	B	$\neg B$	$A \wedge \neg B$	C
X_1	0	0	1	0	-
X_2	0	1	0	0	-
X_3	1	0	1	1	+
X_4	1	1	0	0	-



(b) Draw the graph of the perceptron. The schematic must include x_1 , x_2 , and all model weights.



(c) Manually determine the weights $w = (w_0, w_1, w_2)$ for the decision boundary you drew in (a).

A suitable weight values would be by introducing a negative bias term w_0

$$w_0 = -0.5; \quad w_1 = 1; \quad w_2 = -1$$

$$y = w_0 + x_1 w_1 + x_2 w_2$$

Considering we are using Heaviside function.

Example	X1	X2	y	C
X ₁	0	0	-0.5 < 0	-
X ₂	0	1	-1.5 < 0	-
X ₃	1	0	0.5 > 0	+
X ₄	1	1	-0.5 < 0	-

The results reflect the behavior of the perceptron.

- (d) Now determine w using the perceptron training algorithm (PT). Use a learning rate n of 0.3 and initialize the weights with $w_0 = -0.5$ and $w_1 = w_2 = 0.5$. Instead of selecting examples randomly, use the following examples in the given order (stop after those four examples):

Start with iteration of following examples and check the updated value after each example

X1	X2	C
0	0	0
0	1	0
1	0	1
1	1	0

$$w = [-0.5, 0.5, 0.5] \rightarrow n = 0.3$$

Example	x1	x2	c	y	ΔW_0	ΔW_1	ΔW_2		W0	W1	W2
X1	0	0	0	0	0	0	0		-0.5	0.5	0.5
X2	0	1	0	1	-0.3	0	-0.3		-0.8	0.5	0.2
X3	1	0	1	0	0.3	0.3	0		-0.5	0.8	0.2
X4	1	1	0	1	-0.3	-0.3	-0.3		-0.8	0.5	-0.1

$$y = \text{heaviside}(w_0 + w_1 * x_1 + w_2 * x_2)$$

$$\Delta w = n(c - y) * x \rightarrow x = [x_0, x_1, x_2]$$

$$w = w + \Delta w$$

Final weights

$$W_0 = -0.8$$

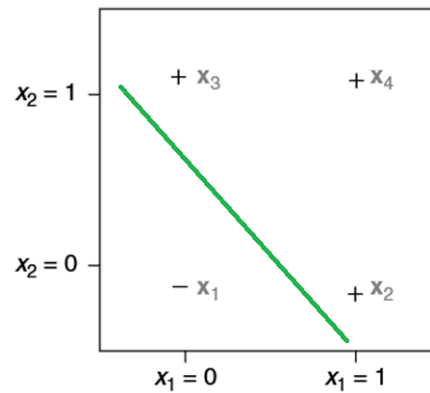
$$W_1 = 0.5$$

$$W_2 = -0.1$$

Draw the decision boundary after every weight update into the coordinate system of (a).

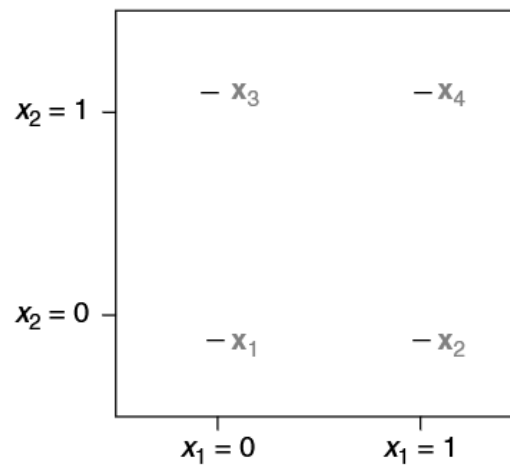
First calculation

Example	x1	x2	c	y
X1	0	0	0	0
X2	0	1	0	1
X3	1	0	1	1
X4	1	1	0	1



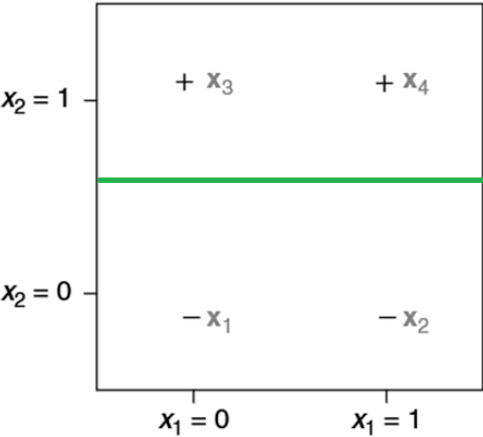
First update

Example	x1	x2	c	y
X1	0	0	0	0
X2	0	1	0	0
X3	1	0	1	0
X4	1	1	0	0



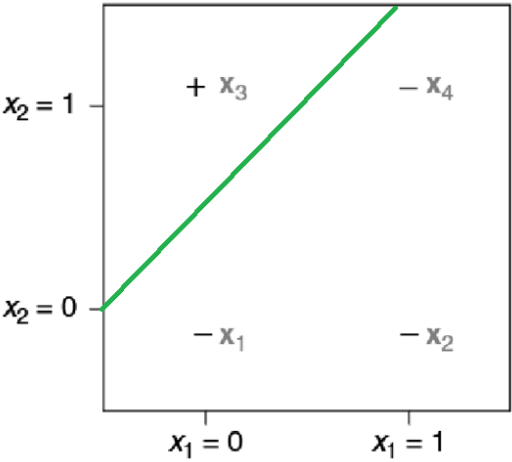
Second update

Example	x1	x2	c	y
X1	0	0	0	0
X2	0	1	0	0
X3	1	0	1	1
X4	1	1	0	1



Third / last update

Example	x1	x2	c	y
X1	0	0	0	0
X2	0	1	0	0
X3	1	0	1	1
X4	1	1	0	0



(e) Briefly describe one effect of changing the learning rate η on the learning progress.

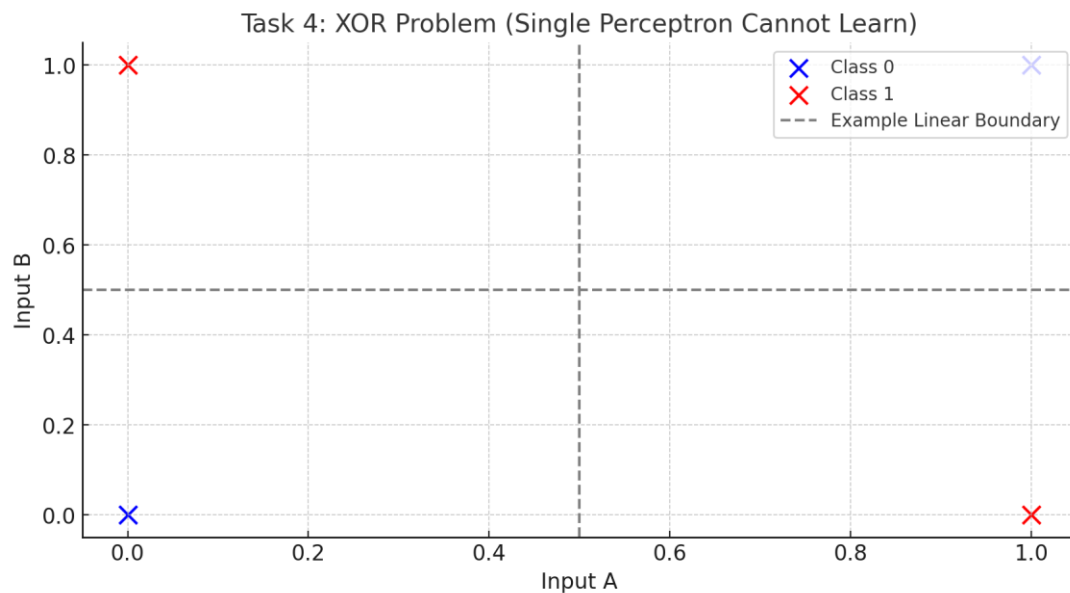
A large learning rate means we are making larger changes to the weights every time we do a learning update. However, if the changes are too large, we might oscillate back and forth rather than focusing on a solution.

- Higher η : Faster convergence but risks overshooting optimal weights.
- Lower η : Slower convergence but smoother and more stable learning.

Exercise 4: Perceptron Learning why a single perceptron cannot learn the Boolean XOR function A XOR B?

XOR function outputs 1 if exactly one of the inputs is 1, otherwise 0.

A	B	A XOR B
0	0	0
0	1	1
1	0	1
1	1	0



Points (0, 0) and (1, 1) represent one class (e.g., blue, 0), while (0, 1) and (1, 0) represent the other class (e.g., red, 1).

Notice that the data is not linearly separable. A single straight line cannot divide the output classes into distinct groups.

No single straight line can separate the two classes (0 and 1) correctly. This is because XOR is not linearly separable.

Exercise 5 : Multilayer Perceptron

Consider the minimum multilayer perceptron that can handle the XOR problem (slide ML:IV-59), but with an activation function $f(z) = (f(z_1), \dots, f(z_d))^T$ instead of Heaviside(). This is a two-layer perceptron with $p = 2$ attributes, a hidden layer with $l = 2$ units, and an output layer with $k = 1$ units.

Therefore, as per slide ML:IV-82:

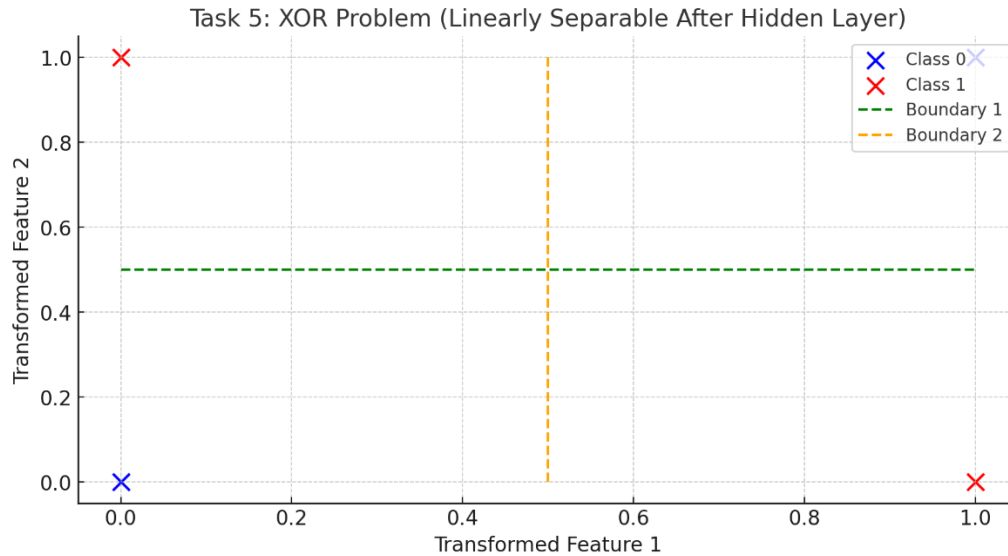
$$y(x) = f \left(W^o \begin{pmatrix} 1 \\ f(W^h x) \end{pmatrix} \right) = f \left((w_{10}^o, w_{11}^o, w_{12}^o) \left(\begin{pmatrix} 1 \\ f \left(\begin{pmatrix} w_{10}^h & w_{11}^h & w_{12}^h \\ w_{20}^h & w_{21}^h & w_{22}^h \end{pmatrix} \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} \end{pmatrix} \right) \right) \right)$$

Confirm for yourself that thus

$$y(x) = f(w_{10}^o + w_{11}^o \cdot f(w_{10}^h + w_{11}^h \cdot x_1 + w_{12}^h \cdot x_2) + w_{12}^o \cdot f(w_{20}^h + w_{21}^h \cdot x_1 + w_{22}^h \cdot x_2))$$

1. Understand the architecture:
XOR can be solved by a 2-layer perceptron:
 - Input layer: Two neurons for A and B.
 - Hidden layer: Two neurons to map the problem to a linearly separable representation.
 - Output layer: One neuron to combine results from the hidden layer.
2. Activation function:
The problem states the activation function is the identity function $f(z) = z$.
3. Key steps:
Write the equations for the hidden layer:
$$h_1 = f(w_{10} + w_{11}A + w_{12}B)$$
$$h_2 = f(w_{20} + w_{21}A + w_{22}B)$$

Write the equation for the output layer:
$$y = f(w_{o0} + w_{o1}h_1 + w_{o2}h_2)$$
4. Proof for failure with $f(z) = z$:
 - Show that with the identity activation, the perceptron behaves linearly.
 - XOR still requires non-linear boundaries, so the network fails.



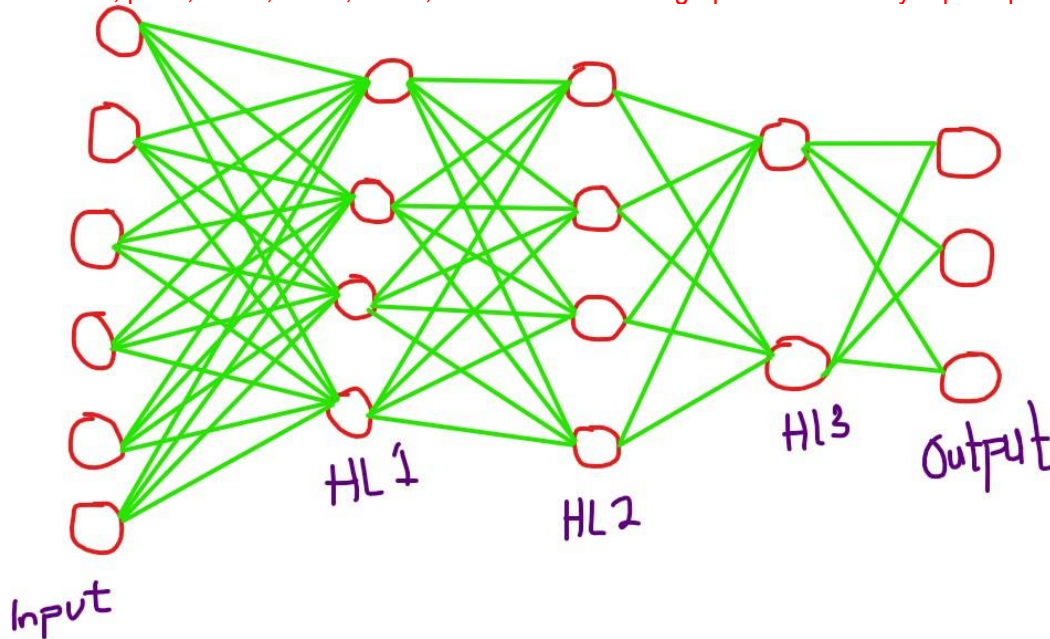
Graph Transformation:

The hidden layer transforms the input space so that the points become separable in a new space. The final output layer combines these to classify XOR correctly.

Exercise 6 : Parameters of the Multilayer Perceptrons (2+1+0+0=3 Points)

In this exercise, you analyze the number of weights (parameters) of multilayer perceptrons. We use the notation from the lecture (e.g., slide ML:IV-104), where multilayer perceptrons have d layers, p attributes, hidden layer i with l_i units, and an output layer with k units.

- a. Let $d = 4$, $p = 6$, $l_1 = 4$, $l_2 = 4$, $l_3 = 2$, and $k = 3$. Draw the graph of the multilayer perceptron.



- b. Calculate the number of weights in the multilayer perceptron of (a).

- Input – HL1 = $6 \times 4 = 24$
- HL1 – HL2 = $4 \times 4 = 16$

- HL 2 – HL 3 = $4 \times 2 = 8$
- HL3 – Output = $3 \times 2 = 6$

Total = $24 + 16 + 8 + 6 = 54$

- c. Calculate the number of weights in the multilayer perceptron of (a) but with each l_i doubled, i.e., $l_1 = 8$, $l_2 = 8$, $l_3 = 4$. Has the number of weights doubled as well?

- Input – HL1 = $6 \times 8 = 48$
- HL 1 – HL 2 = $8 \times 8 = 64$
- HL 2 – HL 3 = $8 \times 4 = 32$
- HL3 – Output = $3 \times 4 = 12$

Total = $48 + 64 + 32 + 12 = 156$

The number of weights has not doubled because if the number of neurons in each layer is increased, this is a multiplication relationship that exists between the layers meaning that the weights exceed double.

- d. Let $f(p, l_1, \dots, l_{d-1}, k)$ be a function that computes the number of weights in the general case. Write down an expression for f .

- Weights of the First layer: $p \times l_1$
- Weights of Hidden Layers: $l_1 \times l_2 + l_2 \times l_3 + \dots + l_{d-2} \times l_{d-1}$
 $= \sum_{i=1}^{d-2} l_i \times l_{i+1}$
- Weights of Output Layers: $l_{d-1} \times k$

Hence, the resulting expression is written as

$$f = p \times l_1 + \sum_{i=1}^{d-2} l_i \times l_{i+1} + l_{d-1} \times k$$

Exercise 7: P Classification with Neural Networks (1+1+1+1+1+1=6 Points)

In this exercise, you will implement a two-layer perceptron for predicting whether a given text was written by a human or generated by a language model. Download and use these files from Moodle (the tsv files are the same as in the last sheet):

- features-train.tsv: Feature vectors for each example in the training set.
- features-test.tsv: Feature vectors for each example in the test set.

- labels-train.tsv: Labels for each example in the training set indicating the class is_human ($C = \{\text{True}, \text{False}\}$)
- programming_exercise_neural_networks.py: Template containing function stubs for the exercise. Use it with the files above as follows:
 - python3 programming_exercise_neural_networks.py
 - features-train.tsv labels-train.tsv features-test.tsv
 - predictions-test.tsv
 - requirements.txt: Requirements file for the template; can be used to install dependencies.
- a. Implement a function `encode_class_values` to encode class values as vectors $c \in \{0, 1\}^k$ (see slide ML:IV-100).
- b. Implement a function `predict_probabilities` that, given W_h and W_o , predicts the class probabilities for each example of a dataset.¹
- c. Implement a function `predict` that, given W_h and W_o , predicts the class for each example of a dataset (as the one with the highest probability).
- d. Implement a function `train_multilayer_perceptron` that fits a multilayer perceptron using the IGD Algorithm (slide ML:IV-100). Like in the last exercise sheet, make sure to return the misclassification rate on both training and validation set for plotting, but this time also return the weights after each iteration (for task (e)).¹
- e. Select the model (iteration) that achieved the best misclassification rate on the validation set and use it to predict the label for each example in the test set (features-test.tsv). Write the predictions as one column to a file `predictions-test.tsv` and submit that along with your other solutions.
- f. Which lines of your code would you need to change from two classes to all three? Mark each of them with the comment `"# change for 3 classes"`.

Check for the answers to this question in the folder “Final Program” you will find in it all necessary files -> Like the code and “predictions-test.tsv”