

Cognitive Algorithms Test Exam

WS 2023/24

Please fill in below your full name, your matriculation number and field of studies.

I hereby confirm that I feel capable to participate in this exam.

Signature

Name	
Field of study	
Matriculation number	
Registration	<input type="checkbox"/> Via Moses <input type="checkbox"/> Other

Section	Points	Score
Overview questions	14	
Linear Classification	25	
Kernel methods and Kernel Ridge Regression	7	
Unsupervised Learning	13	
MLP	14	
Cross-Validation	8	
Total	81	

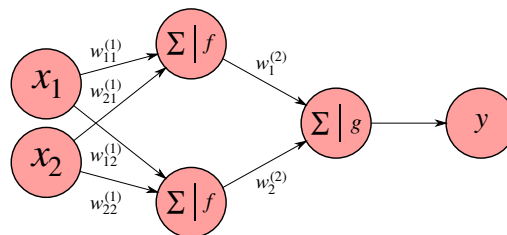
1. Overview questions [14 Points total]

1.1 [1 Point] What do we mean when we say that our machine learning model generalizes well? (1-2 sentences)

test error

Solution: A model generalizes well when its error on new/unseen data is small.

1.2 [2 Points] The neural network displayed below is _____ training algorithm with ____ hidden layer(s). Choose one Answer from the first column to fill in the first gap and one answer from the second column to fill in the second gap.



☒ a supervised

☐ an unsupervised

☒ 1

☐ 2

1.3 [1 Point] The neural network as displayed in the question above is a linear method.

Hint! When we talk about linear methods (including linear regression), we mean methods that are linear in \mathbf{w} , i.e. $f(x) = \mathbf{w}^\top x$ where \mathbf{w} can also contain a non-linear transformation of the data and the offset/bias β .

☐ True

☐ False

☒ depends on f

☐ depends on g

☐ depends on f and g

g just defines what the output is.
even if g is non-linear the network
can be linear.

for example, NNC is a linear model,
but the output uses sin function

1.4 [2 Points] Given a fixed regularization parameter $\lambda > 0$ ridge regression...

...always has a ~~lower~~ ^{higher} training error than linear regression.

☐ True

☒ False

岭回归不总是有比线性回归更低的训练误差。这是因为岭回归通过正则化降低了模型复杂度，防止过拟合，因此在训练集上可能无法达到线性回归那样低的误差，特别是当线性回归模型已经能够很好地拟合训练数据时。

...always has a lower ~~test~~ error than linear regression.

☐ True

☒ False

2.

岭回归也不总是有比线性回归更低的测试误差。虽然正则化帮助模型提高了泛化能力，但这并不意味着在所有情况下岭回归都优于线性回归。是否能够减少测试误差取决于很多因素，包括数据的特点、正则化参数的大小等。

1.5 [1.5 Points] Assume the covariance between observations X and their labels Y is less than 0. Which statement(s) are always true?

☒ The regression function given by ordinary least squares will have negative slope.

☐ The correlation between X and Y is -1 .

☐ Either X or Y has negative variance.

1.6 [1.5 Points] Which statement(s) are true for kernel methods (as presented in this course)?

☒ You need to store all the training data to predict labels of new data points.

☒ For every valid kernel function k , there is a feature map ϕ , such that $k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$.

☐ Kernel ridge regression gives better results and is ~~faster~~ to compute than ridge regression.

1.7 [3 Points] Name 3 algorithms we discussed in the lecture that can be used to address classification problems.

NCC, Perceptron, LDA, MLP

1.8 [1 Point] Name one situation in which it makes sense to use a kernelized algorithm? (1-2 sentences)

$d \gg n$

linearly separable in feature space, but not in input space
(if the mapping is expensive to compute)

few data a lot of dimensions

1.9 [1 Point] "When the amount of data increases and the algorithm and hyperparameters stay the same, overfitting is more likely." True or false? Explain in one or two sentences.

false

more training data is better

2. Linear Classification [25 Points total]

2.1 [8 Points] Recall that the error function of the Perceptron was given by:

$$\mathcal{E}(\mathbf{w}) = - \sum_{m \in \mathcal{M}} \mathbf{w}^\top \mathbf{x}_m y_m \quad (1)$$

Briefly define, in words, \mathbf{w} , \mathbf{x}_m , y_m , and \mathcal{M} . [2 Points]

\mathbf{w} is the weight vector that the Perceptron learns
 \mathbf{x}_m are the data points in the training set, $\mathbf{x}_m \in \mathbb{R}^d$
 y_m are the labels $y_m \in \{-1, +1\}$
 \mathcal{M} is the index set of misclassified data points

Why do we include a minus sign in front of the summation? [2 Points]

- Because $\mathbf{w}^\top \mathbf{x}_m$ is the opposite sign of y_m for misclassified data points
therefore $\mathbf{w}^\top \mathbf{x}_m y_m$ is always negative for misclassified data points
- We want our error function to be positive and large when we have a lot of misclassified data points

Assume we do not want to determine \mathcal{M} before calculating the error of \mathbf{w} . Therefore, we rewrite the error function as $\mathcal{E}_{new}(\mathbf{w}) = \sum_{i \in \mathcal{X}} \max(0, -\mathbf{w}^\top \mathbf{x}_i y_i)$, where \mathcal{X} is the index set of all data points.

Show that \mathcal{E}_{new} is equivalent to the function \mathcal{E} given as equation (1) above. [4 Points].

- if \mathbf{x}_i is correctly classified, $-\mathbf{w}^\top \mathbf{x}_i y_i < 0$ and $\max(0, -\mathbf{w}^\top \mathbf{x}_i y_i) = 0$
i.e. correctly classified data points do not contribute to the error
- if \mathbf{x}_i is wrongly classified, $\max(0, -\mathbf{w}^\top \mathbf{x}_i y_i) = -\mathbf{w}^\top \mathbf{x}_i y_i$

2.2 [3 Points] Again recall the error of the Perceptron (see equation (1)). Assuming that we use

the augmented notation, i.e., that $\mathbf{w}^\top \mathbf{x}_m = [\beta \ w_1 \ \dots \ w_d] \begin{bmatrix} -1 \\ x_{m1} \\ \dots \\ x_{md} \end{bmatrix}$, calculate the update that

is applied to β in each iteration of Stochastic Gradient Descent (SGD) during the training of the Perceptron.

Hint! SGD uses the error of a single, randomly chosen misclassified data point, given by $\mathcal{E}_{x_m}(\mathbf{w}) = -\mathbf{w}^\top \mathbf{x}_m y_m$

$$\begin{aligned} \beta_{\text{new}} &= \beta_{\text{old}} - \eta \nabla \mathcal{E}_{x_m}(\mathbf{w}) & -\mathbf{w}^\top \mathbf{x}_m y_m &= -[\beta \ w_1 \ w_2 \ \dots \ w_d] \begin{bmatrix} -1 \\ x_{m1} \\ \vdots \\ x_{md} \end{bmatrix} y_m \\ & & &= -[\beta + w_1 x_{m1} + w_2 x_{m2} + \dots + w_d x_{md}] y_m \\ \nabla \mathcal{E}_{x_m}(\mathbf{w}) &= \frac{\partial \mathcal{E}_{x_m}(\mathbf{w})}{\partial \beta} = y_m \\ \beta_{\text{new}} &= \beta_{\text{old}} - \eta \cdot y_m \end{aligned}$$

2.3 [5 Points] You are given the following dataset consisting of two **infinite** classes C_{+1} (with all labels being +1) and C_{-1} (with all labels being -1):

$$C_{+1} := \left\{ \begin{bmatrix} x \\ x \end{bmatrix} \mid x \in \mathbb{R}, x > 0 \right\}$$

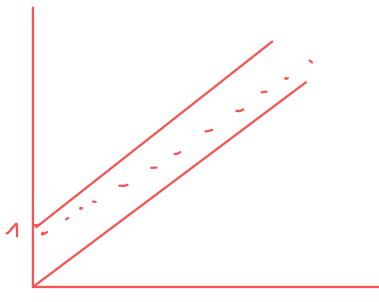
$$C_{-1} := \left\{ \begin{bmatrix} x \\ x+1 \end{bmatrix} \mid x \in \mathbb{R}, x > 0 \right\}$$

$$X := C_{+1} \cup C_{-1} \quad (\text{i.e., } X \text{ is the union of } C_{+1} \text{ and } C_{-1})$$

Which algorithm has a chance of calculating a solution in a finite amount of time? **[1 Point]**

- ☐ Perceptron using Gradient Descent
- ☒ Perceptron using Stochastic Gradient Descent
- ☐ LDA
- ☐ None of the above

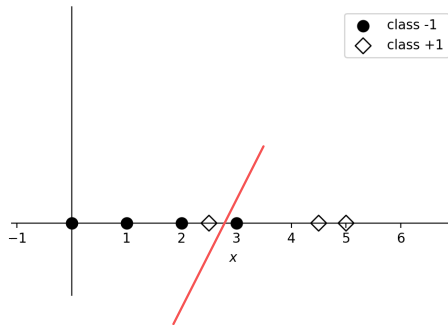
Why did you choose your answer? **[4 Points]**



The update step of GD is impossible to calculate within finite amount of time since we have infinitely many data points

The update step of SGD can be calculated in a finite amount of time

2.4 [9 Points] Consider the following 1-dimensional dataset $X = [0 \ 1 \ 2 \ 2.5 \ 3 \ 4.5 \ 5]$, with its corresponding labels $y = [-1 \ -1 \ -1 \ +1 \ -1 \ +1 \ +1]$. It is also illustrated in the plot below:



$$\begin{aligned} w_{+1} &= \frac{1}{3} (2.5 + 4.5 + 5) = 4 \\ w_{-1} &= \frac{1}{4} (6) = \frac{3}{2} \\ w &= 4 - \frac{3}{2} = \frac{5}{2} \\ \beta &= \frac{1}{2} (w_{+1}^2 - w_{-1}^2) = \frac{1}{2} (16 - \frac{9}{4}) \\ &= \frac{1}{2} \cdot \frac{55}{4} \\ &= \frac{55}{8} \end{aligned}$$

Consider the nearest centroid classifier (NCC). Write down the formulas for w and β that correspond to the NCC [2 Points]. Calculate w and β for the given dataset. [2 Point]. Calculate the **training accuracy** that NCC would achieve for this dataset [1 Point].

$$\begin{aligned} w &= \bar{x}_{+1} - \bar{x}_{-1} = \frac{1}{4} (0 + 1 + 2 + 3) + \frac{1}{3} (2.5 + 4.5 + 5) = 2.5 \\ \beta &= \frac{1}{2} (\bar{x}_{+1}^2 - \bar{x}_{-1}^2) = \frac{1}{2} (4^2 - (\frac{3}{2})^2) = \frac{1}{2} (16 - \frac{9}{4}) = \frac{55}{8} \\ \alpha &= \frac{5}{7} \end{aligned}$$

$$\begin{aligned} w\alpha - \beta &= 0 \\ \alpha &= \frac{\beta}{w} = \frac{\frac{55}{8}}{2.5} = \frac{11}{4} \end{aligned}$$

What is the highest accuracy that a linear discriminator can achieve on this dataset? [1 Point].

$$\frac{6}{7}$$

Denote by n_+ the number of training data points in class +1, and n_- in class -1. Now consider Fisher's linear discriminant analysis (LDA). In the lecture, we derived a β for LDA that depended on n_+ and n_- . What problem is addressed by this dependence? You consider a new training set where n_- is larger than in the first training set, and n_+ is the same. Should β increase, decrease, or stay the same, with respect to the first training set? [3 Points].

$$\begin{aligned} w &= \sum x_i (w_+ - w_-) \\ \beta &= w^T \left(\frac{w_+ + w_-}{2} \right) \left[+ \log \frac{n_-}{n_+} \right] \end{aligned}$$

Having β depend on n_+ , n_- reflects the class imbalance in the training set.

β increase, β moves away from the class that have more points

3. Kernel methods and Kernel Ridge Regression [7 Points total]

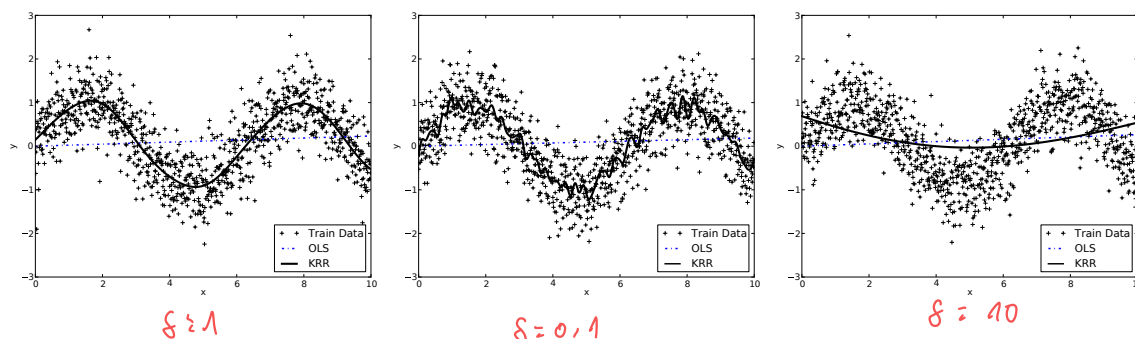
3.1 [4 Points]

1. [3 points] We used a Kernel Ridge Regression with a Gaussian kernel $k(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2\sigma^2}\right)$ on training data that follows a sine-function. Below you find the results for three different kernel widths. Indicate which of the following three labels corresponds to which of the plots:

$$\sigma = 10,$$

$$\sigma = 1,$$

$$\sigma = 0.1$$



2. [1 point] Explain intuitively how the kernel width σ affects the learned model.

Wenn σ is small, $k(x,y)$ rapidly decays to 0
thus, only training data in a close neighborhood of a test data influence the prediction

3.2 [3 Points] You are given the following feature map

$$\phi(x)^\top = (x_1^2, \sqrt{2}x_1x_2, x_2^2), \text{ where } x \in \mathbb{R}^2.$$

Show that the dot product defines a kernel function, i.e. show that

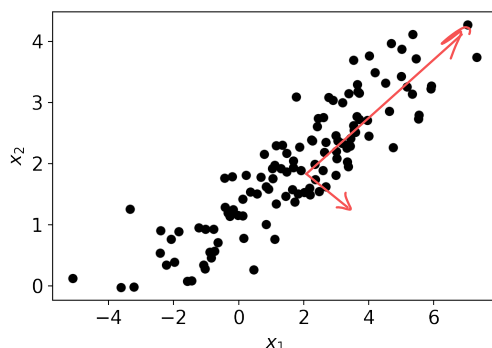
$$\phi(x)^\top \phi(y) = k(x, y) = (x^\top y)^2$$

where again $x, y \in \mathbb{R}^2$.

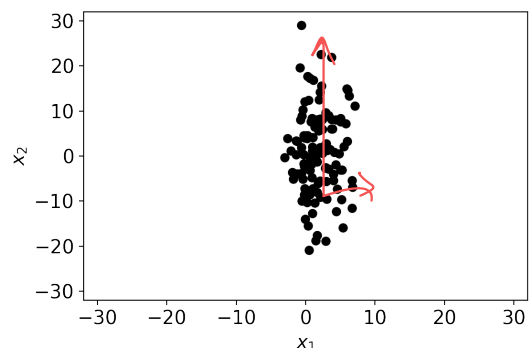
$$\begin{aligned} (x_1^2, \sqrt{2}x_1x_2, x_2^2) \begin{pmatrix} y_1^2 \\ \sqrt{2}y_1y_2 \\ y_2^2 \end{pmatrix} &= x_1^2 y_1^2 + 2x_1x_2 y_1 y_2 + x_2^2 y_2^2 \\ &= (x_1 y_1 + x_2 y_2)^2 \\ &= \left[(x_1 \ x_2) \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \right]^2 \\ &= (x^\top y)^2 = k(x, y) \end{aligned}$$

4. Unsupervised Learning [13 Points total]

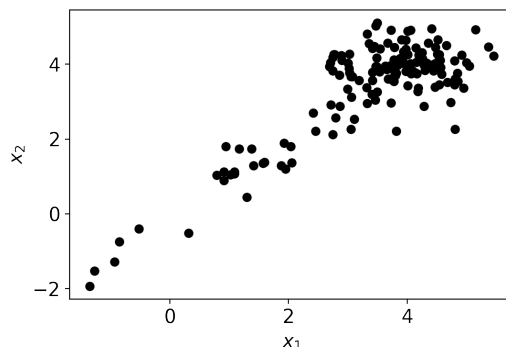
4.1 [7 Points] Below you see three plots of data sets (called “A”, “B”, and “C”).



Data Set “A”



Data Set “B”



Data Set “C”

Which one is most likely not stemming from a Gaussian distribution? [1 Point]

C

For two of the datasets above, draw plausible first and second principal components into their corresponding plots [2 Points]. Also give a **plausible** ratio $\frac{\lambda_1}{\lambda_2}$ for each of the two datasets, where λ_1 and λ_2 are the eigenvalues of the covariance matrix [2 Points].

Dataset 1: $\frac{\lambda_1}{\lambda_2} = 5$

Dataset 2: $\frac{\lambda_1}{\lambda_2} = 6$

(> 1 3/4)

You want to calculate the variance of each dataset along each dimension. For that you wrote your own python function `var`, which measures the average difference of a single dimension of a dataset to its corresponding mean:

```
def var(X : np.ndarray) -> float:
    variance = 0
    mean = np.mean(X)
    n_datapoints = X.shape[0]
    for data_index in range(n_datapoints):
        variance += X[data_index] - mean
    return variance/n_datapoints
```

You have called this function on each dimension of the datasets "B" and "C" and stored them in vectors, where each row corresponds to the estimated variance of that dimension. Unfortunately you do not get very plausible results:

Variance of data set "B": $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$

Variance of data set "C": $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$

Give a plausible output of the function on dataset "A" [1 Points]:

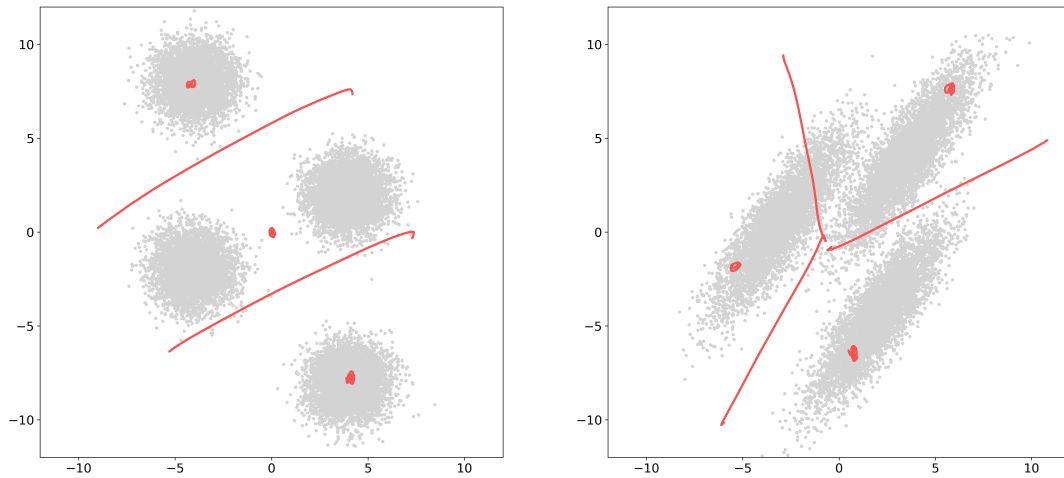
$\begin{bmatrix} 0 \\ 0 \end{bmatrix}$

Why does your function `var` give these wrong results? [2 Points]

You forgot to include a square around the difference between the data and the mean

4.2 [6 Points]

1. [4 points] In the Figure below you find two exemplary unlabeled data sets. Draw for each a plausible result of the k-means algorithm after 50 iterations, i.e. mark plausible **clusters** and **cluster centers** found by k-means for $k = 3$. Initial cluster centers were randomly drawn from the set of data points.



2. [2 points] For your bachelor thesis, you apply a classification algorithm to very high dimensional data you have recorded. Your supervisor is concerned that your features are still very correlated and suggests applying an unsupervised algorithm to reduce the dimensionality of your data before applying the classification algorithm. What does she mean by that? Explain briefly in 1–2 sentences and state one possible algorithm which we discussed in the lecture.

Correlation between features can have negative effect on classification.

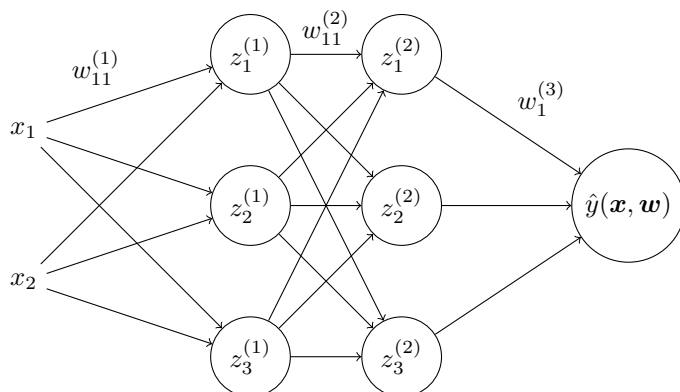
→ redundancy

→ reduce dimension without losing information

PCA

5. MLP [14 Points total]

5.1 [10 Points] Consider the following MLP architecture:



In this MLP, $z_1^{(1)}$ is calculated as

$$z_1^{(1)} = f(a_1^{(1)}) = f(w_{11}^{(1)}x_1 + w_{21}^{(1)}x_2),$$

where f is a non-linear activation function used to calculate all activations. The other $z_i^{(j)}$ are defined analogously, and \hat{y} does not use any activation function.

For an input \mathbf{x} with label y , let

$$\mathcal{E}(\mathbf{x}, \mathbf{w}, y) := \frac{1}{2} (\hat{y}(\mathbf{x}, \mathbf{w}) - y)^2$$

be the error function of the MLP.

- [2 points] You would like to update the weights in the MLP to decrease its error. Given a learning rate $\eta > 0$, which of the following examples of weight updates are correct?

- ☒ $w_{11}^{(1)} \leftarrow w_{11}^{(1)} - \eta \frac{\partial \mathcal{E}}{\partial w_{11}^{(1)}}$
- ☐ $w_{11}^{(1)} \leftarrow w_{11}^{(1)} - \eta \frac{\partial \mathcal{E}}{\partial y}$
- ☐ None of the above.

- [8 points] You would like to update the weights $w_1^{(3)}$ and $w_{11}^{(2)}$. Calculate the weight updates, i.e. derive \mathcal{E} with respect to both these weights. Carry out all partial derivatives iteratively. For the derivative of f , you can simply write f' .

Hint! Chain rule for nested functions:

$$\frac{\partial}{\partial x_j} f(g_1(\mathbf{x}), \dots, g_n(\mathbf{x})) = \sum_{i=1}^n \left(\frac{\partial}{\partial g_i} f(g_1, \dots, g_n) \right) \cdot \left(\frac{\partial}{\partial x_j} g_i(\mathbf{x}) \right).$$

$$\hat{y} = z_1^{(2)} w_1^{(3)} + z_2^{(2)} w_2^{(3)} + z_3^{(2)} w_3^{(3)}$$

$$\frac{\partial \mathcal{E}}{\partial w_1^{(3)}} = \frac{\partial \mathcal{E}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1^{(3)}}$$

$$= (\hat{y} - y) \cdot z_1^{(2)}$$

$$\frac{\partial \mathcal{E}}{\partial w_{11}^{(2)}} = \frac{\partial \mathcal{E}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_{11}^{(2)}}$$

$$= (\hat{y} - y) \cdot w_1^{(3)} \cdot \frac{\partial z_1^{(2)}}{\partial w_{11}^{(2)}}$$

$$= (\hat{y} - y) \cdot w_1^{(3)} \cdot \frac{\partial f(a_1^{(2)})}{\partial w_{11}^{(2)}}$$

$$= (\hat{y} - y) \cdot w_1^{(3)} \cdot f'(a_1^{(2)}) \cdot \frac{\partial a_1^{(2)}}{\partial w_{11}^{(2)}}$$

$$= (\hat{y} - y) \cdot w_1^{(3)} \cdot f'(a_1^{(2)}) \cdot z_1^{(1)}$$

$$\frac{\partial \mathcal{E}}{\partial w_{11}^{(2)}} = \frac{\partial \frac{1}{2}(\hat{y} - y)^2}{\partial (\hat{y} - y)} \cdot \frac{\partial (\hat{y} - y)}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial w_{11}^{(2)}}$$

5.2 [4 Points] Consider the same MLP as above. In this subtask, we will be using the tanh activation function, i.e., $f(z_i^{(k)}) = \tanh(z_i^{(k)})$. Further imagine that we have determined well-working values for all of the weights of the network. Describe alternative weights $w_{ij}'^{(2)}$ and $w_i'^{(3)}$ that are different to $w_{ij}^{(2)}$ and $w_i^{(3)}$ but lead to the **exact same** overall network output.

Hint! Consider the fact that tanh is an odd function, i.e., that $\tanh(-a) = -\tanh(a)$

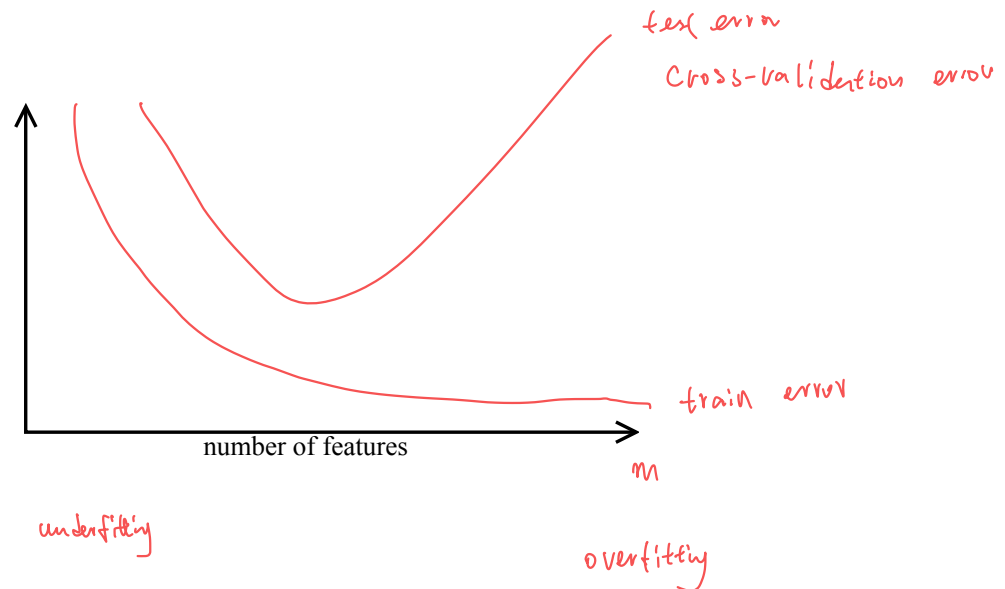
$$w_{ij}'^{(2)} \leftarrow -w_{ij}^{(2)}$$

$$w_i'^{(3)} \leftarrow -w_i^{(3)}$$

6. Cross-Validation [8 Points total]

6.1 [3 Points] Suppose you model the non-linear relationship between a one-dimensional input x and a one-dimensional output y as an m th order polynomial, i.e. $y = w_0 + w_1x + w_2x^2 + \dots + w_mx^m$. The number of training points is fixed, and you estimate the parameters w_0, w_1, \dots, w_m by linear regression.

Draw a graph showing two curves: training error vs. the number of features m and cross-validation error vs. the number of features m , annotate both curves. Include cases where over- and under-fitting occur and mark them in the plot.



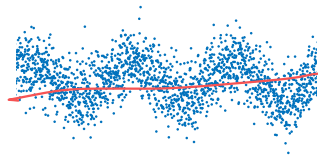
6.2 [2 Points] Find the bugs in the cross-validation algorithm below and correct them.

Algorithm 1 Cross-Validation

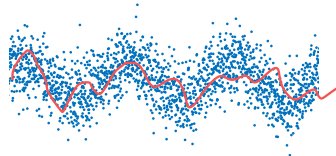
Require: Data $\{(x_1, y_1), \dots, (x_N, y_N)\}$, Number of CV folds F

- 1: Split data in F overlapping folds *disjunct*
 - 2: **for** Fold $f = 1, \dots, F$ **do**
 - 3: Train model on folds $\{1, \dots, F\} \setminus f$.
 - 4: Compute prediction on fold f
 - 5: **end for**
 - 6: **return** average prediction error
-

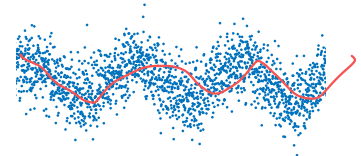
6.3 [3 Points] Below you find 3 equal plots with data points. Sketch possible solutions from a polynomial regression, one that under-fits (a), over-fits (b) and one good fit (c)



(a) under-fitting



(b) over-fitting



(c) good fit