

Note:

- During the attendance check a sticker containing a unique code will be put on this exam.
- This code contains a unique number that associates this exam with your registration number.
- This number is printed both next to the code and to the signature field in the attendance check list.

Maschinelles Lernen

Exam: IN2064 / Endterm
Examiner: Prof. Dr. Stephan Günnemann

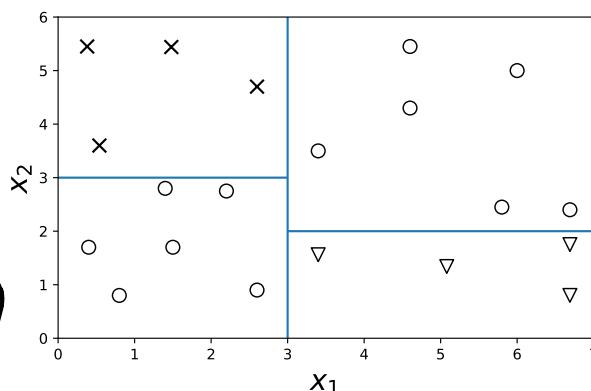
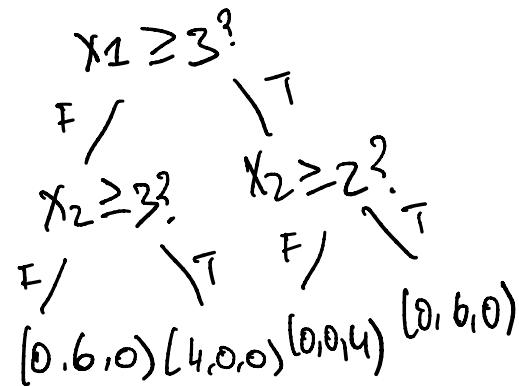
Date: Wednesday 13th April, 2022
Time: 08:00 – 10:00

Working instructions

- This graded exercise consists of **48 pages** with a total of **11** problems and four versions of each problem.
Please make sure now that you received a complete copy of the graded exercise.
- Use the problem versions specified in your personalized submission sheet on TUMExam. Different problems may have different versions: e.g. Problem 1 (Version A), Problem 5 (Version C), etc. If you solve the wrong version you get **zero** points.
- The total amount of achievable credits in this graded exercise is **82**.
- This document is copyrighted and it is **illegal** for you to distribute it or upload it to any third-party websites.
- Do **not** submit the problem descriptions (this document) to TUMexam
- You can ignore the “student sticker” box above.

Problem 1: Decision Trees (Version A) (8 credits)

You are given a dataset with points from three different classes and want to classify them using a decision tree. The plot below illustrates the data points (class labels are indicated by the symbols [\times , \circ , ∇]) and the decision boundaries of a decision tree of depth 2. Each decision boundary corresponds to a specific decision node.



- 0 a) Draw the corresponding decision tree. Make sure that you include the feature (x_1 or x_2) and threshold of each inner node (i.e. decision node). Also, for each inner node and leaf node, indicate the number of samples of each class that pass the node.

- 0 b) Compute the Gini index of each node of your decision tree.

1
2

- 0 c) Compute the misclassification rate of each node of your decision tree.

1
2

- 0 d) Could this tree have been obtained from greedy optimization (as discussed in the lecture) using the Gini index as the impurity measure? Briefly justify your answer.

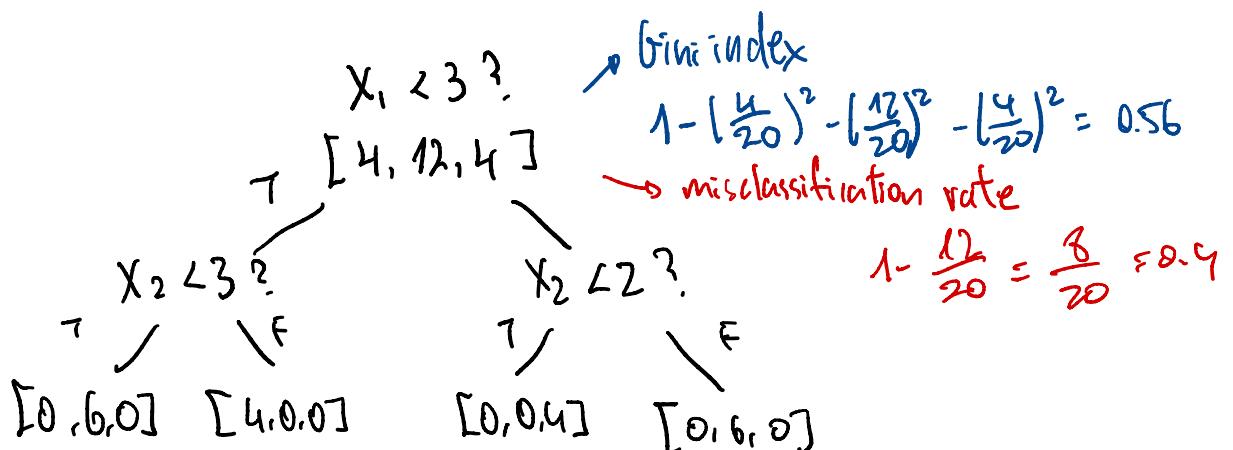
Note: You may refer to the solution of subproblem b).

Yes, since each node decreases the impurity

- 0 e) Could this tree have been obtained from greedy optimization (as discussed in the lecture) using the misclassification rate as the impurity measure? Briefly justify your answer.

Note: You may refer to the solution of subproblem c).

No, because the root node does not decrease impurity, hence it would not have been added by the algorithm



Problem 2: Probabilistic Inference (Version A) (8 credits)

Consider the probabilistic model composed of the two following densities

$$\begin{aligned} \mathbb{P}(\theta | \alpha) &= \frac{\Gamma\left(\sum_{c=1}^C \alpha_c\right)}{\prod_{c=1}^C \Gamma(\alpha_c)} \prod_{c=1}^C \theta_c^{\alpha_c-1} \\ \mathbb{P}(x | \theta) &= \prod_{c=1}^C \theta_c^{\mathbb{I}(x=c)} \text{ with } \sum_{c=1}^C \theta_c = 1 \end{aligned}$$

with probabilistic parameters $\theta \in [0, 1]^C$, a set of observations $\mathcal{D} = \{x_1, \dots, x_N\}$ consisting of N samples $x_i \in \{1, \dots, C\}$, and some given $\alpha \in \mathbb{R}^C$. Here, Γ is the Gamma function.

Derive the maximum a posteriori estimate of the parameter θ denoted θ_{MAP} .

Note: The maximum a posteriori estimate must also fulfill $\sum_{c=1}^C \theta_{MAP,c} = 1$

$$\begin{aligned} \log \mathbb{P}(\theta | \{x_1, \dots, x_N\}, \alpha) &= \sum_{i=1}^N \log \mathbb{P}(x_i | \theta) + \log \mathbb{P}(\theta | \alpha) + D \\ &= \sum_{i=1}^N \sum_{c=1}^C \mathbb{I}(x_i=c) \log \theta_c + \sum_{c=1}^C (\alpha_{c-1}) \log \theta_c - \log B(\alpha) + D \\ &= \sum_{c=1}^C N_c \log \theta_c + \sum_{c=1}^C (\alpha_{c-1}) \log \theta_c - \log B(\alpha) + D \\ &= \sum_{c=1}^C (N_c + \alpha_{c-1}) \log \theta_c + D' \end{aligned}$$

$B, D, D' \in \mathbb{R}$ are constant w.r.t. θ and can be ignored in our optimization problem

Use the normalization constraint $\sum_{c=1}^C \theta_c = 1$ to compute the lagrangian:

$$L = \sum_{c=1}^C (N_c + \alpha_{c-1}) \log \theta_c + \lambda \left(1 - \sum_{c=1}^C \theta_c\right)$$

Set the derivative to 0:

$$\frac{\partial L}{\partial \theta_c} = \frac{(N_c + \alpha_{c-1})}{\theta_c} - \lambda = 0$$

$$\theta_c = \frac{(N_c + \alpha_{c-1})}{\lambda}$$

$$\theta_{MAP,c} = \frac{(N_c + \alpha_{c-1})}{\sum_c (N_c + \alpha_{c-1})}$$

Problem 3: Linear Regression (Version A) (8 credits)

We want to perform regression on a dataset consisting of N samples $\mathbf{x}_i \in \mathbb{R}^D$ ($D > 1$) with corresponding targets $y_i \in \mathbb{R}$ (represented compactly as $\mathbf{X} \in \mathbb{R}^{N \times D}$ and $\mathbf{y} \in \mathbb{R}^N$).

Assume that we have a deep neural network with L layers such that $f(\mathbf{x}, \mathbb{W}) = \mathbf{W}^{(L)}\sigma(\mathbf{W}^{(L-1)} \dots (\sigma(\mathbf{W}^{(0)}\mathbf{x})))$ where $\mathbb{W} = \{\mathbf{W}^{(0)}, \dots, \mathbf{W}^{(L)}\}$. Here, σ is the sigmoid activation function applied element-wise.

- 0 a) Suppose the weights $\mathbf{W}^{(0)}, \dots, \mathbf{W}^{(L-1)}$ are non-zero, known and fixed and we want to solve the following optimization problem over the weights of the *last* layer:

$$2 \\ 3 \\ 4 \\ \arg \min_{\mathbf{W}^{(L)}} \frac{1}{2} \sum_{i=1}^N (f(\mathbf{x}_i, \mathbb{W}) - y_i)^2$$

Can you compute the optimal weights $\mathbf{W}^{*(L)}$ of this optimization problem in closed-form? Justify your response. If yes, provide the closed-form solution $\mathbf{W}^{*(L)}$. If no, explain how we can approximately solve this problem in practice.

- 0 b) Suppose the weights $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L)}$ are non-zero, known and fixed and we want to solve the following optimization problem over the weights of the *first* layer:

$$2 \\ 3 \\ 4 \\ \arg \min_{\mathbf{W}^{(0)}} \frac{1}{2} \sum_{i=1}^N (f(\mathbf{x}_i, \mathbb{W}) - y_i)^2$$

Can you compute the optimal weights $\mathbf{W}^{*(0)}$ of this optimization problem in closed-form? Justify your response. If yes, provide the closed-form solution $\mathbf{W}^{*(0)}$. If no, explain how we can approximately solve this problem in practice.

a) Yes, the problem is linear wrt the weights $\mathbf{W}^{(L)}$. By setting $\Phi = \sigma(\mathbf{W}^{(L-1)} \dots \sigma(\mathbf{W}^{(0)}\mathbf{x}))$, the optimal solution is $\mathbf{W}_L = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$

b) No, the problem is highly non-linear wrt the weights $\mathbf{W}^{(0)}$. Note that we cannot invert the network either since $D > 1$ and the weight matrices are non-zero. We could approximate the optimal solution by optimizing with SGD

Problem 4: Linear classification (Version A) (9 credits)

A friend of yours owns an ice cream truck. To improve his profitability, he wants to predict whether it is profitable for him to drive to the lake and sell ice cream. For this reason, he created a dataset about his days at the lake. Whenever your friend was not able to collect the data, he noted down a “-”. The last row represents the data for today.

x_1 Rain probability	x_2 Temperature forecast	x_3 # Visitors at 9 am	x_4 Public holiday	x_5 Weekday	y Profitable day
10%	25.1	51	True	False	True
20%	-	2	False	True	True
90%	37.8	5	False	False	False
-	23.0	27	True	True	True
45%	-2.1	21	True	True	False
0%	-10.2	18	True	False	?

0
1
2

a)

Your friend is a big fan of logistic regression. Can you name a reason why logistic regression might not be such a good idea here? Justify your answer.

We have missing data and logistic regression treats the data as \mathbb{R}^5 and Gaussian distributed

0
1
2
3
4

b)

You, finally, convinced your friend to use a Naïve Bayes classifier. However, he is unsure which distribution is a good choice for each feature. Which of the following class-conditional distributions should he use for each of the features x_1, x_2, x_3, x_4 and x_5 ? (a) Bernoulli, (b) Normal distribution, (c) Poisson distribution, (d) Beta distribution, and (e) exponential distribution. Justify your answer.

$x_2 \rightarrow$ Normal (continuous values)

$x_3 \rightarrow$ Poisson (integer values)

c) $x_4, x_5 \rightarrow$ Bernoulli, the data is binary

$x_1 \rightarrow$ d [density between 0,1]

0
1
2
3

After you have fitted the model, your friend is still not happy with the predictive performance of the model. Name a reason why the dataset might be not ideal for a Naïve Bayes model (assuming the data follows perfectly the chosen distributions). Justify your answer.

- the features are likely correlated

Problem 5: Optimization – Convexity (Version A) (8 credits)

Consider the two functions

$$f_a(\mathbf{x}) = \max [0, (\mathbf{x} - \mathbf{a})^\top \mathbf{A}(\mathbf{x} - \mathbf{a})]$$

$$f_b(\mathbf{x}) = 2^{\max_i |\mathbf{x}_i - \mathbf{a}_i|}$$

with $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{a} \in \mathbb{R}^n$ and positive semidefinite $\mathbf{A} \in \mathbb{R}^{n \times n}$.

For your reference, here are the convexity-preserving rules from the lecture. Let $f_1 : \mathbb{R}^d \rightarrow \mathbb{R}$ and $f_2 : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex functions, and $g : \mathbb{R}^d \rightarrow \mathbb{R}$ be a concave function, then:

1. $h(\mathbf{x}) = f_1(\mathbf{x}) + f_2(\mathbf{x})$ is convex
2. $h(\mathbf{x}) = \max \{f_1(\mathbf{x}), f_2(\mathbf{x})\}$ is convex
3. $h(\mathbf{x}) = c \cdot f_1(\mathbf{x})$ is convex if $c \geq 0$
4. $h(\mathbf{x}) = c \cdot g(\mathbf{x})$ is convex if $c \leq 0$
5. $h(\mathbf{x}) = f_1(\mathbf{Ax} + \mathbf{b})$ is convex (\mathbf{A} matrix, \mathbf{b} vector)
6. $h(\mathbf{x}) = m(f_1(\mathbf{x}))$ is convex if $m : \mathbb{R} \rightarrow \mathbb{R}$ is convex and nondecreasing.

0
1
2
3
4

a)

Prove or disprove that $f_a(\mathbf{x})$ is convex in \mathbf{x} . If you use a convexity preserving operation (stated above), clearly refer to the used rule. Any intermediate steps that are not direct applications of the given convexity-preserving operations must be proven.

0
1
2
3
4

b)

Prove or disprove that $f_b(\mathbf{x})$ is convex in \mathbf{x} . If you use a convexity preserving operation (stated above), clearly refer to the used rule. Any intermediate steps that are not direct applications of the given convexity-preserving operations must be proven.

a) $f_a(\mathbf{x})$ is convex if $(\mathbf{x} - \mathbf{a})^\top \mathbf{A}(\mathbf{x} - \mathbf{a})$ is convex (rule 2)

$$(\mathbf{x} - \mathbf{a})^\top \mathbf{A}(\mathbf{x} - \mathbf{a}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} - 2\mathbf{x}^\top \mathbf{A}\mathbf{a} + \mathbf{a}^\top \mathbf{A}\mathbf{a}$$

$-2\mathbf{x}^\top \mathbf{A}\mathbf{a} + \mathbf{a}^\top \mathbf{A}\mathbf{a}$ is linear and, therefore, convex in \mathbf{x}

$\mathbf{x}^\top \mathbf{A} \mathbf{x}$ is convex since $\nabla_{\mathbf{x}}^2 [\mathbf{x}^\top \mathbf{A} \mathbf{x}] = 2\mathbf{A}$ and the fact that \mathbf{A} is positive semidefinite

Last, $\mathbf{x}^\top \mathbf{A} \mathbf{x} - 2\mathbf{x}^\top \mathbf{A}\mathbf{a} + \mathbf{a}^\top \mathbf{A}\mathbf{a}$ is a sum of convex functions and, thus, also convex

In conclusion, $f_a(\mathbf{x})$ is convex

b) Since 2^x is a convex and nondecreasing function, $f_b(\mathbf{x})$ is convex if $h(\mathbf{x}) = \max_i |\mathbf{x}_i - \mathbf{a}_i|$ is convex (rule 6)

Problem 6: Deep learning (Version A) (8 credits)

0
1
2
3
4
5
6
7
8

Given the input $x \in \mathbb{R}$, output $y \in \mathbb{R}$, and learnable parameters $\alpha, \beta \in \mathbb{R}$ we define the following four models:

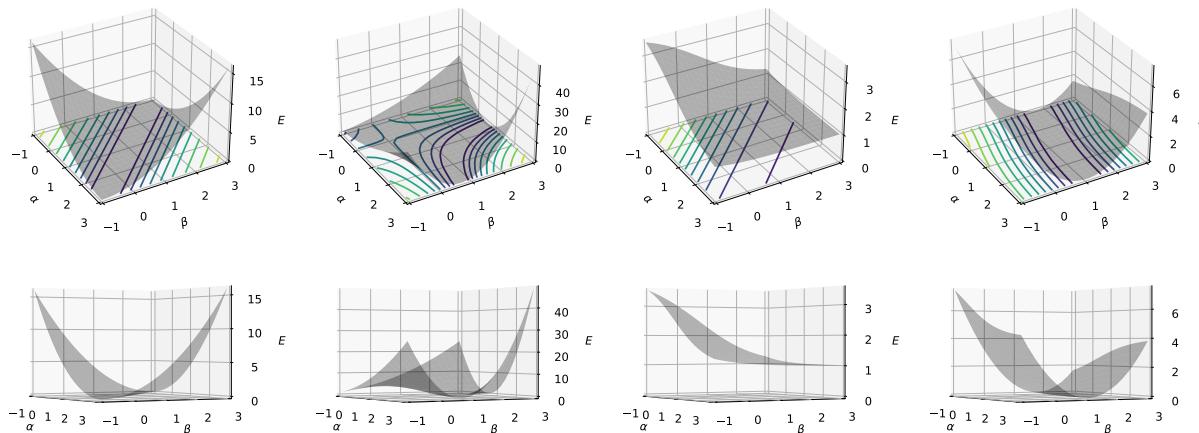
1. $f(x) = \sigma(\alpha x + \beta)$
2. $f(x) = \alpha \beta x$
3. $f(x) = \alpha x + \beta$
4. $f(x) = \sigma(\alpha x) + \beta$

where $\sigma(z) = \frac{1}{1+e^{-z}}$ is the sigmoid function. Additionally, for the true output y and the model output $f(x)$, we define the error as $E = (y - f(x))^2$.

In the four plots below, we fix the input and output to $x = 1$ and $y = 2$ and plot the error curve of each model (1.-4.) as a function of α and β .

Connect the plots (a)-(d) with the models 1.-4. Justify your answer.

NOTE: The plots on the top and the bottom row are the same but shown at a different angle. For example, the two plots in (a) show the same error curve which corresponds to one of the models. The contour plot that spans the α - β plane shows the projected error contours. Darker colors denote lower values and the values along a single contour are constant.



(a)
 $f(x) = \alpha x + \beta$

(b)
 $f(x) = \alpha \beta x$

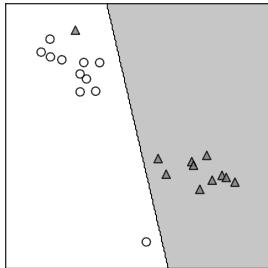
(c)
 $f(x) = \sigma(\alpha x + \beta)$

(d)
 $f(x) = \sigma(\alpha x) + \beta$

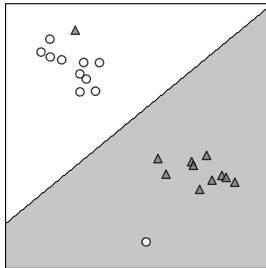
Problem 7: Support Vector Machines (Version A) (4 credits)

0
1
2
3
4

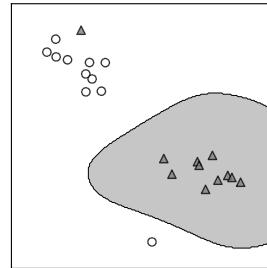
Below are shown four different decision boundaries plotted for four different models trained on the same dataset. Connect the decision boundaries (a)-(d) with the models 1-4. Justify your answer.



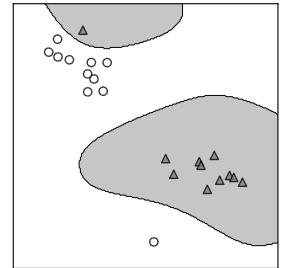
(a)



(b)



(c)



(d)

1. LogisticRegression(lambda=1.0, penalty='l2') *a*

2. SVM(C=0.01, kernel='rbf') *c*

3. SVM(C=0.01, kernel='linear') *b*

4. SVM(C=1000, kernel='rbf') *d*

Note: rbf denotes radial basis (Gaussian) kernel. kernel='linear' corresponds to the plain SVM discussed in Section 2 of lecture 9. Parameters lambda and C appear in the loss for the respective models, as defined in the lecture.

a → separates more points than *b* and if has smaller margin meaning it belongs to logistic

Problem 8: Kernels (Version A) (8 credits)

Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$ and let $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ be a neural network with two hidden layers and ReLU activations between them. Let $p(\mathbf{x})$ denote the probability density function on \mathbb{R}^N . Prove or disprove that the following functions are valid kernels.

0 a) $k(\mathbf{x}, \mathbf{y}) = \sum_i^M f(\mathbf{x})_i + \sum_i^M f(\mathbf{y})_i$

function f together with the summation

1
2
3
4

0 b) $k(\mathbf{x}, \mathbf{y}) = \log p(\mathbf{x}) \log p(\mathbf{y})$

1
2
3
4

Problem 8: Kernels (Version B) (8 credits)

Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$ and let $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ be a neural network with two hidden layers and ReLU activations between them. Let $p(\mathbf{x})$ denote the probability density function on \mathbb{R}^N . Prove or disprove that the following functions are valid kernels.

a) $k(\mathbf{x}, \mathbf{y}) = \left(\sum_i^M f(\mathbf{x})_i \right) \left(\sum_i^M f(\mathbf{y})_i \right)$

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3
<input type="checkbox"/>	4

b) $k(\mathbf{x}, \mathbf{y}) = \sum_z p(\mathbf{x}|z)p(\mathbf{y}|z)p(z)$, where $z \in \{1, \dots, Z\}$

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3
<input type="checkbox"/>	4

Problem 8: Kernels (Version C) (8 credits)

Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$ and let $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ be a neural network with two hidden layers and ReLU activations between them. Let $p(\mathbf{x})$ denote the probability density function on \mathbb{R}^N . Prove or disprove that the following functions are valid kernels.

0 a) $k(\mathbf{x}, \mathbf{y}) = \sum_i^M f(\mathbf{x})_i + \sum_i^M f(\mathbf{y})_i$

1

2

3

4

0 b) $k(\mathbf{x}, \mathbf{y}) = \sum_z p(\mathbf{x}|z)p(\mathbf{y}|z)p(z)$, where $z \in \{1, \dots, Z\}$

1

2

3

4

Problem 8: Kernels (Version D) (8 credits)

Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$ and let $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ be a neural network with two hidden layers and ReLU activations between them. Let $p(\mathbf{x})$ denote the probability density function on \mathbb{R}^N . Prove or disprove that the following functions are valid kernels.

a) $k(\mathbf{x}, \mathbf{y}) = \log p(\mathbf{x}) \log p(\mathbf{y})$

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3
<input type="checkbox"/>	4

b) $k(\mathbf{x}, \mathbf{y}) = \sum_i^M f(\mathbf{x})_i + \sum_i^M f(\mathbf{y})_i$

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3
<input type="checkbox"/>	4

Problem 9: Dimensionality Reduction (Version A) (6 credits)

Consider the following centered data matrix, where the data points are stored as rows

$$\begin{pmatrix} 0 & -5 & 0 & -2 \\ 0 & 0 & 2 & 6 \\ 2 & 1 & 0 & -2 \\ 0 & 0 & -2 & 6 \\ 1 & 1 & -1 & -2 \\ -2 & 1 & 0 & -2 \\ -2 & 1 & 0 & -2 \\ 0 & -1 & 0 & 0 \\ 1 & 1 & 1 & -2 \\ 0 & 1 & 0 & 0 \end{pmatrix}$$

We performed eigenvalue decomposition of the respective covariance matrix $\Sigma_x = \Gamma^T \Lambda \Gamma$. We obtained the following eigenvalue matrix

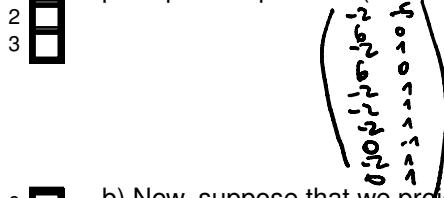
$$\begin{pmatrix} 9.6 & 0 & 0 & 0 \\ 0 & 3.2 & 0 & 0 \\ 0 & 0 & 1.4 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

and the following eigenvector matrix

$$\begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{pmatrix}$$

where the eigenvectors are stored as columns and sorted according to the eigenvalues in decreasing order.

- 0 a) Compute the coordinates of all data points when projected onto the subspace spanned by the top-2 principal components (i.e. the eigenvectors with the two largest eigenvalues). Justify your answer.



- 0 b) Now, suppose that we projected the original data points on the subspace defined by the THIRD principal component (i.e. the eigenvector with the third-largest eigenvalues). Compute the variance of the projected data. Justify your answer.

Problem 9: Dimensionality Reduction (Version B) (6 credits)

Consider the following centered data matrix, where the data points are stored as rows

$$\begin{pmatrix} 0 & 2 & 0.5 & -1 \\ 2 & 1 & -1 & 1 \\ 0 & -2 & -1 & 1 \\ 0 & 0 & 3.5 & 1 \\ -2 & 0 & 0 & -1 \\ 1 & -1 & 0 & -1 \\ -2 & 1 & -2 & 1 \\ -1 & -1 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & -1 & -1 & -1 \end{pmatrix}$$

We performed eigenvalue decomposition of the respective covariance matrix $\Sigma_X = \Gamma^T \Lambda \Gamma$. We obtained the following eigenvalue matrix

$$\begin{pmatrix} 2.05 & 0 & 0 & 0 \\ 0 & 1.6 & 0 & 0 \\ 0 & 0 & 1.4 & 0 \\ 0 & 0 & 0 & 0.8 \end{pmatrix}$$

and the following eigenvector matrix

$$\begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

where the eigenvectors are stored as columns and sorted according to the eigenvalues in decreasing order.

- a) Compute the coordinates of all data points when projected onto the subspace spanned by the top-2 principal components (i.e. the eigenvectors with the two largest eigenvalues). Justify your answer.

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3

- b) Now, suppose that we projected the original data points on the subspace defined by the THIRD principal component (i.e. the eigenvector with the third-largest eigenvalues). Compute the variance of the projected data. Justify your answer.

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3

Problem 9: Dimensionality Reduction (Version C) (6 credits)

Consider the following centered data matrix, where the data points are stored as rows

$$\begin{pmatrix} 1 & 1 & 0 & -4 \\ 2 & 2 & 1 & 2 \\ 0 & -4.5 & 0 & 1 \\ 0 & 1.5 & -2 & 1 \\ 1 & -2 & 0 & -1 \\ 1 & 0 & 0 & 1 \\ -2 & 1 & 1 & 0 \\ -1 & 0 & 0 & 1 \\ -1 & 0 & 0 & -2 \\ -1 & 1 & 0 & 1 \end{pmatrix}$$

We performed eigenvalue decomposition of the respective covariance matrix $\Sigma_x = \Gamma^T \Lambda \Gamma$. We obtained the following eigenvalue matrix

$$\begin{pmatrix} 3.35 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 1.4 & 0 \\ 0 & 0 & 0 & 0.6 \end{pmatrix}$$

and the following eigenvector matrix

$$\begin{pmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{pmatrix}$$

where the eigenvectors are stored as columns and sorted according to the eigenvalues in decreasing order.

- 0 a) Compute the coordinates of all data points when projected onto the subspace spanned by the top-2 principal components (i.e. the eigenvectors with the two largest eigenvalues). Justify your answer.
1
2
3

- 0 b) Now, suppose that we projected the original data points on the subspace defined by the THIRD principal component (i.e. the eigenvector with the third-largest eigenvalues). Compute the variance of the projected data. Justify your answer.
1
2
3

Problem 9: Dimensionality Reduction (Version D) (6 credits)

Consider the following centered data matrix, where the data points are stored as rows

$$\begin{pmatrix} 0 & 3 & 0 & 3 \\ -1 & 1 & 1 & 1 \\ 1 & -1 & 4 & 3 \\ 1 & 1 & 1 & -7 \\ -2 & -1 & 1 & -2 \\ 1 & -1 & -2 & 1 \\ 0 & -2 & -1 & 1 \\ 0 & 1 & -2 & 1 \\ 0 & 0 & -1 & -1 \\ 0 & -1 & -1 & 0 \end{pmatrix}$$

We performed eigenvalue decomposition of the respective covariance matrix $\Sigma_X = \Gamma^T \Lambda \Gamma$. We obtained the following eigenvalue matrix

$$\begin{pmatrix} 7.6 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0.8 \end{pmatrix}$$

and the following eigenvector matrix

$$\begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}$$

where the eigenvectors are stored as columns and sorted according to the eigenvalues in decreasing order.

- a) Compute the coordinates of all data points when projected onto the subspace spanned by the top-2 principal components (i.e. the eigenvectors with the two largest eigenvalues). Justify your answer.

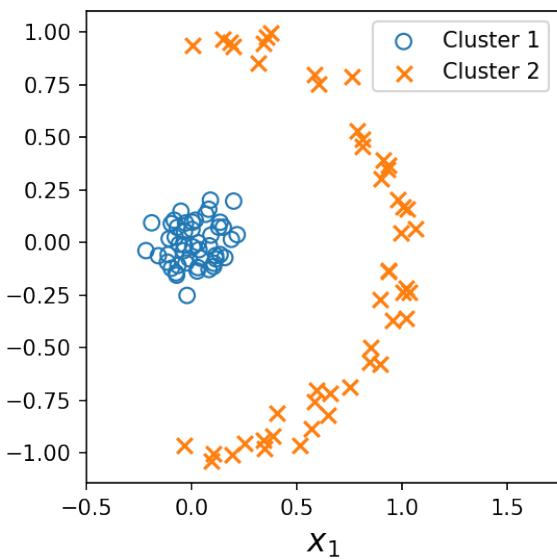
<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3

- b) Now, suppose that we projected the original data points on the subspace defined by the THIRD principal component (i.e. the eigenvector with the third-largest eigenvalues). Compute the variance of the projected data. Justify your answer.

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3

Problem 10: Clustering (Version A) (8 credits)

You would like to perform clustering on the following dataset.



- 0 a) Is K-means, in theory, able to recover the true clusters on the above dataset? Justify your answer.

1 No, because K-means with 2 clusters has linear decision boundaries

- 0 b) Is GMM with FULL covariance matrix SHARED across all clusters, in theory, able to recover the true clusters on the above dataset? Justify your answer.

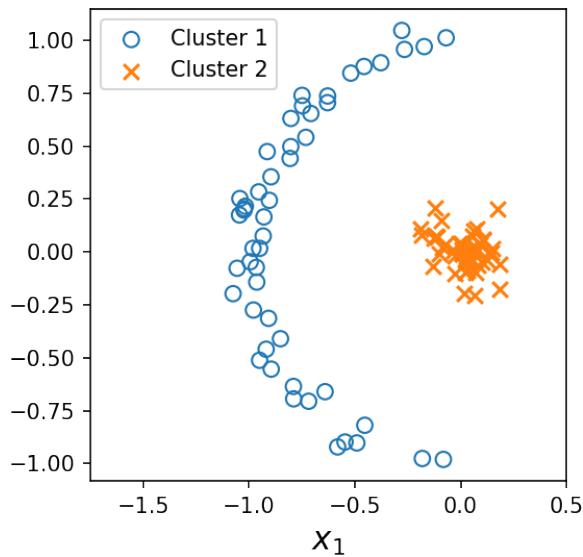
1 No

- 0 c) Is GMM with SEPARATE DIAGONAL covariance matrix for each cluster, in theory, able to recover the true clusters on the above dataset? Justify your answer.

1 Yes

Problem 10: Clustering (Version B) (8 credits)

You would like to perform clustering on the following dataset.



a) Is K -means, in theory, able to recover the true clusters on the above dataset? Justify your answer.

<input type="checkbox"/>	0
<input checked="" type="checkbox"/>	1
<input type="checkbox"/>	2

b) Is GMM with FULL covariance matrix SHARED across all clusters, in theory, able to recover the true clusters on the above dataset? Justify your answer.

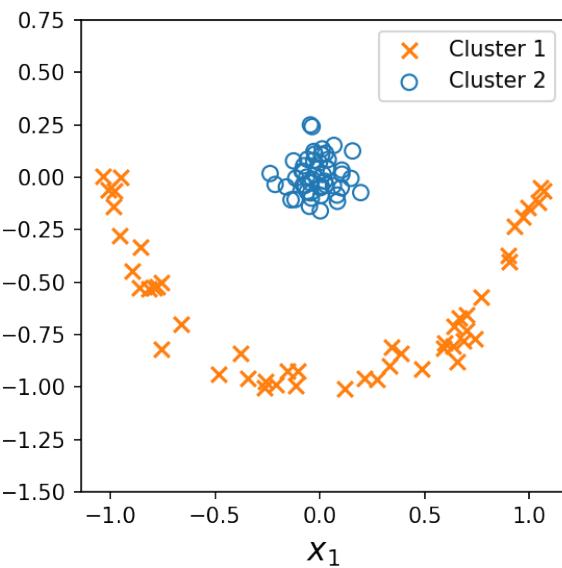
<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3

c) Is GMM with SEPARATE DIAGONAL covariance matrix for each cluster, in theory, able to recover the true clusters on the above dataset? Justify your answer.

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3

Problem 10: Clustering (Version C) (8 credits)

You would like to perform clustering on the following dataset.



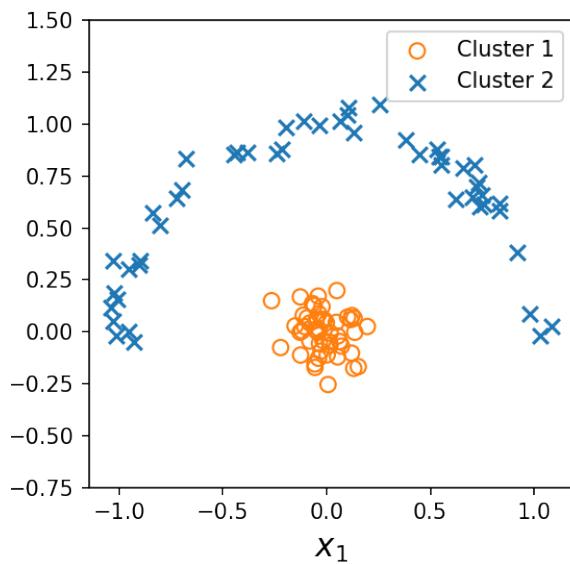
0 a) Is K -means, in theory, able to recover the true clusters on the above dataset? Justify your answer.
1
2

0 b) Is GMM with FULL covariance matrix SHARED across all clusters, in theory, able to recover the true clusters on the above dataset? Justify your answer.
1
2
3

0 c) Is GMM with SEPARATE DIAGONAL covariance matrix for each cluster, in theory, able to recover the true clusters on the above dataset? Justify your answer.
1
2
3

Problem 10: Clustering (Version D) (8 credits)

You would like to perform clustering on the following dataset.



a) Is K -means, in theory, able to recover the true clusters on the above dataset? Justify your answer.

<input type="checkbox"/>	0
<input checked="" type="checkbox"/>	1
<input type="checkbox"/>	2

b) Is GMM with FULL covariance matrix SHARED across all clusters, in theory, able to recover the true clusters on the above dataset? Justify your answer.

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3

c) Is GMM with SEPARATE DIAGONAL covariance matrix for each cluster, in theory, able to recover the true clusters on the above dataset? Justify your answer.

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3

Additional space for solutions—clearly mark the (sub)problem your answers are related to and strike out invalid solutions.

A large grid of squares, approximately 20 columns by 30 rows, intended for writing solutions. The grid is composed of thin black lines on a white background.

