

Eexam

Place student sticker here

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.

Machine Learning

Graded Exercise: IN2064 / Endterm
Examiner: Prof. Dr. Stephan Günnemann

Date: Tuesday 16th February, 2021
Time: 11:00 – 13:00

Working instructions

- This graded exercise consists of **36 pages** with a total of **10 problems**.
Please make sure now that you received a complete copy of the graded exercise.
- The total amount of achievable credits in this graded exercise is 48.
- 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

Left room from _____ to _____ / Early submission at _____

Problem 1: Probabilistic Inference (Version A) (4 credits)

Imagine a disease spreading across a small village. To make accurate forecasts for the necessary hospital beds, you want to estimate the severity of the disease with the following model. Let s be a measure for the severity of a disease. We assume a priori that s follows a standard normal distribution.

$$s \sim \mathcal{N}(0, 1) \propto \exp(-s^2)$$

The severity probabilistically influences the required days of hospital care t for a patient according to

$$t | s \sim \text{Exp}(\lambda) = \lambda \exp(-\lambda t) \text{ where } \lambda = s^2. \quad s^2 \propto \exp(-s^2 t_i)$$

After some time, you were able to collect $N = 5$ data points. The respective patients left the hospital after 3, 7, 3, 2 and 4 days.

Derive an a-posteriori most likely value of the severity s considering these observations. Justify your answer.

Note: You can assume that $s \neq 0$, that is, you can safely divide by s if necessary.

$$\begin{aligned} p(s | D) &\propto p(D | s) \cdot p(s) \\ &\propto \prod p(t_i | s) \cdot p(s) \\ &\propto \prod s^2 \exp(-s^2 t_i) \cdot \exp(-s^2) \end{aligned}$$

$$\begin{aligned} \ln p(s | D) &\propto \sum \ln [s^2 \exp(-s^2 t_i)] + \ln \exp(-s^2) \\ &\propto \sum \ln s^2 + \sum \ln \exp(-s^2 t_i) - s^2 \\ &\propto N \ln s^2 + \sum -s^2 t_i - s^2 \\ &\propto N \ln s^2 - s^2 \sum t_i - s^2 \end{aligned}$$

$$\frac{\partial (\ln p(s | D))}{\partial s} = \frac{2N}{s} - 2s \sum t_i - 2s \stackrel{!}{=} 0$$

$$N = s^2 (\sum t_i + 1)$$

$$s_{\text{ML}}^2 = \sqrt{\frac{N}{\sum t_i + 1}}$$

$$s_{\text{ML}}^* = \sqrt{\frac{5}{11 + 1}} = \sqrt{\frac{1}{2}} = \frac{1}{\sqrt{2}}$$

Problem 1: Probabilistic Inference (Version B) (4 credits)

Imagine a disease spreading across a small village. To make accurate forecasts for the necessary hospital beds, you want to estimate the severity of the disease with the following model. Let s be a measure for the severity of a disease. We assume a priori that s follows a standard normal distribution.

$$s \sim \mathcal{N}(0, 1) \propto \exp(-s^2)$$

The severity probabilistically influences the required days of hospital care t for a patient according to

$$t \mid s \sim \text{Exp}(\lambda) = \lambda \exp(-\lambda t) \text{ where } \lambda = s^2.$$

After some time, you were able to collect $N = 3$ data points. The respective patients left the hospital after 10, 11 and 5 days.

Derive an a-posteriori most likely value of the severity s considering these observations. Justify your answer.

Note: You can assume that $s \neq 0$, that is, you can safely divide by s if necessary.

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

Problem 1: Probabilistic Inference (Version C) (4 credits)

0 ☐ Imagine a disease spreading across a small village. To make accurate forecasts for the necessary hospital beds, you want to estimate the severity of the disease with the following model. Let s be a measure for the severity of a disease. We assume a priori that s follows a standard normal distribution.

1 ☐

2 ☐

$$s \sim \mathcal{N}(0, 1) \propto \exp(-s^2)$$

3 ☐

The severity probabilistically influences the required days of hospital care t for a patient according to

4 ☐

$$t \mid s \sim \text{Exp}(\lambda) = \lambda \exp(-\lambda t) \text{ where } \lambda = s^2.$$

After some time, you were able to collect $N = 4$ data points. The respective patients left the hospital after 8, 12, 9 and 6 days.

Derive an a-posteriori most likely value of the severity s considering these observations. Justify your answer.

Note: You can assume that $s \neq 0$, that is, you can safely divide by s if necessary.

Problem 2: k-nearest neighbors (Version A) (4 credits)

Consider a k -nearest neighbor classifier using Euclidean distance on a binary classification task. The prediction for an instance is the *majority* class of its k -nearest neighbors. The training set is shown on Figure 4.1.

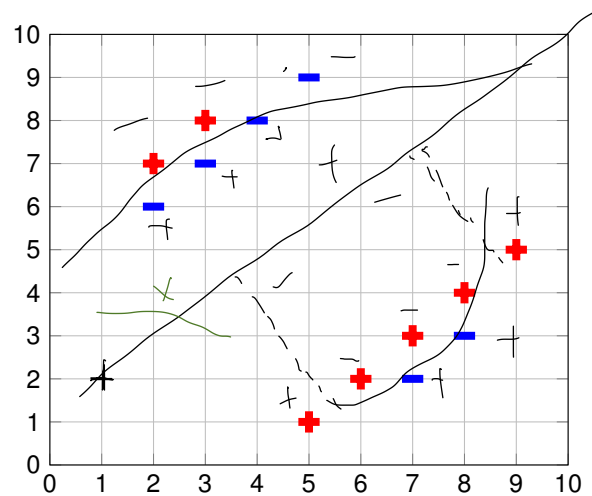


Figure 4.1: A red $+$ denotes instances from class 1, and a blue $-$ denotes instances from class 2.

a) Specify one value of k that minimizes the leave-one-out cross-validation (LOOCV) error. Please consider only *odd* values of k (e.g. 1, 3, 5, 7, ...) to avoid ties. What is the resulting error, i.e. the number of misclassified data points?

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

b) Imagine that the training dataset contains one additional point with coordinates (1, 2) and label $+$. Would this change the decision boundary for a 1-nearest neighbor classifier? Why or why not?

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

c) If your answer above was *no* what is the *shortest* distance that you need to move (1, 2) such that it changes the decision boundary? If your answer above was *yes* what is the *shortest* distance that you need to move (1, 2) so it does not change the decision boundary?

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

Write down the new coordinates after moving and specify the shortest distance.

必需把点放回刻 + 上

Problem 2: k-nearest neighbors (Version B) (4 credits)

Consider a k -nearest neighbor classifier using Euclidean distance on a binary classification task. The prediction for an instance is the *majority* class of its k -nearest neighbors. The training set is shown on Figure 5.1.

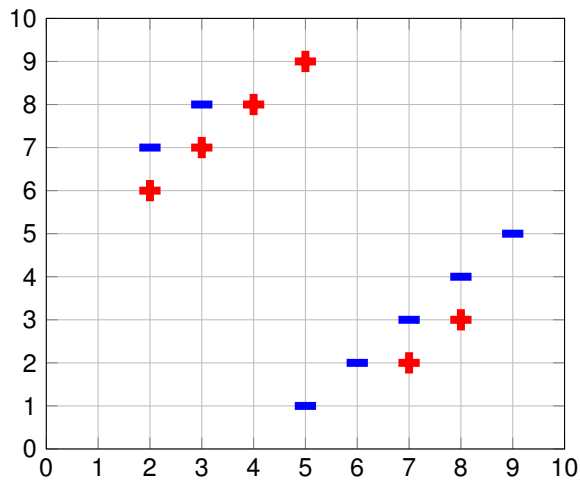


Figure 5.1: A red + denotes instances from class 1, and a blue - denotes instances from class 2.

- 0 ☐
- 1 ☐
- 2 ☐
- a) Specify one value of k that minimizes the leave-one-out cross-validation (LOOCV) error. Please consider only *odd* values of k (e.g. 1, 3, 5, 7, ...) to avoid ties. What is the resulting error, i.e. the number of misclassified data points?

- 0 ☐
- 1 ☐
- b) Imagine that the training dataset contains one additional point with coordinates (1, 2) and label +. Would this change the decision boundary for a 1-nearest neighbor classifier? Why or why not?

- 0 ☐
- 1 ☐
- c) If your answer above was *no* what is the *shortest* distance that you need to move (1, 2) such that it changes the decision boundary? If your answer above was *yes* what is the *shortest* distance that you need to move (1, 2) so it does not change the decision boundary?

Write down the new coordinates after moving and specify the shortest distance.

Problem 2: k-nearest neighbors (Version C) (4 credits)

Consider a k -nearest neighbor classifier using Euclidean distance on a binary classification task. The prediction for an instance is the *majority* class of its k -nearest neighbors. The training set is shown on Figure 6.1.

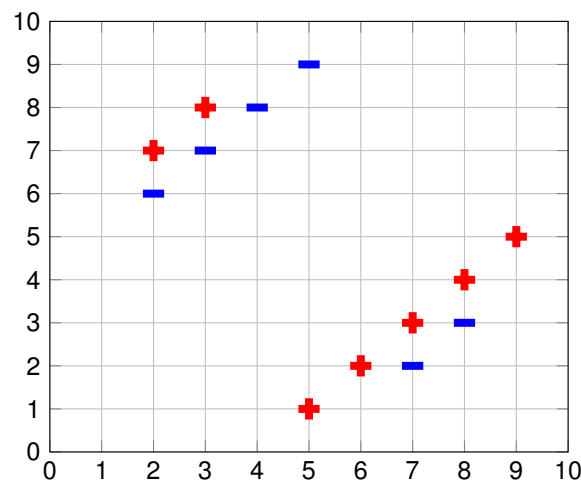


Figure 6.1: A red + denotes instances from class 1, and a blue - denotes instances from class 2.

a) Specify one value of k that minimizes the leave-one-out cross-validation (LOOCV) error. Please consider only *odd* values of k (e.g. 1, 3, 5, 7, ...) to avoid ties. What is the resulting error, i.e. the number of misclassified data points?

 0
 1
 2

b) Imagine that the training dataset contains one additional point with coordinates (1, 2) and label -. Would this change the decision boundary for a 1-nearest neighbor classifier? Why or why not?

 0
 1

c) If your answer above was *no* what is the *shortest* distance that you need to move (1, 2) such that it changes the decision boundary? If your answer above was *yes* what is the *shortest* distance that you need to move (1, 2) so it does not change the decision boundary?

 0
 1

Write down the new coordinates after moving and specify the shortest distance.

Problem 2: k-nearest neighbors (Version D) (4 credits)

Consider a k -nearest neighbor classifier using Euclidean distance on a binary classification task. The prediction for an instance is the *majority* class of its k -nearest neighbors. The training set is shown on Figure 7.1.

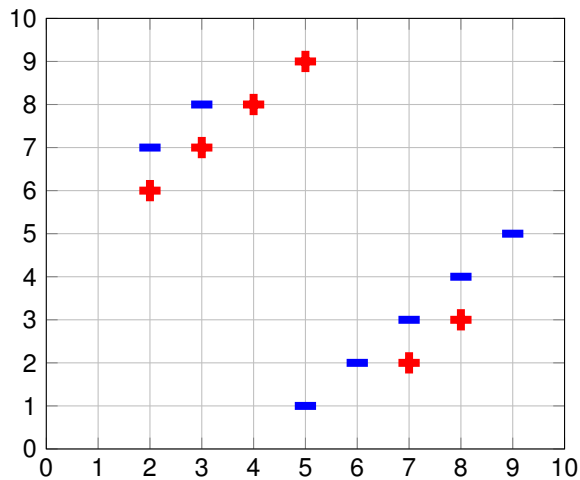


Figure 7.1: A red + denotes instances from class 1, and a blue - denotes instances from class 2.

- 0 ☐ a) Specify one value of k that minimizes the leave-one-out cross-validation (LOOCV) error. Please consider only *odd* values of k (e.g. 1, 3, 5, 7, ...) to avoid ties. What is the resulting error, i.e. the number of misclassified data points?
- 1 ☐
- 2 ☐

- 0 ☐ b) Imagine that the training dataset contains one additional point with coordinates (1, 2) and label -. Would this change the decision boundary for a 1-nearest neighbor classifier? Why or why not?
- 1 ☐

- 0 ☐ c) If your answer above was *no* what is the *shortest* distance that you need to move (1, 2) such that it changes the decision boundary? If your answer above was *yes* what is the *shortest* distance that you need to move (1, 2) so it does not change the decision boundary?
- 1 ☐

Write down the new coordinates after moving and specify the shortest distance.

Problem 3: Linear Regression (Version A) (6 credits)

$$E(x^{(i)}) = 0$$

Consider a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, with $\mathbf{x}^{(i)} \in \mathbb{R}^D$, $y^{(i)} \in \mathbb{R}$ and centered features so that $\frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} = \mathbf{0}$. During preprocessing, we absorb the bias as a constant feature to produce the transformed dataset $\tilde{\mathcal{D}} = \{(\tilde{\mathbf{x}}^{(i)}, \tilde{y}^{(i)})\}_{i=1}^N$ where we map each $(\mathbf{x}^{(i)}, y^{(i)})$ to

$$\tilde{\mathbf{x}}^{(i)} = \begin{pmatrix} \mathbf{x}^{(i)} \\ 1 \end{pmatrix} \in \mathbb{R}^{D+1} \quad \text{and} \quad \tilde{y}^{(i)} = y^{(i)}.$$

$$\tilde{\mathbf{w}}^T = (w_1 \dots w_{D+1})$$

$$\dots w^T \cdot x^{(i)} + w_{D+1}$$

We want to perform ridge regression with regularization strength λ to find the optimal weight vector $\tilde{\mathbf{w}}^* \in \mathbb{R}^{D+1}$.

Reminder: In ridge regression we optimize the loss function

$$\mathcal{L}(\tilde{\mathbf{w}}) = \frac{1}{2} \sum_{i=1}^N (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^{(i)} - y^{(i)})^2 + \frac{\lambda}{2} \|\tilde{\mathbf{w}}\|_2^2.$$

$$\tilde{\mathbf{x}}^{(i)} = \begin{pmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_D^{(i)} \\ 1 \end{pmatrix} \quad \tilde{\mathbf{w}} = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_{D+1} \end{pmatrix}$$

a) Derive the closed form solution for the last element of the weight vector $\tilde{\mathbf{w}}_{D+1}^*$ corresponding to the absorbed bias obtained from ridge regression on $\tilde{\mathcal{D}}$.

$$\mathcal{L}(\tilde{\mathbf{w}}) = \frac{1}{2} \sum_{i=1}^N (w^T x^{(i)} + \tilde{w}_{D+1} - y^{(i)})^2 + \frac{\lambda}{2} [w_1^2 + \dots + \tilde{w}_{D+1}^2]$$

$$\frac{\partial \mathcal{L}(\tilde{\mathbf{w}})}{\partial \tilde{w}_{D+1}} = \sum_{i=1}^N (w^T x^{(i)} + \tilde{w}_{D+1} - y^{(i)}) + \lambda \tilde{w}_{D+1} \stackrel{!}{=} 0$$

$$w^T \sum_{i=1}^N x^{(i)} + \sum_{i=1}^N \tilde{w}_{D+1} - \sum_{i=1}^N y^{(i)} + \lambda \tilde{w}_{D+1} \stackrel{!}{=} 0$$

$$N \tilde{w}_{D+1} - \sum_{i=1}^N y_i + \lambda \tilde{w}_{D+1} \stackrel{!}{=} 0$$

$$\tilde{w}_{D+1} (N + \lambda) = \sum_{i=1}^N y_i$$

$$\tilde{w}_{D+1}^* = \frac{1}{N + \lambda} \sum_{i=1}^N y_i$$

0
1
2
3

b) Propose an alternative *preprocessing step*, i.e. define an alternative transformed dataset $\hat{\mathcal{D}} = \{(\hat{\mathbf{x}}^{(i)}, \hat{y}^{(i)})\}_{i=1}^N$, such that the optimal ridge regression vector $\hat{\mathbf{w}}^* \in \mathbb{R}^D$ on $\hat{\mathcal{D}}$ is equivalent to $\tilde{\mathbf{w}}^*$ obtained on $\tilde{\mathcal{D}}$. Justify that ridge regression on your preprocessed data $\hat{\mathcal{D}}$ finds an optimal $\hat{\mathbf{w}}^* \in \mathbb{R}^D$ that is equal to $\tilde{\mathbf{w}}^*$ in the first D elements, i.e. $\hat{\mathbf{w}}_i^* = \tilde{\mathbf{w}}_i^*$ for $i \in \{1, \dots, D\}$. You do not need to derive the closed-form solution for $\hat{\mathbf{w}}^*$.

$$\mathcal{L}(\hat{\mathbf{w}}) = \frac{1}{2} \sum_{i=1}^N \left(\hat{\mathbf{w}}^T \hat{\mathbf{x}}^{(i)} + \frac{1}{N + \lambda} \sum_{i=1}^N y_i - y^{(i)} \right)^2 + \frac{\lambda}{2} \|\hat{\mathbf{w}}\|_2^2 + \frac{\lambda}{2} \left(\frac{1}{N + \lambda} \sum_{i=1}^N y_i \right)^2$$

derivative is 0 after nothing

$$\mathcal{L}(\hat{\mathbf{w}}) = \frac{1}{2} \sum_{i=1}^N \left(\hat{\mathbf{w}}^T \hat{\mathbf{x}}^{(i)} + \left(\frac{1}{N + \lambda} \sum_{i=1}^N y_i \right) - y^{(i)} \right)^2 + \frac{\lambda}{2} \|\hat{\mathbf{w}}\|_2^2$$

so $\hat{\mathbf{w}}$ is equal to $\tilde{\mathbf{w}}^*$ if y_i are shifted by $\frac{1}{N + \lambda} \sum_{i=1}^N y_i$

0
1
2
3

Problem 3: Linear Regression (Version B) (6 credits)

Consider a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, with $\mathbf{x}^{(i)} \in \mathbb{R}^D$, $y^{(i)} \in \mathbb{R}$ and centered features so that $\frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} = \mathbf{0}$. During preprocessing, we absorb the bias as a constant feature to produce the transformed dataset $\tilde{\mathcal{D}} = \{(\tilde{\mathbf{x}}^{(i)}, \tilde{y}^{(i)})\}_{i=1}^N$ where we map each $(\mathbf{x}^{(i)}, y^{(i)})$ to

$$\tilde{\mathbf{x}}^{(i)} = \begin{pmatrix} \mathbf{x}^{(i)} \\ 1 \end{pmatrix} \in \mathbb{R}^{D+1} \quad \text{and} \quad \tilde{y}^{(i)} = y^{(i)}.$$

We want to perform ridge regression with regularization strength λ to find the optimal weight vector $\tilde{\mathbf{w}}^* \in \mathbb{R}^{D+1}$.

Reminder: In ridge regression we optimize the loss function

$$\mathcal{L}(\tilde{\mathbf{w}}) = \frac{1}{2} \sum_{i=1}^N (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}^{(i)} - y^{(i)})^2 + \frac{\lambda}{2} \|\tilde{\mathbf{w}}\|_2^2.$$

- 0 ☐ a) Derive the closed form solution for the last element of the weight vector $\tilde{\mathbf{w}}_{D+1}^*$ corresponding to the absorbed bias obtained from ridge regression on $\tilde{\mathcal{D}}$.

1 ☐

2 ☐

3 ☐

- 0 ☐ b) Propose an alternative *preprocessing step*, i.e. define an alternative transformed dataset $\hat{\mathcal{D}} = \{(\hat{\mathbf{x}}^{(i)}, \hat{y}^{(i)})\}_{i=1}^N$, such that the optimal ridge regression vector $\hat{\mathbf{w}}^* \in \mathbb{R}^D$ on $\hat{\mathcal{D}}$ is equivalent to $\tilde{\mathbf{w}}^*$ obtained on $\tilde{\mathcal{D}}$. Justify that ridge regression on your preprocessed data $\hat{\mathcal{D}}$ finds an optimal $\hat{\mathbf{w}}^* \in \mathbb{R}^D$ that is equal to $\tilde{\mathbf{w}}^*$ in the first D elements, i.e. $\hat{\mathbf{w}}_i^* = \tilde{\mathbf{w}}_i^*$ for $i \in \{1, \dots, D\}$. You do not need to derive the closed-form solution for $\hat{\mathbf{w}}^*$.

1 ☐

2 ☐

3 ☐

Problem 4: Naive Bayes (Version A) (6 credits)

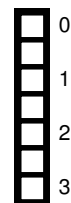
We have collected the dataset shown in the table below with the continuous feature x_1 , and the discrete feature x_2 that can take either of the values yes or no. Each data point is labeled as one of three classes 1, 2 or 3 denoted by y .

Table 10.1: Naive Bayes Data (each column is one data point)

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
x_1	0.0	2.0	-2.0	-1.0	3.0	4.0	6.0
x_2	yes	no	no	yes	no	yes	yes
y	1	1	2	2	2	3	3

a) Set up a naive Bayes classifier (choose likelihoods and their parameterization) for the data in Table 10.1 using Normal distributions with fixed variance of 1 and Bernoulli distributions. Compute the maximum likelihood estimate of all parameters θ required for naive Bayes.

$$\begin{aligned}
 x_1 | y=c &\sim \mathcal{N}(\mu_c, 1) & \mu_1 &= 1 & \mu_2 &= 0 & \mu_3 &= 5 \\
 x_2 | y=c, \theta &\sim \text{Ber}(\theta) & \theta_1 &= \frac{1}{2} & \theta_2 &= \frac{1}{3} & \theta_3 &= 1 \\
 \pi_1 &= \frac{3}{7} & \pi_2 &= \frac{3}{7} & \pi_3 &= \frac{2}{7}
 \end{aligned}$$



b) You observe a new data point $\mathbf{x}^{(b)} = (1 \text{ yes})$. Compute the unnormalized posterior over classes $p(y^{(b)} | \mathbf{x}^{(b)}, \theta)$ for $\mathbf{x}^{(b)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?

$$\begin{aligned}
 p(y^{(b)} | \mathbf{x}^{(b)}, \theta) &\propto p(x_1^{(b)} | y^{(b)}, \theta) \cdot p(y^{(b)} | \theta) \\
 &\propto p(x_1^{(b)} | y^{(b)}, \theta) \cdot p(x_2^{(b)} | y^{(b)}, \theta) \cdot p(y^{(b)} | \theta) \\
 y=1 &\propto \exp(-\frac{1}{2}(0)^2) \cdot \frac{1}{2} \cdot \frac{3}{7} = \frac{1}{7} \\
 y=2 &\propto \exp(-\frac{1}{2}(1)^2) \cdot \frac{1}{3} \cdot \frac{3}{7} = \frac{1}{7} \cdot e^{-\frac{1}{2}} \\
 y=3 &\propto \exp(-\frac{1}{2}(1)^2) \cdot 1 \cdot \frac{2}{7} = \frac{2}{7} e^{-\frac{1}{2}}
 \end{aligned}$$

y=1.



c) Next, you get another data point $\mathbf{x}^{(d)} = (n/a \text{ n/a})$, where all values are missing. Compute the unnormalized posterior over classes $p(y^{(d)} | \mathbf{x}^{(d)}, \theta)$ for $\mathbf{x}^{(d)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?

know nothing about data point posterior distribution is prior distribution



d) Finally, you see a data point $\mathbf{x}^{(c)} = (n/a \text{ no})$, i.e. the features are only partially known. Compute the unnormalized posterior over classes $p(y^{(c)} | \mathbf{x}^{(c)}, \theta)$ for $\mathbf{x}^{(c)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?

$$\begin{aligned}
 y=1 &\propto \exp(-\frac{1}{2}(x_1-1)^2) \cdot \frac{1}{2} \cdot \frac{2}{7} = \frac{1}{7} \\
 y=2 &\propto \exp(-\frac{1}{2}x_1^2) \cdot \frac{1}{3} \cdot \frac{3}{7} = \frac{2}{7} \\
 y=3 &\propto 0
 \end{aligned}$$



Problem 4: Naive Bayes (Version B) (6 credits)

We have collected the dataset shown in the table below with the continuous feature x_1 , and the discrete feature x_2 that can take either of the values yes or no. Each data point is labeled as one of three classes 1, 2 or 3 denoted by y .

Table 11.1: Naive Bayes Data (each column is one data point)

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
x_1	0.0	2.0	-2.0	-1.0	3.0	4.0	6.0
x_2	yes	no	no	yes	no	yes	yes
y	1	1	2	2	2	3	3

- 0 ☐ a) Set up a naive Bayes classifier (choose likelihoods and their parameterization) for the data in Table 11.1 using Normal distributions with fixed variance of 1 and Bernoulli distributions. Compute the maximum likelihood estimate of *all* parameters θ required for naive Bayes.
- 1 ☐
- 2 ☐
- 3 ☐

- 0 ☐ b) You observe a new data point $\mathbf{x}^{(b)} = (2 \text{ yes})$. Compute the unnormalized posterior over classes $p(y^{(b)} | \mathbf{x}^{(b)}, \theta)$ for $\mathbf{x}^{(b)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?
- 1 ☐

- 0 ☐ c) Next, you get another data point $\mathbf{x}^{(d)} = (\text{n/a} \text{ n/a})$, where all values are missing. Compute the unnormalized posterior over classes $p(y^{(d)} | \mathbf{x}^{(d)}, \theta)$ for $\mathbf{x}^{(d)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?
- 1 ☐

- 0 ☐ d) Finally, you see a data point $\mathbf{x}^{(c)} = (\text{n/a} \text{ no})$, i.e. the features are only partially known. Compute the unnormalized posterior over classes $p(y^{(c)} | \mathbf{x}^{(c)}, \theta)$ for $\mathbf{x}^{(c)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?
- 1 ☐

Problem 4: Naive Bayes (Version C) (6 credits)

We have collected the dataset shown in the table below with the continuous feature x_1 , and the discrete feature x_2 that can take either of the values yes or no. Each data point is labeled as one of three classes 1, 2 or 3 denoted by y .

Table 12.1: Naive Bayes Data (each column is one data point)

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
x_1	-3.0	-1.0	2.0	2.0	3.0	3.0	6.0
x_2	yes	no	no	no	no	yes	yes
y	1	1	2	2	3	3	3

a) Set up a naive Bayes classifier (choose likelihoods and their parameterization) for the data in Table 12.1 using Normal distributions with fixed variance of 1 and Bernoulli distributions. Compute the maximum likelihood estimate of *all* parameters θ required for naive Bayes.

☐ 0
☐ 1
☐ 2
☐ 3

b) You observe a new data point $\mathbf{x}^{(b)} = (\text{yes} \quad 1)$. Compute the unnormalized posterior over classes $p(y^{(b)} | \mathbf{x}^{(b)}, \theta)$ for $\mathbf{x}^{(b)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?

☐ 0
☐ 1

c) Next, you get another data point $\mathbf{x}^{(d)} = (\text{n/a} \quad \text{n/a})$, where all values are missing. Compute the unnormalized posterior over classes $p(y^{(d)} | \mathbf{x}^{(d)}, \theta)$ for $\mathbf{x}^{(d)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?

☐ 0
☐ 1

d) Finally, you see a data point $\mathbf{x}^{(c)} = (\text{n/a} \quad \text{no})$, i.e. the features are only partially known. Compute the unnormalized posterior over classes $p(y^{(c)} | \mathbf{x}^{(c)}, \theta)$ for $\mathbf{x}^{(c)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?

☐ 0
☐ 1

Problem 4: Naive Bayes (Version D) (6 credits)

We have collected the dataset shown in the table below with the continuous feature x_1 , and the discrete feature x_2 that can take either of the values yes or no. Each data point is labeled as one of three classes 1, 2 or 3 denoted by y .

Table 13.1: Naive Bayes Data (each column is one data point)

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
x_1	-3.0	-1.0	2.0	2.0	3.0	3.0	6.0
x_2	yes	no	no	no	no	yes	yes
y	1	1	2	2	3	3	3

- 0 ☐ a) Set up a naive Bayes classifier (choose likelihoods and their parameterization) for the data in Table 13.1 using Normal distributions with fixed variance of 1 and Bernoulli distributions. Compute the maximum likelihood estimate of *all* parameters θ required for naive Bayes.
- 1 ☐
- 2 ☐
- 3 ☐

- 0 ☐ b) You observe a new data point $\mathbf{x}^{(b)} = (\text{yes} \quad 2)$. Compute the unnormalized posterior over classes $p(y^{(b)} | \mathbf{x}^{(b)}, \theta)$ for $\mathbf{x}^{(b)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?
- 1 ☐

- 0 ☐ c) Next, you get another data point $\mathbf{x}^{(d)} = (\text{n/a} \quad \text{n/a})$, where all values are missing. Compute the unnormalized posterior over classes $p(y^{(d)} | \mathbf{x}^{(d)}, \theta)$ for $\mathbf{x}^{(d)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?
- 1 ☐

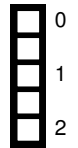
- 0 ☐ d) Finally, you see a data point $\mathbf{x}^{(c)} = (\text{n/a} \quad \text{no})$, i.e. the features are only partially known. Compute the unnormalized posterior over classes $p(y^{(c)} | \mathbf{x}^{(c)}, \theta)$ for $\mathbf{x}^{(c)}$. Simplify as far as possible without evaluating exponential functions, square roots, logarithms, etc. Briefly justify how you arrive at your solution. What is the most likely class for this data point?
- 1 ☐

Problem 5: Optimization (Version A) (2 credits)

Suppose that $\mathbf{a} \in \mathbb{R}^d$ is some fixed vector. We define the function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ as

$$f(\mathbf{x}) = \exp(e^{\mathbf{a}^T \mathbf{x}} + e^{-\mathbf{a}^T \mathbf{x}}).$$

Prove or disprove that f is convex on \mathbb{R}^d .



e^x is convex due to $(e^x)'' = e^x > 0$

let $H(x) = e^x$ $(\mathbf{a}^T \mathbf{x})' = (\mathbf{a}^T)' = 0$.

$\mathbf{a}^T \mathbf{x}$ and $-\mathbf{a}^T \mathbf{x}$ is convex in \mathbf{x} .

rule 5 $H(\mathbf{a}^T \mathbf{x})$ and $H(-\mathbf{a}^T \mathbf{x})$ is convex

rule 1 $H(\mathbf{a}^T \mathbf{x}) + H(-\mathbf{a}^T \mathbf{x})$ is convex

rule 5 $H(H(\mathbf{a}^T \mathbf{x}) + H(-\mathbf{a}^T \mathbf{x}))$ is convex.

Problem 5: Optimization (Version B) (2 credits)

- 0 ☐ Suppose that $\mathbf{a} \in \mathbb{R}^d$ is some fixed vector. We define the function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ as
- 1 ☐
$$f(\mathbf{x}) = \exp(e^{\mathbf{a}^T \mathbf{x}} + e^{-\mathbf{a}^T \mathbf{x}}).$$
- 2 ☐ Prove or disprove that f is convex on \mathbb{R}^d .

Problem 5: Optimization (Version C) (2 credits)

Suppose that $\mathbf{a} \in \mathbb{R}^d$ is some fixed vector. We define the function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ as

$$f(\mathbf{x}) = \exp(e^{\mathbf{a}^\top \mathbf{x}} + e^{-\mathbf{a}^\top \mathbf{x}}).$$

Prove or disprove that f is convex on \mathbb{R}^d .

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

Problem 5: Optimization (Version D) (2 credits)

- 0 ☐ Suppose that $\mathbf{a} \in \mathbb{R}^d$ is some fixed vector. We define the function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ as
- 1 ☐
$$f(\mathbf{x}) = \exp(e^{\mathbf{a}^T \mathbf{x}} + e^{-\mathbf{a}^T \mathbf{x}}).$$
- 2 ☐ Prove or disprove that f is convex on \mathbb{R}^d .

Problem 6: Convolutional Neural Networks (Version A) (3 credits)

Recall that a convolutional layer is defined by the following parameters:

- C_{in} — number of input channels
- C_{out} — number of output channels
- K — kernel (sliding window) size
- P — padding size
- S — stride

Suppose that \mathbf{x} is an image with height 64 and width 32 (in pixels) that has 3 channels, which we can represent as a tensor (three dimensional array) of shape $[3, 64, 32]$.

We passed \mathbf{x} through a neural network consisting of two convolutional layers conv1 and conv2 (in this order). As output we obtained a tensor of shape $[16, 16, 8]$ (i.e., 16 channels, height 16 and width 8).

We know that the first convolutional layer conv1 has the following parameters

- $C_{in} = 3$
- $C_{out} = 32$
- $K = 3$
- $P = 1$
- $S = 4$

We also know that the second convolutional layer conv2 has kernel size $K = 3$.

What are the remaining parameters C_{in} , C_{out} , P , S of the second convolutional layer conv2? Justify your answer.

$$\left\lfloor \frac{N + 2P - K}{S} \right\rfloor + 1$$

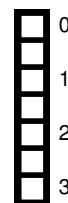
$$\left\lfloor \frac{64 + 2 \times 1 - 3}{4} \right\rfloor + 1 = 16$$

$$\left\lfloor \frac{32 + 2 \times 1 - 3}{4} \right\rfloor + 1 = 8$$

$$[32, 16, 8]$$

$$\left\lfloor \frac{16 + 2 \times P - 3}{S} \right\rfloor + 1 = 16$$

$$\left\lfloor \frac{8 + 2P - 3}{S} \right\rfloor + 1 = 8$$



Problem 6: Convolutional Neural Networks (Version B) (3 credits)

Recall that a convolutional layer is defined by the following parameters:

- C_{in} — number of input channels
- C_{out} — number of output channels
- K — kernel (sliding window) size
- P — padding size
- S — stride

Suppose that \mathbf{x} is an image with height 64 and width 32 (in pixels) that has 3 channels, which we can represent as a tensor (three dimensional array) of shape $[3, 64, 32]$.

We passed \mathbf{x} through a neural network consisting of two convolutional layers `conv1` and `conv2` (in this order). As output we obtained a tensor of shape $[16, 16, 8]$ (i.e., 16 channels, height 16 and width 8).

We know that the first convolutional layer `conv1` has the following parameters

- $C_{\text{in}} = 3$
- $C_{\text{out}} = 32$
- $K = 3$
- $P = 1$
- $S = 1$

We also know that the second convolutional layer `conv2` has kernel size $K = 3$.

0 ☐ What are the remaining parameters $C_{\text{in}}, C_{\text{out}}, P, S$ of the second convolutional layer `conv2`? Justify your answer.

1 ☐

2 ☐

3 ☐

Problem 6: Convolutional Neural Networks (Version C) (3 credits)

Recall that a convolutional layer is defined by the following parameters:

- C_{in} — number of input channels
- C_{out} — number of output channels
- K — kernel (sliding window) size
- P — padding size
- S — stride

Suppose that \mathbf{x} is an image with height 64 and width 32 (in pixels) that has 3 channels, which we can represent as a tensor (three dimensional array) of shape $[3, 64, 32]$.

We passed \mathbf{x} through a neural network consisting of two convolutional layers `conv1` and `conv2` (in this order). As output we obtained a tensor of shape $[16, 16, 8]$ (i.e., 16 channels, height 16 and width 8).

We know that the first convolutional layer `conv1` has the following parameters

- $C_{\text{in}} = 3$
- $C_{\text{out}} = 8$
- $K = 3$
- $P = 1$
- $S = 4$

We also know that the second convolutional layer `conv2` has kernel size $K = 3$.

What are the remaining parameters $C_{\text{in}}, C_{\text{out}}, P, S$ of the second convolutional layer `conv2`? Justify your answer.

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

Problem 6: Convolutional Neural Networks (Version D) (3 credits)

Recall that a convolutional layer is defined by the following parameters:

- C_{in} — number of input channels
- C_{out} — number of output channels
- K — kernel (sliding window) size
- P — padding size
- S — stride

Suppose that \mathbf{x} is an image with height 64 and width 32 (in pixels) that has 3 channels, which we can represent as a tensor (three dimensional array) of shape $[3, 64, 32]$.

We passed \mathbf{x} through a neural network consisting of two convolutional layers `conv1` and `conv2` (in this order). As output we obtained a tensor of shape $[16, 16, 8]$ (i.e., 16 channels, height 16 and width 8).

We know that the first convolutional layer `conv1` has the following parameters

- $C_{\text{in}} = 3$
- $C_{\text{out}} = 8$
- $K = 3$
- $P = 1$
- $S = 1$

We also know that the second convolutional layer `conv2` has kernel size $K = 3$.

0 ☐ What are the remaining parameters $C_{\text{in}}, C_{\text{out}}, P, S$ of the second convolutional layer `conv2`? Justify your answer.

1 ☐

2 ☐

3 ☐

Problem 7: SVMs (Version A) (5 credits)

Consider a soft-margin SVM trained on a binary classification task with some fixed and finite penalty C , i.e. $0 < C < \infty$. Assume that the training set contains at least two instances from each class. After training you observe that the optimal slack variables are zero for all instances except for a single instance q . Specifically, $\xi_q > 2$ for q and $\xi_i = 0$ for all $i \neq q$. Let $m_{\text{soft}} = \frac{2}{\|w_{\text{soft}}\|}$ be the optimal margin where w_{soft} are the optimal weights.

a) Now you *remove* the instance q from the dataset and you train a new *hard-margin* SVM. Let m_{hard} be the resulting optimal margin. Which one of the following holds: $m_{\text{hard}} \leq m_{\text{soft}}$, $m_{\text{hard}} = m_{\text{soft}}$, $m_{\text{hard}} \geq m_{\text{soft}}$. Justify your answer.

only q point violate the margin, all other are properly classified

If removed q , w_{soft} has already fit the hard sum, so it could be better or equal

$$m_{\text{hard}} \geq m_{\text{soft}}$$

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

b) Now instead of removing it, you *relabel* the instance q (if before it was class -1 you set it to class $+1$ and vice versa) and you train a new *hard-margin* SVM. Let m_{hard} be the resulting optimal margin. Which one of the following holds: $m_{\text{hard}} \leq m_{\text{soft}}$, $m_{\text{hard}} = m_{\text{soft}}$, $m_{\text{hard}} \geq m_{\text{soft}}$. Justify your answer.

$$\xi_i = 1 - y_q (w^T x_q + b) > 2$$

$$y_q (w^T x_q + b) < -1$$

$$\text{Set } y_q = -y_q$$

$$y_q (w^T x_q + b) > 1 \Leftarrow \text{correct classify.}$$

so w is fit hard sum

and it can at least equal or better

$$m_{\text{hard}} \geq m_{\text{soft}}$$

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

Problem 7: SVMs (Version B) (5 credits)

Consider a soft-margin SVM trained on a binary classification task with some fixed and finite penalty C , i.e. $0 < C < \infty$. Assume that the training set contains at least two instances from each class. After training you observe that the optimal slack variables are zero for all instances except for a *single* instance q . Specifically, $\xi_q > 2$ for q and $\xi_i = 0$ for all $i \neq q$. Let $m_{\text{soft}} = \frac{2}{\|\mathbf{w}_{\text{soft}}\|}$ be the optimal margin where \mathbf{w}_{soft} are the optimal weights.

0 ☐ a) Now you *remove* the instance q from the dataset and you train a new *hard-margin* SVM. Let m_{hard} be the resulting optimal margin. Which one of the following holds: $m_{\text{hard}} \leq m_{\text{soft}}$, $m_{\text{hard}} = m_{\text{soft}}$, $m_{\text{hard}} \geq m_{\text{soft}}$. Justify your answer.

1 ☐

2 ☐

3 ☐

0 ☐ b) Now instead of removing it, you *relabel* the instance q (if before it was class -1 you set it to class $+1$ and vice versa) and you train a new *hard-margin* SVM. Let m_{hard} be the resulting optimal margin. Which one of the following holds: $m_{\text{hard}} \leq m_{\text{soft}}$, $m_{\text{hard}} = m_{\text{soft}}$, $m_{\text{hard}} \geq m_{\text{soft}}$. Justify your answer.

1 ☐

2 ☐

Problem 8: Low-rank Approximation and Regression (Version A) (6 credits)

Consider a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, with $\mathbf{x}^{(i)} \in \mathbb{R}^D$, $y^{(i)} \in \mathbb{R}$. Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the matrix of features and $\mathbf{y} \in \mathbb{R}^N$ be the vector of regression targets.

We construct the matrix $\mathbf{M} = [\mathbf{X}, \mathbf{y}] \in \mathbb{R}^{N \times (D+1)}$ where \mathbf{y} has been appended as an additional column to \mathbf{X} . Let $\mathbf{M}' \in \mathbb{R}^{N \times (D+1)}$ be the best rank K approximation of \mathbf{M} for some chosen constant K . Define $\mathbf{M}' = [\mathbf{X}', \mathbf{y}']$ where \mathbf{X}' is the matrix containing the first D columns of \mathbf{M}' and \mathbf{y}' is the last column of \mathbf{M}' .

We will fit a standard linear regression model using \mathbf{X}' as the feature matrix and \mathbf{y}' as the target, i.e. we find the optimal weight vector $\mathbf{w}^* \in \mathbb{R}^D$ and the bias $b^* \in \mathbb{R}$. *Careful: We do regression on \mathbf{X}' and \mathbf{y}' not on \mathbf{X} and \mathbf{y} !*

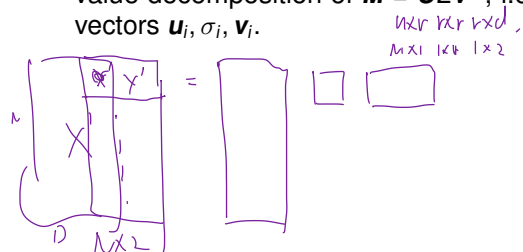
a) Consider the case where $D = 1$ and $K = 1$, i.e. our features are one dimensional and \mathbf{M}' is the the best rank 1 approximation of \mathbf{M} . Assuming the features \mathbf{X}' contain at least two distinct elements, what is the *training error* achieved by the optimal \mathbf{w}^* and b^* ? Justify your answer.

Only rank 1, all rows are linearly dependent to one row

so it can be perfectly expressed \Rightarrow error is equal to zero

0
1
2

b) For the same case as in (a) where $D = 1$ and $K = 1$, write down the optimal \mathbf{w}^* and b^* in terms of the singular value decomposition of $\mathbf{M} = \mathbf{U}\Sigma\mathbf{V}^T$, i.e. write down \mathbf{w}^* and b^* as a function of the singular values and singular vectors $\mathbf{u}_i, \sigma_i, \mathbf{v}_i$.



$$\mathbf{M}' = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T$$

$$\mathbf{X}' = \sigma_1 \mathbf{u}_1 \mathbf{v}_{11}$$

$$\mathbf{y}' = \sigma_1 \mathbf{u}_1 \mathbf{v}_{12}$$

$$\mathbf{w}^* \mathbf{X}' + b^* = \mathbf{y}'$$

$$\mathbf{w}^* \sigma_1 \mathbf{u}_1 \mathbf{v}_{11} + b^* = \sigma_1 \mathbf{u}_1 \mathbf{v}_{12}$$

\nwarrow linearly dependent.

$$\mathbf{w}^* = \frac{\mathbf{v}_{12}}{\mathbf{v}_{11}}$$

0
1
2
3

c) Consider the general case of D -dimensional ($D > 1$) feature vectors and a rank K approximation of \mathbf{M} . Assuming that \mathbf{X}' is full rank, what is the *training error* achieved by the optimal \mathbf{w}^* and b^* as a function of K ?

$$K = D \text{ or } K = D + 1$$

If $K = D$ \mathbf{y} is linear combination of \mathbf{X}

If $K = D + 1$ depends on dataset, error ≥ 0 .

0
1

Problem 8: Low-rank Approximation and Regression (Version B) (6 credits)

Consider a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, with $\mathbf{x}^{(i)} \in \mathbb{R}^D, y^{(i)} \in \mathbb{R}$. Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the matrix of features and $\mathbf{y} \in \mathbb{R}^N$ be the vector of regression targets.

We construct the matrix $\mathbf{M} = [\mathbf{X}, \mathbf{y}] \in \mathbb{R}^{N \times (D+1)}$ where \mathbf{y} has been appended as an additional column to \mathbf{X} . Let $\mathbf{M}' \in \mathbb{R}^{N \times (D+1)}$ be the *best* rank K approximation of \mathbf{M} for some chosen constant K . Define $\mathbf{M}' = [\mathbf{X}', \mathbf{y}']$ where \mathbf{X}' is the matrix containing the first D columns of \mathbf{M}' and \mathbf{y}' is the last column of \mathbf{M}' .

We will fit a standard linear regression model using \mathbf{X}' as the feature matrix and \mathbf{y}' as the target, i.e. we find the optimal weight vector $\mathbf{w}^* \in \mathbb{R}^D$ and the bias $b^* \in \mathbb{R}$. *Careful: We do regression on \mathbf{X}' and \mathbf{y}' not on \mathbf{X} and \mathbf{y} !*

- 0 ☐
- 1 ☐
- 2 ☐
- a) Consider the case where $D = 1$ and $K = 1$, i.e. our features are one dimensional and \mathbf{M}' is the the best rank 1 approximation of \mathbf{M} . Assuming the features \mathbf{X}' contain at least two distinct elements, what is the *training error* achieved by the optimal \mathbf{w}^* and b^* ? Justify your answer.

- 0 ☐
- 1 ☐
- 2 ☐
- 3 ☐
- b) For the same case as in (a) where $D = 1$ and $K = 1$, write down the optimal \mathbf{w}^* and b^* in terms of the singular value decomposition of $\mathbf{M} = \mathbf{U}\Sigma\mathbf{V}^T$, i.e. write down \mathbf{w}^* and b^* as a function of the singular values and singular vectors $\mathbf{u}_i, \sigma_i, \mathbf{v}_i$.

- 0 ☐
- 1 ☐
- c) Consider the general case of D -dimensional ($D > 1$) feature vectors and a rank K approximation of \mathbf{M} . Assuming that \mathbf{X}' is full rank, what is the *training error* achieved by the optimal \mathbf{w}^* and b^* as a function of K ?

Problem 8: Low-rank Approximation and Regression (Version C) (6 credits)

Consider a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, with $\mathbf{x}^{(i)} \in \mathbb{R}^D, y^{(i)} \in \mathbb{R}$. Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the matrix of features and $\mathbf{y} \in \mathbb{R}^N$ be the vector of regression targets.

We construct the matrix $\mathbf{M} = [\mathbf{X}, \mathbf{y}] \in \mathbb{R}^{N \times (D+1)}$ where \mathbf{y} has been appended as an additional column to \mathbf{X} . Let $\mathbf{M}' \in \mathbb{R}^{N \times (D+1)}$ be the *best* rank K approximation of \mathbf{M} for some chosen constant K . Define $\mathbf{M}' = [\mathbf{X}', \mathbf{y}']$ where \mathbf{X}' is the matrix containing the first D columns of \mathbf{M}' and \mathbf{y}' is the last column of \mathbf{M}' .

We will fit a standard linear regression model using \mathbf{X}' as the feature matrix and \mathbf{y}' as the target, i.e. we find the optimal weight vector $\mathbf{w}^* \in \mathbb{R}^D$ and the bias $b^* \in \mathbb{R}$. *Careful: We do regression on \mathbf{X}' and \mathbf{y}' not on \mathbf{X} and \mathbf{y} !*

a) Consider the case where $D = 1$ and $K = 1$, i.e. our features are one dimensional and \mathbf{M}' is the the best rank 1 approximation of \mathbf{M} . Assuming the features \mathbf{X}' contain at least two distinct elements, what is the *training error* achieved by the optimal \mathbf{w}^* and b^* ? Justify your answer.

☐ 0
☐ 1
☐ 2

b) For the same case as in (a) where $D = 1$ and $K = 1$, write down the optimal \mathbf{w}^* and b^* in terms of the singular value decomposition of $\mathbf{M} = \mathbf{U}\Sigma\mathbf{V}^T$, i.e. write down \mathbf{w}^* and b^* as a function of the singular values and singular vectors $\mathbf{u}_i, \sigma_i, \mathbf{v}_i$.

☐ 0
☐ 1
☐ 2
☐ 3

c) Consider the general case of D -dimensional ($D > 1$) feature vectors and a rank K approximation of \mathbf{M} . Assuming that \mathbf{X}' is full rank, what is the *training error* achieved by the optimal \mathbf{w}^* and b^* as a function of K ?

☐ 0
☐ 1

Problem 8: Low-rank Approximation and Regression (Version D) (6 credits)

Consider a dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$, with $\mathbf{x}^{(i)} \in \mathbb{R}^D, y^{(i)} \in \mathbb{R}$. Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the matrix of features and $\mathbf{y} \in \mathbb{R}^N$ be the vector of regression targets.

We construct the matrix $\mathbf{M} = [\mathbf{X}, \mathbf{y}] \in \mathbb{R}^{N \times (D+1)}$ where \mathbf{y} has been appended as an additional column to \mathbf{X} . Let $\mathbf{M}' \in \mathbb{R}^{N \times (D+1)}$ be the *best* rank K approximation of \mathbf{M} for some chosen constant K . Define $\mathbf{M}' = [\mathbf{X}', \mathbf{y}']$ where \mathbf{X}' is the matrix containing the first D columns of \mathbf{M}' and \mathbf{y}' is the last column of \mathbf{M}' .

We will fit a standard linear regression model using \mathbf{X}' as the feature matrix and \mathbf{y}' as the target, i.e. we find the optimal weight vector $\mathbf{w}^* \in \mathbb{R}^D$ and the bias $b^* \in \mathbb{R}$. *Careful: We do regression on \mathbf{X}' and \mathbf{y}' not on \mathbf{X} and \mathbf{y} !*

- 0 ☐ a) Consider the case where $D = 1$ and $K = 1$, i.e. our features are one dimensional and \mathbf{M}' is the the best rank 1 approximation of \mathbf{M} . Assuming the features \mathbf{X}' contain at least two distinct elements, what is the *training error* achieved by the optimal \mathbf{w}^* and b^* ? Justify your answer.
- 1 ☐
- 2 ☐

- 0 ☐ b) For the same case as in (a) where $D = 1$ and $K = 1$, write down the optimal \mathbf{w}^* and b^* in terms of the singular value decomposition of $\mathbf{M} = \mathbf{U}\Sigma\mathbf{V}^T$, i.e. write down \mathbf{w}^* and b^* as a function of the singular values and singular vectors $\mathbf{u}_i, \sigma_i, \mathbf{v}_i$.
- 1 ☐
- 2 ☐
- 3 ☐

- 0 ☐ c) Consider the general case of D -dimensional ($D > 1$) feature vectors and a rank K approximation of \mathbf{M} . Assuming that \mathbf{X}' is full rank, what is the *training error* achieved by the optimal \mathbf{w}^* and b^* as a function of K ?
- 1 ☐

Problem 9: K-means (Version A) (6 credits)

Consider a variant of K-means that uses the (squared) Mahalanobis distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \Sigma^{-1} (\mathbf{x}_i - \mathbf{x}_j), \quad (1)$$

with the covariance matrix Σ , instead of the L_2 distance.

Hint: $\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{a}} = (\mathbf{X} + \mathbf{X}^T) \mathbf{a}$

a) Derive the K-means cluster assignment and centroid updates for the (squared) Mahalanobis distance.

$$J(\mathbf{X}, \mathbf{Z}, \mathbf{M}) = \sum_{i=1}^N \sum_{k=1}^K z_{ik} (\mathbf{x}_i - \mathbf{m}_k)^T \Sigma^{-1} (\mathbf{x}_i - \mathbf{m}_k)$$

$$z_{ik} = \begin{cases} 1 & \text{if } k = \arg \min_l (\mathbf{x}_i - \mathbf{m}_l)^T \Sigma^{-1} (\mathbf{x}_i - \mathbf{m}_l) \\ 0 & \text{else} \end{cases}$$

$$\mathbf{m}_k^* = \arg \min_{\mathbf{m}_k} J(\mathbf{X}, \mathbf{Z}, \mathbf{M})$$

$$\frac{\partial J}{\partial \mathbf{m}_k} = \frac{\sum_{i=1}^N z_{ik} [2 \mathbf{x}_i^T \Sigma^{-1} - 2 \mathbf{x}_i^T \Sigma^{-1} \mathbf{m}_k + \mathbf{m}_k^T \Sigma^{-1} \mathbf{m}_k]}{\partial \mathbf{m}_k}$$

$$\sum_{i=1}^N z_{ik} (-2 \mathbf{x}_i^T \Sigma^{-1}) + (\Sigma^{-1} + \Sigma^{-1}) \mathbf{m}_k^T$$

b) Instead of considering data samples we will now consider how the cluster centroids μ_i (which are given and fixed) partition the space into K parts. Every part is defined as the space of all points $\mathbf{x} \in \mathbb{R}^d$ that would be assigned to the corresponding cluster k using the (squared) Mahalanobis distance.

Consider the following result using the (squared) Mahalanobis distance with Σ^{-1} , where the parts are denoted by the colors orange, pink, and green. Specify a possible Σ^{-1} that leads to this partitioning. Justify your answer.

$$\Sigma^{-1} = \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix}$$

$$(\mathbf{x}_1 - \mathbf{m}_1)^T \Sigma^{-1} (\mathbf{x}_1 - \mathbf{m}_1)$$

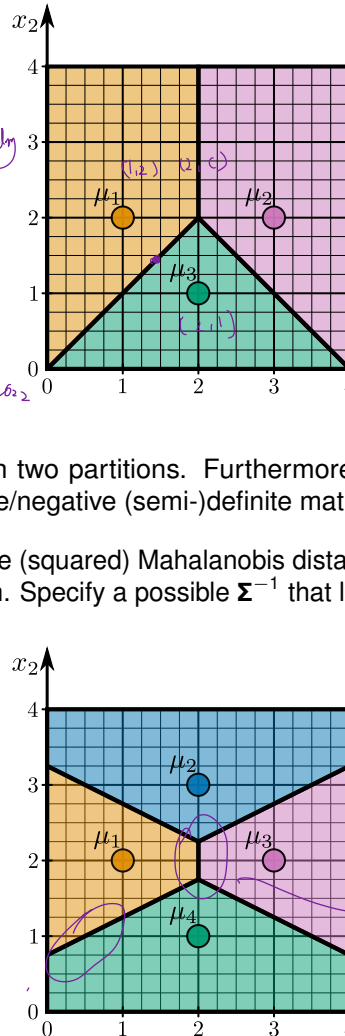
$$d(\mathbf{x}_1, \mathbf{m}_1) = d(\mathbf{x}_1, \mathbf{m}_2) \text{ for } \mathbf{x}_1 \text{ on boundary}$$

$$\begin{pmatrix} 1 & c \end{pmatrix} \Sigma^{-1} \begin{pmatrix} 1 \\ c \end{pmatrix} = \begin{pmatrix} -1 & c \end{pmatrix} \Sigma^{-1} \begin{pmatrix} -1 \\ c \end{pmatrix}$$

$$\begin{pmatrix} b_{11} + c b_{21} & b_{12} + c b_{22} \end{pmatrix} \begin{pmatrix} 1 \\ c \end{pmatrix} = \begin{pmatrix} -b_{11} + c b_{21} & -b_{12} + c b_{22} \end{pmatrix} \begin{pmatrix} -1 \\ c \end{pmatrix}$$

$$b_{11} + c b_{21} - b_{12} + c^2 b_{22} = -b_{11} + c b_{21} - b_{12} + c^2 b_{22}$$

$$b_{21} = 0$$



$$\sum_{i=1}^N z_{ik} \Sigma^{-1} (\mathbf{m}_k^T - \mathbf{x}_i^T) = 0$$

$$\sum_{i=1}^N z_{ik} \mathbf{m}_k^T = \sum_{i=1}^N z_{ik} \mathbf{x}_i^T$$

$$\mathbf{m}_k = \frac{\sum_{i=1}^N z_{ik} \mathbf{x}_i}{\sum_{i=1}^N z_{ik} = N_k}$$

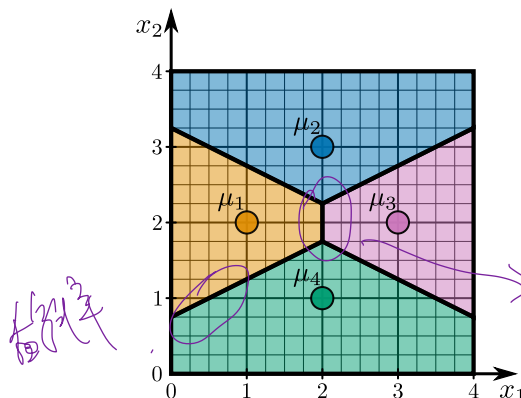
$$\begin{pmatrix} 1 & c \\ 1-c & 2-c \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} \begin{pmatrix} 1-c \\ 2-c \end{pmatrix} = \begin{pmatrix} 2-c & 1-c \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} \begin{pmatrix} 2-c \\ 1-c \end{pmatrix}$$

$$\begin{pmatrix} (1-c)b_{11} + (2-c)b_{21} & (1-c)b_{12} + (2-c)b_{22} \end{pmatrix} \begin{pmatrix} 1-c \\ 2-c \end{pmatrix}$$

$$\begin{pmatrix} (1-c)^2 b_{11} + (1-c)(2-c)b_{21} + (1-c)(2-c)b_{12} + (2-c)^2 b_{22} & (2-c)b_{11} + (1-c)b_{21} & (2-c)b_{12} + (1-c)b_{22} & (2-c)^2 b_{11} + (1-c)(2-c)b_{21} + (2-c)(1-c)b_{12} + (1-c)^2 b_{22} \end{pmatrix}$$

Hint: Consider the boundary between two partitions. Furthermore, the inverse of a symmetric matrix is also symmetric, and the inverse of a positive/negative (semi-)definite matrix is also positive/negative (semi-)definite.

c) Consider the following result using the (squared) Mahalanobis distance with Σ^{-1} , where the parts are denoted by the colors blue, orange, pink, and green. Specify a possible Σ^{-1} that leads to this partitioning. Justify your answer.



$$(4 - 2c + c^2 - 1 + 2c - c^2) b_{22}$$

$$= [4 - 2c + c^2 - 1 + 2c - c^2] b_{11}$$

$$b_{22} = b_{11}$$

$$\Sigma^{-1} = \begin{pmatrix} a & 0 \\ 0 & a \end{pmatrix}$$

$$\sigma_{12} = 0$$

Problem 9: K-means (Version B) (6 credits)

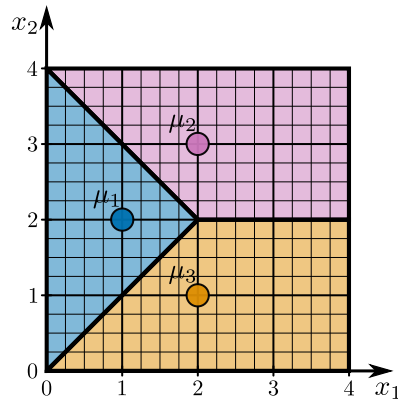
Consider a variant of K-means that uses the (squared) Mahalanobis distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \mathbf{x}_j), \quad (2)$$

with the covariance matrix $\boldsymbol{\Sigma}$, instead of the L_2 distance.

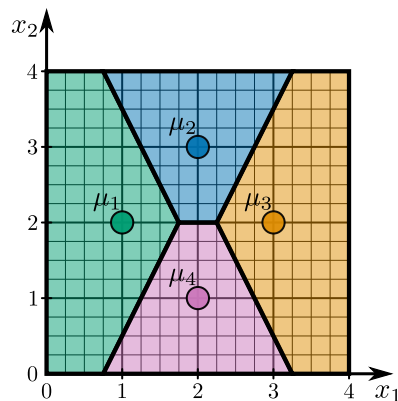
Hint: $\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{a}} = (\mathbf{X} + \mathbf{X}^T) \mathbf{a}$

- 0 ☐
- 1 ☐
- 2 ☐
- a) Derive the K-means cluster assignment and centroid updates for the (squared) Mahalanobis distance.
- 0 ☐
- 1 ☐
- 2 ☐
- b) Instead of considering data samples we will now consider how the cluster centroids μ_i (which are given and fixed) partition the space into K parts. Every part is defined as the space of all points $\mathbf{x} \in \mathbb{R}^d$ that would be assigned to the corresponding cluster k using the (squared) Mahalanobis distance.
- 0 ☐
- 1 ☐
- 2 ☐
- Consider the following result using the (squared) Mahalanobis distance with $\boldsymbol{\Sigma}^{-1}$, where the parts are denoted by the colors blue, orange, and pink. Specify a possible $\boldsymbol{\Sigma}^{-1}$ that leads to this partitioning. Justify your answer.



Hint: Consider the boundary between two partitions. Furthermore, the inverse of a symmetric matrix is also symmetric, and the inverse of a positive/negative (semi-)definite matrix is also positive/negative (semi-)definite.

- 0 ☐
- 1 ☐
- 2 ☐
- c) Consider the following result using the (squared) Mahalanobis distance with $\boldsymbol{\Sigma}^{-1}$, where the parts are denoted by the colors blue, orange, pink, and green. Specify a possible $\boldsymbol{\Sigma}^{-1}$ that leads to this partitioning. Justify your answer.



Problem 9: K-means (Version C) (6 credits)

Consider a variant of K-means that uses the (squared) Mahalanobis distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \mathbf{x}_j), \quad (3)$$

with the covariance matrix $\boldsymbol{\Sigma}$, instead of the L_2 distance.

Hint: $\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{a}} = (\mathbf{X} + \mathbf{X}^T) \mathbf{a}$

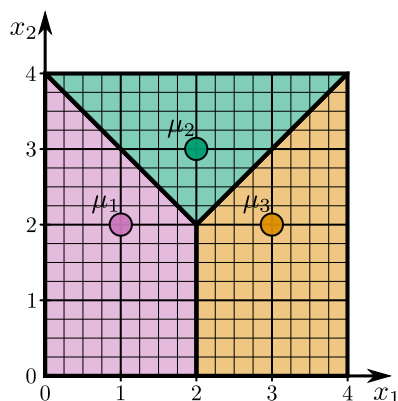
a) Derive the K-means cluster assignment and centroid updates for the (squared) Mahalanobis distance.

	0
	1
	2

b) Instead of considering data samples we will now consider how the cluster centroids μ_i (which are given and fixed) partition the space into K parts. Every part is defined as the space of all points $\mathbf{x} \in \mathbb{R}^d$ that would be assigned to the corresponding cluster k using the (squared) Mahalanobis distance.

	0
	1
	2

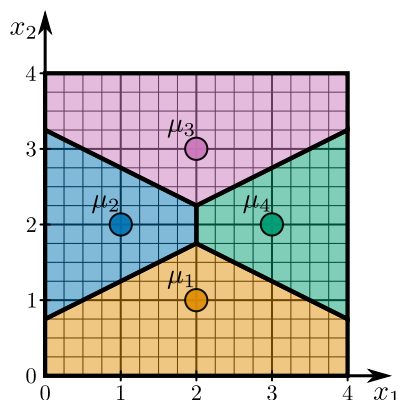
Consider the following result using the (squared) Mahalanobis distance with $\boldsymbol{\Sigma}^{-1}$, where the parts are denoted by the colors green, orange, and pink. Specify a possible $\boldsymbol{\Sigma}^{-1}$ that leads to this partitioning. Justify your answer.



Hint: Consider the boundary between two partitions. Furthermore, the inverse of a symmetric matrix is also symmetric, and the inverse of a positive/negative (semi-)definite matrix is also positive/negative (semi-)definite.

c) Consider the following result using the (squared) Mahalanobis distance with $\boldsymbol{\Sigma}^{-1}$, where the parts are denoted by the colors blue, orange, pink, and green. Specify a possible $\boldsymbol{\Sigma}^{-1}$ that leads to this partitioning. Justify your answer.

	0
	1
	2



Problem 9: K-means (Version D) (6 credits)

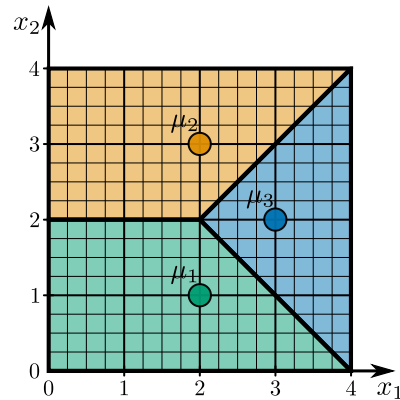
Consider a variant of K-means that uses the (squared) Mahalanobis distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \mathbf{x}_j), \quad (4)$$

with the covariance matrix $\boldsymbol{\Sigma}$, instead of the L_2 distance.

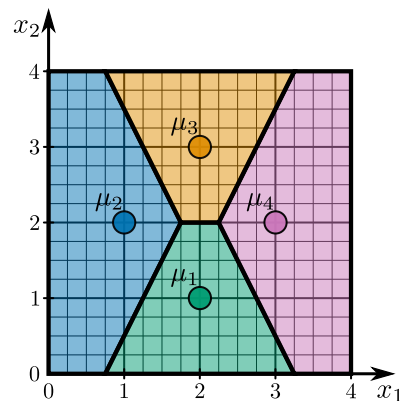
Hint: $\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{a}} = (\mathbf{X} + \mathbf{X}^T) \mathbf{a}$

- 0 ☐
- 1 ☐
- 2 ☐
- a) Derive the K-means cluster assignment and centroid updates for the (squared) Mahalanobis distance.
- 0 ☐
- 1 ☐
- 2 ☐
- b) Instead of considering data samples we will now consider how the cluster centroids μ_i (which are given and fixed) partition the space into K parts. Every part is defined as the space of all points $\mathbf{x} \in \mathbb{R}^d$ that would be assigned to the corresponding cluster k using the (squared) Mahalanobis distance.
- 0 ☐
- 1 ☐
- 2 ☐
- Consider the following result using the (squared) Mahalanobis distance with $\boldsymbol{\Sigma}^{-1}$, where the parts are denoted by the colors orange, blue, and green. Specify a possible $\boldsymbol{\Sigma}^{-1}$ that leads to this partitioning. Justify your answer.



Hint: Consider the boundary between two partitions. Furthermore, the inverse of a symmetric matrix is also symmetric, and the inverse of a positive/negative (semi-)definite matrix is also positive/negative (semi-)definite.

- 0 ☐
- 1 ☐
- 2 ☐
- c) Consider the following result using the (squared) Mahalanobis distance with $\boldsymbol{\Sigma}^{-1}$, where the parts are denoted by the colors blue, orange, pink, and green. Specify a possible $\boldsymbol{\Sigma}^{-1}$ that leads to this partitioning. Justify your answer.



Problem 10: Differential Privacy (Version A) (6 credits)

You are given a dataset with n instances $\{x_1, \dots, x_n\}$, with $x_i \in \mathcal{X}$. The instances are randomly split into disjoint groups G_1, G_2, \dots, G_m , each of size $\frac{n}{m}$ (assume that m divides n , i.e. $\frac{n}{m}$ is an integer).

a) First you apply an *arbitrary* function $f : \mathcal{X}^{\frac{n}{m}} \rightarrow [a, b]$ (where a and b are given constants) to each of the groups, i.e. you compute $g_1 = f(G_1), g_2 = f(G_2), \dots, g_m = f(G_m)$. Then you compute the final output by aggregating the per-group outputs by computing $\text{mean}(g_1, \dots, g_m)$. Derive the global Δ_1 sensitivity of the function $f' := \text{mean}(f(G_1), \dots, f(G_m))$?

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

b) How does the global sensitivity of f' change if we increase n keeping m fixed? How does the global sensitivity of f' change if we increase m keeping n fixed?

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

c) Can you make the function f' differentially private for any function f ? If yes, specify the noise distribution from which we have to sample to obtain an ϵ -DP private mechanism. If no, why not?

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

Problem 10: Differential Privacy (Version B) (6 credits)

You are given a dataset with n instances $\{x_1, \dots, x_n\}$, with $x_i \in \mathcal{X}$. The instances are randomly split into disjoint groups G_1, G_2, \dots, G_m , each of size $\frac{n}{m}$ (assume that m divides n , i.e. $\frac{n}{m}$ is an integer).

- 0 ☐ a) First you apply an *arbitrary* function $f : \mathcal{X}^{\frac{n}{m}} \rightarrow [a, b]$ (where a and b are given constants) to each of the groups, i.e. you compute $g_1 = f(G_1), g_2 = f(G_2), \dots, g_m = f(G_m)$. Then you compute the final output by aggregating the per-group outputs by computing $\text{median}(g_1, \dots, g_m)$. Derive the global Δ_1 sensitivity of the function $f' := \text{median}(f(G_1), \dots, f(G_m))$?

1 ☐

2 ☐

3 ☐

- 0 ☐ b) How does the global sensitivity of f' change if we increase n keeping m fixed? How does the global sensitivity of f' change if we increase m keeping n fixed?

1 ☐

2 ☐

- 0 ☐ c) Can you make the function f' differentially private for any function f ? If yes, specify the noise distribution from which we have to sample to obtain an ϵ -DP private mechanism. If no, why not?

1 ☐

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

This image shows a full page of blank graph paper. The grid consists of small, equal-sized squares formed by thin gray lines. There are 20 columns and 20 rows of squares, creating a total of 400 square units. The grid covers the entire area of the page, leaving no margins or other markings.

