

## Machine Learning — Repeat Exam

1	2	3	4	5	6	7	8	9	10	11	$\Sigma$
4	5	6	4	6	5	4	7	4	2	7	54

*Do not write anything above this line*

Name:

Student ID:

Signature:

- Only write on the sheets given to you by supervisors. If you need more paper, ask the supervisors.
- Pages 16-18 can be used as scratch paper.
- All sheets (including scratch paper) have to be returned at the end.
- **Do not unstaple the sheets!**
- Wherever answer boxes are provided, please write your answers in them.
- Please write your student ID (*Matrikelnummer*) on every sheet you hand in.
- **Only use a black or a blue pen (no pencils, red or green pens!).**
- You are allowed to use your A4 sheet of handwritten notes (two sides). **No other materials (e.g. books, cell phones, calculators) are allowed!**
- Exam duration - 120 minutes.
- This exam consists of 18 pages, 11 problems. You can earn 54 points.

### Probability distributions

For your reference, we provide the following probability distribution.

- Univariate normal distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

- Bernoulli distribution

$$\text{Bern}(x|\theta) = \theta^x(1-\theta)^{1-x}$$

## Decision Trees



**Problem 1 [(2+2)=4 points]** Assume you want to build a decision tree. Your data set consists of  $N$  samples, each with  $k$  features ( $k \leq N$ ).

- a) If the features are binary, what is the maximum possible number of leaf nodes and the maximum depth of your decision tree?

minimum leaf nod.  ~~$2^{k-1} + 2^{k-2} + \dots + 2^0$~~   $\min(2^k, N)$

maximum depth.  ~~$k \leq k$~~

完全分类. (但是  $k$  已经很大)

- b) If the features are continuous, what is the maximum possible number of leaf nodes and the maximum depth of your decision tree?

minimum leaf node are  $\underline{\underline{N}}$

maximum depth are  $\underline{\underline{N}}$

## Regression

**Problem 2 [(1+4)=5 points]** We want to perform regression on a dataset consisting of  $N$  samples  $\mathbf{x}_i \in \mathbb{R}^D$  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 fitted an  $L_2$ -regularized linear regression model and obtained the optimal weight vector  $\mathbf{w}^* \in \mathbb{R}^D$  as

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

Note that there is no bias term.

Now, assume that we obtained a new data matrix  $\mathbf{X}_{new}$  by scaling all samples by the same positive factor  $a \in (0, \infty)$ . That is,  $\mathbf{X}_{new} = a\mathbf{X}$  (and respectively  $\mathbf{x}_i^{new} = a\mathbf{x}_i$ ).

- a) Find the weight vector  $\mathbf{w}_{new}$  that will produce the same predictions on  $\mathbf{X}_{new}$  as  $\mathbf{w}^*$  produces on  $\mathbf{X}$ .

$$\begin{aligned} \hat{y} &= \mathbf{x}^T \mathbf{w} = \mathbf{x}_{new}^T \mathbf{w}_{new} \\ \mathbf{x}^T \mathbf{w} &= a \mathbf{x}^T \mathbf{w}_{new} \\ \mathbf{w}_{new} &= \frac{1}{a} \mathbf{w}^* \end{aligned}$$

✓

- b) Find the regularization factor  $\lambda_{new} \in \mathbb{R}$ , such that the solution  $\mathbf{w}_{new}^*$  of the new  $L_2$ -regularized linear regression problem

$$\mathbf{w}_{new}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i^{new} - y_i)^2 + \frac{\lambda_{new}}{2} \mathbf{w}^T \mathbf{w}$$

will produce the same predictions on  $\mathbf{X}_{new}$  as  $\mathbf{w}^*$  produces on  $\mathbf{X}$ .

Provide a mathematical justification for your answer.

$$\begin{aligned} L &= \frac{1}{2} (\mathbf{x}_{new}^T \mathbf{w} - y)^T (\mathbf{x}_{new}^T \mathbf{w} - y) + \frac{\lambda_{new}}{2} \mathbf{w}^T \mathbf{w} \\ \frac{\partial L}{\partial \mathbf{w}} &= (\mathbf{x}_{new}^T \mathbf{x}_{new} \mathbf{w} - \mathbf{x}_{new}^T y) + \lambda_{new} \mathbf{w} = 0 \\ \mathbf{w}_{new}^* &= (\mathbf{x}_{new}^T \mathbf{x}_{new} + \lambda_{new})^{-1} \mathbf{x}_{new}^T y \quad \mathbf{w}_{new}^* = \frac{1}{a} \mathbf{w}^* \\ &= a (\mathbf{x}^T \mathbf{x} + \lambda)^{-1} \mathbf{x}^T y \\ &= \frac{1}{a} (\mathbf{x}^T \mathbf{x} + \boxed{\frac{\lambda_{new}}{a^2}})^{-1} \mathbf{x}^T y \\ &\quad \boxed{\lambda} = \frac{\lambda_{new}}{a^2} \Leftrightarrow \lambda_{new} = a^2 \lambda \end{aligned}$$

✓

## Classification

**Problem 3 [(1+2+3)=6 points]** We would like to perform binary classification on multivariate binary data. That is, the data points  $\mathbf{x}_i \in \{0, 1\}^D$  are binary vectors of length  $D$ , and each sample belongs to one of two classes  $y_i \in \{1, 2\}$ .

Consider the following generative classification model. We place a categorical prior on  $y$

$$p(y = 1) = \pi_1 \quad p(y = 2) = \pi_2.$$

The class-conditional distributions are products of independent Bernoulli distributions

$$p(\mathbf{x} | y = 1, \boldsymbol{\alpha}) = \prod_{j=1}^D \text{Bern}(x_j | \alpha_j), \quad \alpha_j^y (1 - \alpha_j)^{1-y}$$

$$\text{Bern}(x_j | \theta) = \theta^x (1 - \theta)^{1-x}$$

$$p(\mathbf{x} | y = 2, \boldsymbol{\beta}) = \prod_{j=1}^D \text{Bern}(x_j | \beta_j),$$

where  $\boldsymbol{\alpha} \in [0, 1]^D$  and  $\boldsymbol{\beta} \in [0, 1]^D$  are the respective parameter vectors for both classes.

That is, each component  $x_j$  is distributed as  $x_j \sim \text{Bern}(\alpha_j)$  if  $y = 1$  or  $x_j \sim \text{Bern}(\beta_j)$  if  $y = 2$ .

- a) Write down the expression for the posterior distribution  $p(y | \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi})$ .

$$\begin{aligned} p(y | \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) &\propto \frac{p(\mathbf{x}|y, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) p(y|\boldsymbol{\pi})}{p(\boldsymbol{\alpha}|\boldsymbol{\beta}, \boldsymbol{\pi})} \\ &= \frac{p(\mathbf{x}|y, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) p(y|\boldsymbol{\pi})}{p(\mathbf{x}|y=1, \boldsymbol{\alpha}) p(y=1|\boldsymbol{\pi}) + p(\mathbf{x}|y=2, \boldsymbol{\beta}) p(y=2|\boldsymbol{\pi})} \end{aligned}$$

- b) Assume that  $D = 3$ ,  $\boldsymbol{\alpha} = [1/3, 1/3, 3/4]$ ,  $\boldsymbol{\beta} = [2/3, 1/2, 1/2]$ ,  $\pi_1 = 1/3$  and  $\pi_2 = 2/3$ .

Write down a data point  $\mathbf{x}_1 \in \{0, 1\}^3$  that will be classified as class 1 by our model. Additionally, compute the posterior probability  $p(y = 1 | \mathbf{x}_1, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi})$ .

$$\begin{aligned} \mathbf{x}_1 &= (0, 1, 1)^\top \\ p(y=1 | \mathbf{x}_1, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) &= \frac{p(\mathbf{x}_1 | y=1, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) p(y=1|\boldsymbol{\pi})}{p(\mathbf{x}_1 | y=1, \boldsymbol{\alpha}) p(y=1|\boldsymbol{\pi}) + p(\mathbf{x}_1 | y=2, \boldsymbol{\beta}) p(y=2|\boldsymbol{\pi})} \\ &= \frac{(1-\alpha_1)(1-\alpha_2)\alpha_3 \cdot \pi_1}{(1-\alpha_1)(1-\alpha_2)\alpha_3 \pi_1 + (1-\beta_1)(1-\beta_2)\beta_3 \pi_2} \\ &= \frac{\frac{2}{3} \cdot \frac{1}{2} \cdot \frac{1}{4} \cdot \frac{1}{3}}{\frac{2}{3} \cdot \frac{1}{2} \cdot \frac{1}{4} \cdot \frac{1}{3} + \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{2}{3}} = \frac{\frac{1}{24}}{\frac{1}{24} + \frac{1}{12}} = \frac{1}{3} \end{aligned}$$

- c) Consider the case when  $D = 2$ ,  $\pi_1 = \pi_2 = 1/2$ , and  $\boldsymbol{\alpha} \in [0, 1]^2$  and  $\boldsymbol{\beta} \in [0, 1]^2$  are known and fixed. Show that the resulting classification rule can be represented as a linear function of  $\mathbf{x}$ . That is, find  $\mathbf{w} \in \mathbb{R}^2$  and  $b \in \mathbb{R}$ , such that

$$\{\mathbf{x} \in \{0, 1\}^2 : \mathbf{w}^T \mathbf{x} + b > 0\} = \{\mathbf{x} \in \{0, 1\}^2 : p(y = 1 | \mathbf{x}) > p(y = 2 | \mathbf{x})\}$$

$$p(y=1|x) = \frac{p(x|y=1)}{p(x|y=1) + p(x|y=2)}$$

$$\frac{p(y=1|x)}{p(y=2|x)} > 1.$$

$\text{Bern}(x|\theta) = \theta^x(1-\theta)^{1-x}$

$$p(x|y=1, \alpha) = \prod_{j=1}^D \text{Bern}(x_j | \alpha_j),$$

$$p(x|y=2, \beta) = \prod_{j=1}^D \text{Bern}(x_j | \beta_j),$$

$$\frac{p(x_1|y=1)}{p(x_1|y=2)} > 1$$

$$p(x_1|y=1) = p(x_1|y=1) \cdot p(x_2|y=1)$$

$$\log p(x_1|y=1) + \log p(x_2|y=1) > \log p(x_1|y=2) + \log p(x_2|y=2)$$

$$x_1 \ln \alpha_1 + (1-x_1) \ln (1-\alpha_1) + x_2 \ln \beta_2 + (1-x_2) \ln (1-\beta_2) - x_1 \ln \beta_1 - (1-x_1) \ln (1-\beta_1) - x_2 \ln \alpha_2 - (1-x_2) \ln (1-\alpha_2) > 0$$

$$\underbrace{[\ln \alpha_1 - \ln (1-\alpha_1) - \ln \beta_1 + \ln (1-\beta_1)]}_{w_1} x_1 + \underbrace{[\ln \beta_2 - \ln (1-\beta_2) - \ln \alpha_2 + \ln (1-\alpha_2)]}_{w_2} x_2 + \underbrace{[\ln (1-\alpha_1) + \ln (1-\alpha_2) - \ln (1-\beta_1) - \ln (1-\beta_2)]}_{b} > 0$$

## Kernels

**Problem 4 [(4)=4 points]** Prove or disprove whether the following operations on sets  $A, B \subseteq \mathcal{X}$ , where  $\mathcal{X}$  is a finite set, define a valid kernel.

- a)  $k(A, B) = |A \times B|$ , where  $A \times B = \{(a, b) : a \in A, b \in B\}$  denotes the cartesian product and  $|S|$  denotes the cardinality of set  $S$ , i.e. the number of elements in  $S$ .

$$|A \times B| = |A| \cdot |B|$$

b)  $k(A, B) = |A \cap B|$

denote  $x_i$  is member of  $X$   $\Leftrightarrow \{x_1, x_2, x_3 \dots x_i\} \in X$

denote  $c_i = \begin{cases} 1, & \text{if } x_i \in A \\ 0, & \text{else} \end{cases}$   $d_i = \begin{cases} 1, & \text{if } x_i \in B \\ 0, & \text{else} \end{cases}$

$$|A \cap B| = C^T D$$

c)  $k(A, B) = |A \cup B|$

Not a kernel

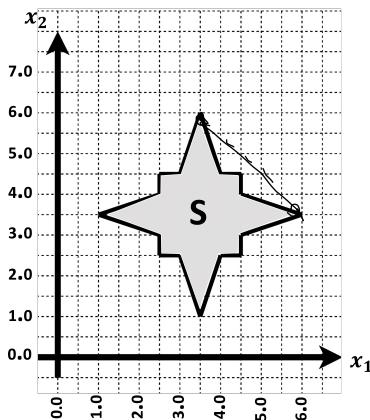
$ A $	$ B $	$ A \cap B $	$ A \cup B $
1	0	0	1
0	1	0	1
0	0	0	0
1	1	1	1

## Optimization

**Problem 5 [(1+3+2)=6 points]** Let  $f$  be the following convex function on  $\mathbb{R}^2$ :

$$f(x_1, x_2) = e^{x_1+x_2} - 5 \cdot \log(x_2)$$

- a) Consider the following shaded region  $S$  in  $\mathbb{R}^2$ . Is this region convex? Why?



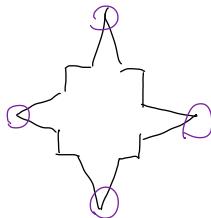
No points on the line is out of region

- b) Find the maximizer  $(x_1^*, x_2^*)$  of  $f$  over the shaded region  $S$ . For your computations, you can pick values from the following table. Justify your answer.

$e^{4.5} = 90.017$	$e^{5.0} = 148.41$	$e^{5.5} = 244.69$	$e^{6.5} = 665.14$
$e^{7.0} = 1096.63$	$e^{7.5} = 1808.04$	$e^{8.0} = 2980.95$	$e^{8.5} = 4914.76$
$e^{9.0} = 8103.08$	$e^{9.5} = 13359.726$	$e^{10.0} = 22026.46$	$e^{10.5} = 36315.50$
$\log(1.0) = 0$	$\log(2.5) = 0.9162$	$\log(3.0) = 1.0986$	$\log(3.5) = 1.2527$
$\log(4.0) = 1.3862$	$\log(4.5) = 1.5040$	$\log(5.0) = 1.6094$	$\log(6.0) = 1.7917$

Few corner

Then find the maximum



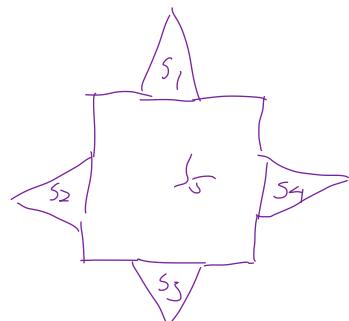
- c) Assume that we are given an algorithm  $\text{ConvOpt}(f, \mathcal{X})$  that takes as input a convex function  $f$  and any convex region  $\mathcal{X}$ , and returns the minimum of  $f$  over  $\mathcal{X}$ .

Using the ConvOpt algorithm, how would you find the global minimum of  $f$  over the shaded region  $S$ ?

Divide the non-convex region  $S$  into convex region.

And calculate the minimum

Then find the smaller one  
the global minima



## SVM

**Problem 6 [(5)=5 points]** Given the data points

$$\mathbf{x}_1 = (1, 1, 0, 1)^T \quad \mathbf{x}_2 = (1, 1, 1, 0)^T \quad \mathbf{x}_3 = (0, 1, 1, 1)^T \quad \mathbf{x}_4 = (0, 0, 1, 1)^T$$

Prove or disprove whether the following combinations of labels  $\mathbf{y}$  and dual variables  $\boldsymbol{\alpha}$  are the optimal solutions of a soft-margin SVM with  $C = 1$ .

a)  $\mathbf{y} = (-1, -1, 1, 1)^T$ ,  $\boldsymbol{\alpha} = (0.6, 0.6, 1, 0)^T$

b)  $\mathbf{y} = (-1, -1, 1, 1)^T$ ,  $\boldsymbol{\alpha} = (\frac{2}{3}, \frac{2}{3}, \frac{4}{3}, 0)^T$

c)  $\mathbf{y} = (-1, 1, -1, 1)^T$ ,  $\boldsymbol{\alpha} = (1, 1, 1, 1)^T$

1) optimal solution for soft SVM.

$$\sum_{i=1}^n \alpha_i y_i = -0.6 - 0.6 + 1 + 0 = -0.2 \neq 0$$

$$\sum_{i=1}^n \alpha_i y_i = -\frac{2}{3} - \frac{2}{3} + \frac{4}{3} = 0 = 0 \quad \text{X} \quad \cancel{\alpha_4 = \frac{4}{3} > C = 1} \quad \text{Wrong}$$

$$\sum_{i=1}^n \alpha_i y_i = -1 + -1 + 1 = 0 = 0 \quad \checkmark$$

## Deep Learning

**Problem 7 [(2+2)=4 points]** You are trying to solve a regression task and you want to choose between two approaches:

1. A simple linear regression model.
2. A feed forward neural network  $f_{\mathbf{W}}(\mathbf{x})$  with  $L$  hidden layers, where each hidden layer  $l \in \{1, \dots, L\}$  has a weight matrix  $\mathbf{W}_l \in \mathbb{R}^{D \times D}$  and a ReLU activation function. The output layer has a weight matrix  $\mathbf{W}_{L+1} \in \mathbb{R}^{D \times 1}$  and no activation function.

In both models, there are no bias terms.

Your dataset  $\mathcal{D}$  contains data points with nonnegative features  $\mathbf{x}_i$  and the target  $y_i$  is continuous:

$$\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N, \quad \mathbf{x}_i \in \mathbb{R}_{\geq 0}^D, \quad y_i \in \mathbb{R}$$

Let  $\mathbf{w}_{LS}^* \in \mathbb{R}^D$  be the optimal weights for the linear regression model corresponding to a global minimum of the following least squares optimization problem:

$$\mathbf{w}_{LS}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^D} \mathcal{L}_{LS}(\mathbf{w}) = \arg \min_{\mathbf{w} \in \mathbb{R}^D} \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$

Let  $\mathbf{W}_{NN}^* = \{\mathbf{W}_1^*, \dots, \mathbf{W}_{L+1}^*\}$  be the optimal weights for the neural network corresponding to a global minimum of the following optimization problem:

$$\mathbf{W}_{NN}^* = \arg \min_{\mathbf{W}} \mathcal{L}_{NN}(\mathbf{W}) = \arg \min_{\mathbf{W}} \frac{1}{2} \sum_{i=1}^N (f_{\mathbf{W}}(\mathbf{x}_i) - y_i)^2$$

- a) Assume that the optimal  $\mathbf{W}_{NN}^*$  you obtain are non-negative.

What will be the relation ( $<, \leq, =, \geq, >$ ) between the neural network loss  $\mathcal{L}_{NN}(\mathbf{W}_{NN}^*)$  and the linear regression loss  $\mathcal{L}_{LS}(\mathbf{w}_{LS}^*)$ ? Provide a mathematical argument to justify your answer.

Non negative

$$\begin{aligned} f(\mathbf{w}_{NN}^*) &= \min_{\mathbf{w}} \frac{1}{2} \ell(\mathcal{U}(\text{ReLU}(\mathbf{x}_i^T \cdot \mathbf{w}_1^*))) \mathcal{U}_2^* \dots \mathcal{U}_{L+1}^* \\ &\approx \mathbf{x}_i^T \mathcal{U}_1^* \dots \mathcal{U}_{L+1}^* \\ &\approx \mathbf{x}_i^T \mathbf{w}^* \end{aligned}$$

$$E_{LS}(\mathbf{w}_{LS}^*) = E_{NN}(\mathbf{w}_{NN}^*)$$

- b) In contrast to (a), now assume that the optimal weights  $\mathbf{w}_{LS}^*$  you obtain are non-negative. What will be the relation ( $<$ ,  $\leq$ ,  $=$ ,  $\geq$ ,  $>$ ) between the linear regression loss  $\mathcal{L}_{LS}(\mathbf{w}_{LS}^*)$  and the neural network loss  $\mathcal{L}_{NN}(\mathbf{W}_{NN}^*)$ ? Provide a mathematical argument to justify your answer.

$$\mathcal{L}_{NN}(\mathbf{W}_{NN}^*) < \mathcal{L}_{LS}(\mathbf{w}_{LS}^*)$$

- (1) Non-negative  $\mathbf{w}_{LS}^*$  mean nothing
- (2) NN must perform better than Linear regression

## Dimensionality Reduction

**Problem 8 [(3+2+2)=7 points]** You are given  $N = 4$  data points:  $\{\mathbf{x}_i\}_{i=1}^4, \mathbf{x}_i \in \mathbb{R}^3$ , represented with the matrix  $\mathbf{X} \in \mathbb{R}^{4 \times 3}$ .

$$\mathbf{X} = \begin{bmatrix} 4 & 3 & 2 \\ 2 & 1 & -2 \\ 4 & -1 & 2 \\ -2 & 1 & 2 \end{bmatrix}$$

*Hint: In this task the results of all (final and intermediate) computations happen to be integers.*

- a) Perform principal component analysis (PCA) of the data  $\mathbf{X}$ , i.e. find the principal components and their associated variances in the transformed coordinate system. Show your work.

$$\begin{aligned} \bar{\mathbf{x}} &= \frac{1}{N} \mathbf{X}^T \cdot \mathbf{1}_N = \frac{1}{4} \begin{bmatrix} 4 & 2 & 4 & -2 \\ 2 & 1 & -1 & 1 \\ 4 & -1 & 1 & 1 \\ -2 & 1 & 2 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} \\ \tilde{\mathbf{x}} &= \mathbf{X} - \bar{\mathbf{x}} = \begin{bmatrix} 2 & 2 & 1 \\ 0 & 0 & -3 \\ 2 & -2 & 1 \\ -4 & 0 & 1 \end{bmatrix} \\ \Sigma_{\tilde{\mathbf{x}}} &= \frac{1}{N} \tilde{\mathbf{x}}^T \tilde{\mathbf{x}} = \frac{1}{4} \begin{bmatrix} 2 & 0 & 2 & -4 \\ 2 & 0 & -2 & 0 \\ 1 & -3 & 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & 2 & 1 \\ 0 & 0 & -3 \\ 2 & -2 & 1 \\ -4 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 5 & 3 \\ 0 & 3 & 3 \end{bmatrix} = \Lambda \\ &\quad \Gamma = \mathbb{I}_3 \end{aligned}$$

- b) Project the data to two dimensions, i.e. write down the transformed data matrix  $\mathbf{Y} \in \mathbb{R}^{4 \times 2}$

using the top-2 principal components you computed in (a). What fraction of variance of  $\mathbf{X}$  is preserved by  $\mathbf{Y}$ ?

$$\begin{aligned} \mathbf{Y} &= \tilde{\mathbf{X}} \cdot \mathbf{P}_{\text{reduced}} \\ &= \begin{bmatrix} 2 & 2 & 1 \\ 0 & 0 & -3 \\ 2 & -2 & 1 \\ -4 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 0 & -3 \\ 2 & 1 \\ -4 & 1 \end{bmatrix} \end{aligned}$$

- c) Let  $\mathbf{x}_5 \in \mathbb{R}^3$  be a new data point. Specify the vector  $\mathbf{x}_5$  such that performing PCA on the data including the new data point  $\{\mathbf{x}_i\}_{i=1}^5$  leads to exactly the same principal components as in (a).

$$\tilde{\mathbf{x}}_5 = \text{mean} = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$$

## Clustering

**Problem 9 [(4)=4 points]** Let  $\mu_1, \dots, \mu_K$  be the centroids computed by the  $K$ -means algorithm. Prove that the set  $\mathcal{X}_j$  of all points in  $\mathbb{R}^D$  assigned during inference to the cluster  $j$  is a convex set.

$$\mathcal{X}_j := \{\mathbf{x} \in \mathbb{R}^D : \mathbf{x} \text{ would be assigned to centroid } \mu_j \text{ by } K\text{-means}\}$$

*Hint: start by thinking about the case with  $K = 2$ .*

--

**Problem 10 [(2)=2 points]** Given three 1-dimensional Gaussian distributions  $\mathcal{N}(\mu_i, \sigma_i^2)$  with parameters

$$\begin{array}{lll} \mu_1 = 1, & \text{(x-1)}^2 & \mu_2 = -1, & \text{(x+1)} \\ \sigma_1 = 1, & & \sigma_2 = 0.5, & \\ & & & \mu_3 = 0, \\ & & & \sigma_3 = 2.5 \end{array}$$

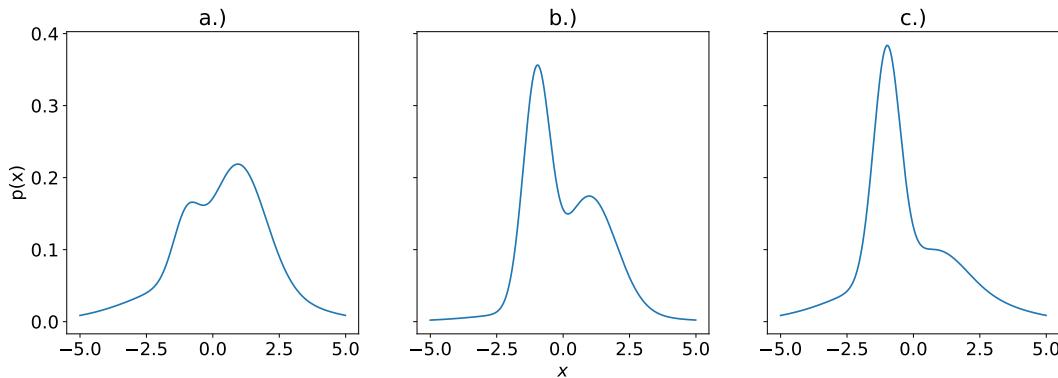
and three different vectors of mixing coefficients  $\pi$  defining categorical cluster priors.

Match the value of  $\pi$  in each row of the following table with one of the probability density functions

$$p(x) = \sum_{i=1}^3 \pi_i \mathcal{N}(x | \mu_i, \Sigma_i) \quad \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

of the resulting GMM showed below. Fill in the last column of the table, no argumentation required.

	$\pi_1$	$\pi_2$	$\pi_3$	PDF (a, b or c)
case 1	0.111...	0.444...	0.444...	c
case 2	0.444...	0.111...	0.444...	a
case 3	0.444...	0.444...	0.111...	b



## Variational Inference

**Problem 11 [(3+1+1+2)=7 points]** Consider the following latent variable probabilistic model

$$\begin{aligned} p(z) &= \mathcal{N}(z | 0, 1) \\ p(x | z) &= \mathcal{N}(x | z, 1) \end{aligned}$$

We want to approximate the posterior distribution  $p(z | x)$  using the following variational family

$$\mathcal{Q} = \{\mathcal{N}(z | \mu, 1) \text{ for } \mu \in \mathbb{R}\}$$

that includes all normal distributions with unit variance.

Questions (a), (b), (c) and (d) are all concerning this setup.

*Hint: Variance of  $p(z | x)$  is equal to 0.5.*

- a) Write down the closed-form expression for ELBO  $\mathcal{L}(q)$  and simplify it. You can ignore all the terms constant in  $\mu$ .

- b) Find the optimal variational distribution  $q^* \in \mathcal{Q}$  that maximizes the ELBO

$$q^* = \arg \max_{q \in \mathcal{Q}} \mathcal{L}(q)$$

i.e. find the mean  $\mu^*$  of the optimal variational distribution  $q^*$ .

c) Assume that the optimal  $q^*$  (i.e., the optimal  $\mu^*$ ) from question (b) is given. Which of the following statements is true?

- (1)  $\text{KL}(q(z | \mu^*) \| p(z | x)) < 0$
- (2)  $\text{KL}(q(z | \mu^*) \| p(z | x)) = 0$
- (3)  $\text{KL}(q(z | \mu^*) \| p(z | x)) > 0$

Justify your answer.

d) For each of the conditions (1), (2), (3) from question (c) above, provide a parametric variational family  $\mathcal{Q}_i$ , such that the optimal  $q_i^*$  from each family would fulfill the respective condition, or explain why it's impossible.

That is, provide  $\mathcal{Q}_1$ , such that for  $q_1^* = \arg \max_{q \in \mathcal{Q}_1} \mathcal{L}(q)$  we have  $\text{KL}(q_1^*(z) \| p(z | x)) < 0$ , for  $q_2^* = \arg \max_{q \in \mathcal{Q}_2} \mathcal{L}(q)$  we have  $\text{KL}(q_2^*(z) \| p(z | x)) = 0$ , and for  $q_3^* = \arg \max_{q \in \mathcal{Q}_3} \mathcal{L}(q)$  we have  $\text{KL}(q_3^*(z) \| p(z | x)) > 0$ .



