

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: Thursday 17th February, 2022

Time: 17:00 – 19:00

Working instructions

- This graded exercise consists of **52 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 **96**.
- 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: Probabilistic inference (Version A) (10 credits)

Consider the following probabilistic model:

$$\mathbb{P}(\theta | \lambda, \alpha) = \begin{cases} \frac{\alpha \lambda^\alpha}{\theta^{\alpha+1}} & \text{if } \lambda \leq \theta \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbb{P}(x | \theta) = \begin{cases} \frac{1}{\theta} & \text{if } 0 \leq x \leq \theta \\ 0 & \text{otherwise} \end{cases}$$

with $\lambda > 0, \alpha > 0$ and a set of observations $\mathcal{D} = \{x_1, \dots, x_N\}$ consisting of N samples $x_i \in \mathbb{R}_+$ generated from the above probabilistic model.

0 Derive the posterior distribution $\mathbb{P}(\theta | \mathcal{D}, \lambda, \alpha)$.
 1
 2
 3
 4
 5
 6
 7
 8
 9
 10

posterior \propto likelihood · prior

$$\mathbb{P}(\theta | \mathcal{D}, \lambda, \alpha) = \prod_{i=1}^N \mathbb{P}(x_i | \theta) \cdot \mathbb{P}(\theta | \lambda, \alpha)$$

$$= \frac{1}{\theta^N} \cdot \frac{1}{\max(x_i) \leq \theta} \cdot \frac{1}{\theta^{\alpha+1}} \cdot \frac{1}{\max(\lambda) \leq \theta} \cdot \alpha \cdot \lambda^\alpha$$

$$\propto \frac{1}{\theta^{N+\alpha+1}} \cdot \frac{1}{\max(x_1, \dots, x_N, \lambda) \leq \theta}$$

We recognize that the posterior distribution is $\mathbb{P}(\theta | \{x_1, \dots, x_N\}, \lambda, \alpha) = \mathbb{P}(\theta | \lambda_{\text{new}}, \alpha_{\text{new}}) = \text{Pareto}(\lambda_{\text{new}} = \max(x_1, \dots, x_N, \lambda), \alpha_{\text{new}} = N + \alpha)$

Problem 2: Linear regression (Version A) (8 credits)

We want to perform regression on a dataset consisting of N samples $\mathbf{x}_i \in \mathbb{R}^D$ 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 a linear regression model and obtained the optimal weight vector $\mathbf{w}^* \in \mathbb{R}^D$ as

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

Note that there is no bias term.

Now, assume that we normalize the target variables to have a variance of 1, i.e. $\mathbf{y}_{\text{new}} = \frac{1}{\sigma} \cdot \mathbf{y}$ with $\sigma = \text{Var}(\mathbf{y})$, where $\text{Var}(\mathbf{y})$ is the sample variance of \mathbf{y} .

0 Find the data matrix $\mathbf{X}_{\text{new}} \in \mathbb{R}^{N \times D}$ such that the solution to the new problem:

$$\mathbf{w}_{\text{new}}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_{\text{new},i} - y_{\text{new},i})^2$$

will be the same as the solution to the previous problem i.e. $\mathbf{w}_{\text{new}}^* = \mathbf{w}^*$. Justify your answer.

Note: $\mathbf{x}_{\text{new},i}$ is row i of \mathbf{X}_{new} , represented as a column vector.

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

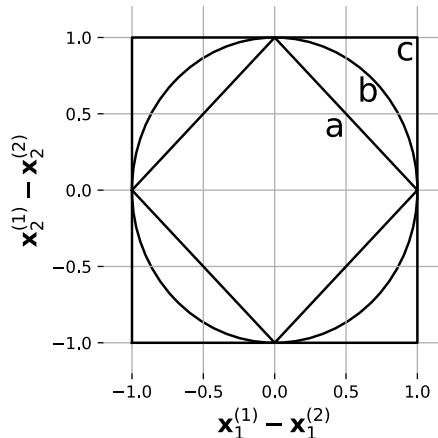
$$\mathbf{w}_{\text{new}}^* = (\mathbf{X}_{\text{new}}^T \mathbf{X}_{\text{new}})^{-1} \mathbf{X}_{\text{new}}^T \mathbf{y}_{\text{new}} = (\mathbf{X}_{\text{new}}^T \mathbf{X}_{\text{new}})^{-1} \mathbf{X}_{\text{new}}^T \frac{1}{\sigma} \mathbf{y}$$

$$\mathbf{w}_{\text{new}}^* = \mathbf{w}^* \Rightarrow (\mathbf{X}_{\text{new}}^T \mathbf{X}_{\text{new}})^{-1} \mathbf{X}_{\text{new}}^T \frac{1}{\sigma} \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{X}_{\text{new}} = \frac{1}{\sigma} \cdot \mathbf{X}$$

Problem 3: k-nearest neighbors (Version A) (3 credits)

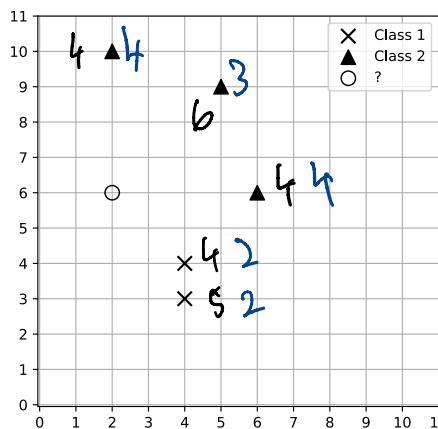
In the following figure we see the unit circles of three distance functions.



- 0 1 a) Assign each of the following three distance functions its corresponding unit circle (letter a-c) from the figure.

- L_2 -distance: $\|\mathbf{x}^{(1)} - \mathbf{x}^{(2)}\|_2 = \sqrt{\sum_i (\mathbf{x}_i^{(1)} - \mathbf{x}_i^{(2)})^2}$ → b
- L_1 -distance: $\|\mathbf{x}^{(1)} - \mathbf{x}^{(2)}\|_1 = \sum_i |\mathbf{x}_i^{(1)} - \mathbf{x}_i^{(2)}|$ → a
- L_∞ -distance: $\|\mathbf{x}^{(1)} - \mathbf{x}^{(2)}\|_\infty = \max_i |\mathbf{x}_i^{(1)} - \mathbf{x}_i^{(2)}|$ → c

In the following figure we see a two-dimensional dataset with two classes. We would like to classify the point (2, 6) marked with a circle using k-nearest-neighbors with $k = 3$.



- 0 1 b) What is the predicted class of the point when using the L_1 distance?

Class 2

- 0 1 c) What is the predicted class of the point when using the L_∞ distance?

Class 1

Problem 4: Classification (Version A) (6 credits)

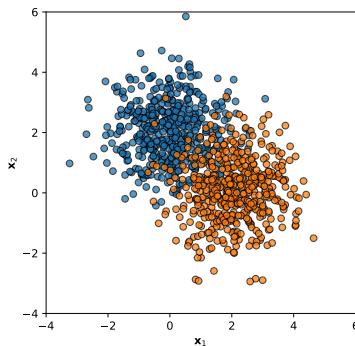
You are given a balanced dataset with two classes, i.e. $p(y = 0) = p(y = 1)$. Assume that the ground truth class conditional distributions are bivariate Gaussian distributions, i.e. $p(\mathbf{x} | c) = \mathcal{N}(\mathbf{x} | \mu_c, \Sigma_c)$ with mean μ_c and covariance Σ_c for each class $c \in \{0, 1\}$.

Further assume that we can choose between two models to fit the data:

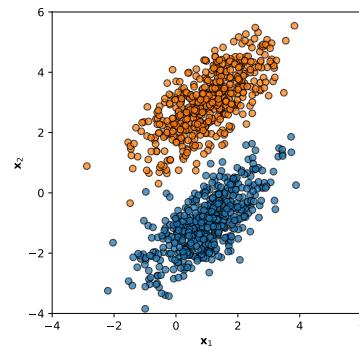
- Linear Discriminant Analysis with Gaussian class conditional distributions
- Naïve Bayes with Gaussian class conditional distributions

For each of the datasets shown below (a, b, c), choose one of the possible options (1,2,3) and justify your answer:

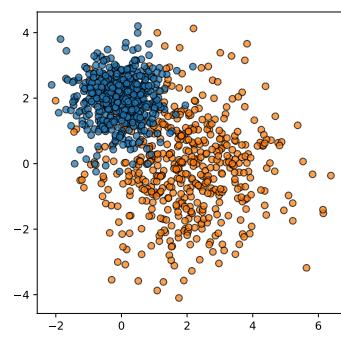
1. We should use Linear Discriminant Analysis.
2. We should use Naïve Bayes.
3. There is no clear reason to prefer one model over the other.



(a)



(b)



(c)

0
1
2

a) Two isotropic bivariate Gaussian distributions with identical covariance matrices

Inconclusive, since both models' assumptions match the data

0
1
2

b) Identical covariance matrices but correlated features

LDA since the covariance matrices match for both classes and the features are correlated

0
1
2

c) Covariance matrices do not match and no apparent correlation between features

Naïve Bayes since the features are conditionally independent and variances mismatch between classes \Rightarrow quadratic decision boundary

Problem 5: Optimization – Convexity (Version A) (10 credits)

Consider the two functions

$$f(\mathbf{x}) = \max_{i=1,\dots,n} x_i - \min_{i=1,\dots,n} x_i$$

$$g(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n |x_i - \text{median}(\mathbf{x})|$$

with $\mathbf{x} \in \mathbb{R}^n$. You may assume that n is odd.

0
1
2
3
4

- a)
Prove or disprove that $f(\mathbf{x})$ is convex in \mathbf{x} .

0
1
2
3
4

- b)
Prove or disprove that $g(\mathbf{x})$ is convex in \mathbf{x} .

Hint: $\text{median}(\mathbf{x}) = \arg \min_{t \in \mathbb{R}} \|\mathbf{x} - t\mathbf{1}\|_1$ with $\|\cdot\|_1$ being the sum over \mathbf{x} 's elements' absolute values.

0
1
2
3
4

a)
Consider the function $e_i(x) = x_i$ which is clearly convex in x since it is constant in all dimensions but the i -th, in which is linear. From the definition of convexity, it is easy to see that $e_i(\lambda x + (1-\lambda)y) \leq \lambda e_i(x) + (1-\lambda)e_i(y)$ holds
(By induction, rule (2) also holds for more than two arguments). Thus, by rule (2) $\max_{i=1,\dots,n} x_i$ is convex (plugging in $e_i(x)$ for $f_i(x)$)

Equivalently, $\min_{i=1,\dots,n} x_i$ is concave. By rule (4) it follows that $-\min_{i=1,\dots,n} x_i$ is convex

last, by rule (1) it follows that $f(\mathbf{x})$ is convex in \mathbf{x}

b)

For two arbitrary $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, by definition of $g(\mathbf{x})$: $t_x^* = \text{median}(\mathbf{x}) = \arg \min_t \|\mathbf{x} - t\mathbf{1}\|_1$ and $t_y^* = \text{median}(\mathbf{y}) = \arg \min_t \|\mathbf{y} - t\mathbf{1}\|_1$. $t_\lambda^* = \text{median}(\lambda \mathbf{x} + (1-\lambda)\mathbf{y})$
 $= \arg \min_t \|\lambda \mathbf{x} + (1-\lambda)\mathbf{y} - t\mathbf{1}\|_1$ for an arbitrary $\lambda \in [0,1]$

$$\lambda g(\mathbf{x}) + (1-\lambda)g(\mathbf{y}) = \lambda \frac{1}{n} \|\mathbf{x} - t_x^*\mathbf{1}\|_1 + (1-\lambda) \frac{1}{n} \|\mathbf{y} - t_y^*\mathbf{1}\|_1.$$

$$= \frac{1}{n} [\|\lambda \mathbf{x} - \lambda t_x^*\mathbf{1}\|_1 + \|(1-\lambda)\mathbf{y} - (1-\lambda)t_y^*\mathbf{1}\|_1]$$

$$\geq \frac{1}{n} \|\lambda \mathbf{x} - \lambda t_x^*\mathbf{1} + (1-\lambda)\mathbf{y} - (1-\lambda)t_y^*\mathbf{1}\|_1$$

$$\geq \frac{1}{n} \|\lambda \mathbf{x} + (1-\lambda)\mathbf{y} - t_\lambda^*\mathbf{1}\|_1$$

$$= g(\lambda \mathbf{x} + (1-\lambda)\mathbf{y})$$

triangle inequality

$$\|\mathbf{a}\| + \|\mathbf{b}\| \geq \|\mathbf{a} + \mathbf{b}\|$$

definition of t_λ^*

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

Suppose $\mathbf{x} \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^N$ are two vectors. We define the functions $f : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}^N$ and $g : \mathbb{R}^N \rightarrow \mathbb{R}$, and use them to compute

$$\begin{aligned}\mathbf{z} &= f(\mathbf{x}, \mathbf{y}) \\ t &= g(\mathbf{z}).\end{aligned}$$

The code below implements the computation of f and g , as well as its gradients using backpropagation. Your task is to complete the missing code fragments.

NOTE: The code is given in Python but you can write the solution in pseudocode as long as it is clear and unambiguous, making sure that the return values have correct shapes.

```
import numpy as np

class F:
    def forward(self, x, y):
        self.cache = (x, y)
        ##### MISSING CODE FRAGMENT #1 #####
        ##### MISSING CODE FRAGMENT #2 #####
        return out

    def backward(self, d_out):
        # x, y are arrays of shape (N,)
        x, y = self.cache
        out = x * np.sin(y) # d_out = x * sin(y)
        d_x = np.sin(y) * d_out
        d_y = x * np.cos(y) * d_out
        return d_x, d_y

class G:
    def forward(self, z):
        self.cache = z
        out = np.mean(z)
        return out

    def backward(self, d_out):
        # z is an array of shape (N,)
        z = self.cache
        ##### MISSING CODE FRAGMENT #2 #####
        # MISSING CODE FRAGMENT #2
        ##### MISSING CODE FRAGMENT #2 #####
        return d_z

# Example usage
f, g = F(), G()
x = np.array([1, 2, 3])
y = np.array([4, 5, 6])

z = f.forward(x, y)
t = g.forward(z)

d_z = g.backward(d_out=1.0)
d_x, d_y = f.backward(d_z)
```

$\text{out} = x \cdot \sin(y)$

$N = \text{len}(z)$

$d_z = \frac{\text{np.ones_like}(z)}{N} \cdot d_{\text{out}}$

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

Suppose $\mathbf{x} \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^N$ are two vectors. We define the functions $f : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}^N$ and $g : \mathbb{R}^N \rightarrow \mathbb{R}$, and use them to compute

$$\begin{aligned}\mathbf{z} &= f(\mathbf{x}, \mathbf{y}) \\ t &= g(\mathbf{z}).\end{aligned}$$

The code below implements the computation of f and g , as well as its gradients using backpropagation. Your task is to complete the missing code fragments.

NOTE: The code is given in Python but you can write the solution in pseudocode as long as it is clear and unambiguous, making sure that the return values have correct shapes.

```
import numpy as np

class F:
    def forward(self, x, y):
        self.cache = (x, y)
        ##### MISSING CODE FRAGMENT #1 #####
        ##### MISSING CODE FRAGMENT #1 #####
        return out

    def backward(self, d_out):
        # x, y are arrays of shape (N,)
        x, y = self.cache
        d_x = np.exp(x) / np.exp(y) * d_out
        d_y = -d_x
        return d_x, d_y

class G:
    def forward(self, z):
        self.cache = z
        out = np.sum(z)
        return out

    def backward(self, d_out):
        # z is an array of shape (N,)
        z = self.cache
        ##### MISSING CODE FRAGMENT #2 #####
        ##### MISSING CODE FRAGMENT #2 #####
        return d_z

# Example usage
f, g = F(), G()
x = np.array([1, 2, 3])
y = np.array([4, 5, 6])

z = f.forward(x, y)
t = g.forward(z)

d_z = g.backward(d_out=1.0)
d_x, d_y = f.backward(d_z)
```

$$d\hat{x} = \frac{e^x}{e^y}$$

$$\cancel{\frac{e^x \cdot e^y}{e^y \cdot e^x}}$$

$$\cancel{\frac{e^x}{e^y}}$$

$$N = \ln(z)$$

$$d_z = \frac{1}{z} \cdot d_{out}$$

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

Suppose $\mathbf{x} \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^N$ are two vectors. We define the functions $f : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}^N$ and $g : \mathbb{R}^N \rightarrow \mathbb{R}$, and use them to compute

$$\begin{aligned}\mathbf{z} &= f(\mathbf{x}, \mathbf{y}) \\ t &= g(\mathbf{z}).\end{aligned}$$

The code below implements the computation of f and g , as well as its gradients using backpropagation. Your task is to complete the missing code fragments.

NOTE: The code is given in Python but you can write the solution in pseudocode as long as it is clear and unambiguous, making sure that the return values have correct shapes.

```
import numpy as np

class F:
    def forward(self, x, y):
        self.cache = (x, y)
        ######
        # MISSING CODE FRAGMENT #1
        #####
        return out

    def backward(self, d_out):
        # x, y are arrays of shape (N,)
        x, y = self.cache
        temp = np.cos(x * y) * d_out
        d_x = y * temp
        d_y = x * temp
        return d_x, d_y

class G:
    def forward(self, z):
        self.cache = z
        out = np.prod(z) # Product of array elements
        return out

    def backward(self, d_out):
        # z is an array of shape (N,)
        z = self.cache
        ######
        # MISSING CODE FRAGMENT #2
        #####
        return d_z

# Example usage
f, g = F(), G()
x = np.array([1, 2, 3])
y = np.array([4, 5, 6])

z = f.forward(x, y)
t = g.forward(z)

d_z = g.backward(d_out=1.0)
d_x, d_y = f.backward(d_z)
```

$out = np.sin(x * y)$

$d_z = \frac{\partial \cdot prod(z)}{\partial z} \cdot d_{out}$

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

Suppose $\mathbf{x} \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^N$ are two vectors. We define the functions $f : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}^N$ and $g : \mathbb{R}^N \rightarrow \mathbb{R}$, and use them to compute

$$\begin{aligned}\mathbf{z} &= f(\mathbf{x}, \mathbf{y}) \\ t &= g(\mathbf{z}).\end{aligned}$$

The code below implements the computation of f and g , as well as its gradients using backpropagation. Your task is to complete the missing code fragments.

NOTE: The code is given in Python but you can write the solution in pseudocode as long as it is clear and unambiguous, making sure that the return values have correct shapes.

```
import numpy as np

class F:
    def forward(self, x, y):
        self.cache = (x, y)
        ##### # MISSING CODE FRAGMENT #1 #####
        ##### # MISSING CODE FRAGMENT #2 #####
        return out

    def backward(self, d_out):
        # x, y are arrays of shape (N,)
        x, y = self.cache
        d_x = (1 + y) * d_out
        d_y = x * d_out
        return d_x, d_y

class G:
    def forward(self, z):
        self.cache = z
        out = np.dot(z, z) # Dot product
        return out

    def backward(self, d_out):
        # z is an array of shape (N,)
        z = self.cache
        ##### # MISSING CODE FRAGMENT #1 #####
        ##### # MISSING CODE FRAGMENT #2 #####
        return d_z

# Example usage
f, g = F(), G()
x = np.array([1, 2, 3])
y = np.array([4, 5, 6])

z = f.forward(x, y)
t = g.forward(z)

d_z = g.backward(d_out=1.0)
d_x, d_y = f.backward(d_z)
```

$\text{out} = \mathbf{x}\mathbf{y} + \mathbf{x}$

$d_z = 2 * z * d_{out}$

Problem 7: Dimensionality reduction (Version A) (12 credits)

We would like to perform binary classification on a dataset (X, y) , where $X \in \mathbb{R}^{N \times D}$ and $y \in \{0, 1\}^N$.

Assume that we first reduce the dimensionality of X via PCA to obtain the matrix $\tilde{X} \in \mathbb{R}^{N \times K}$ (where $K < D$).

- 0 a) Suppose the original dataset (X, y) is linearly separable. Is the dataset (\tilde{X}, y) also guaranteed to be linearly separable? If yes, prove it. If not, provide a counterexample.

1 No



- 0 b) Suppose the original dataset (X, y) is NOT linearly separable. Is the dataset (\tilde{X}, y) guaranteed to NOT be linearly separable either?

- 2 • If yes (i.e. (\tilde{X}, y) is NOT linearly separable), prove it.
- 3 • If no (i.e. (\tilde{X}, y) may be linearly separable), provide a counterexample.

4 Yes (Proof by contradiction)

5 Assume that (X, y) was not linearly separable, but (\tilde{X}, y) was linearly separable

6 Let $T \in \mathbb{R}^{D \times K}$ be the top- K eigenvectors of the covariance matrix (the matrix we use to perform dimensionality reduction)

Since (\tilde{X}, y) is linearly separable, there is a linear classifier $I[w^T \tilde{x} \geq b]$ with parameters $w \in \mathbb{R}^K$ and $b \in \mathbb{R}$ that correctly classifies all points in (\tilde{X}, y)

Based on the definition of \tilde{X} ($\tilde{x} = T^T x$), this means that any point in (X) can be correctly classified by first performing dimensionality reduction with PCA and then applying the linear classifier specified above

However, this procedure in itself is a linear classifier $I[(w^T r^T) x \geq b]$ with weight vector $w' = T \cdot w$

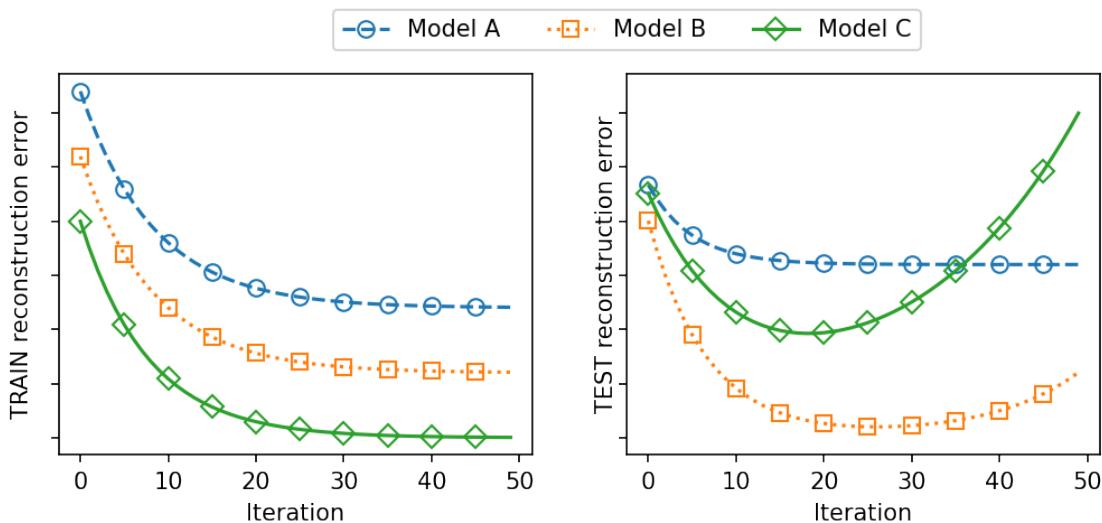
This contradicts our assumption that (X, y) is not linearly separable

Problem 8: Matrix factorization (Version A) (6 credits)

We would like to perform recommendation using matrix factorization. We have trained 3 latent factor models on the same dataset with gradient descent. These models are identical, except using a different value of k (number of latent factors):

- Model 1: $k = 5$
- Model 2: $k = 20$
- Model 3: $k = 50$

The figure below shows the reconstruction error for different models at each optimization step.



Your task is to assign the different models (1, 2, 3) to the loss curves in the figure above (A, B, C). Justify your answer.

A \rightarrow 5 (1) : Both the training and test losses are high - the model underfits (this may happen if k is too low)

B \rightarrow 20 (2) : by exclusion

C \rightarrow 50 (3) : The train loss is low, but the test loss increases sharply after a few iterations - the model overfits (this may happen if k is too high)

Problem 9: Clustering (Version A) (12 credits)

Consider the following mixture model with K components and a uniform prior over the components:

$$p(z_i = k) = \frac{1}{K} \quad p(\mathbf{x}_i | z_i = k, \mu_1, \dots, \mu_K) = \prod_{d=1}^D \frac{(\mu_{kd} x_{id})^{(x_{id}-1)} \exp(-\mu_{kd} x_{id})}{x_{id}!},$$

with parameters $\mu_k = (\mu_{k1}, \dots, \mu_{kD})^T \in [0, 1]^D$ for $k \in \{1, \dots, K\}$.

Suppose we are given a dataset consisting of N data points $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, where each data point is represented by a D -dimensional vector of positive natural number, that is $\mathbf{x}_i = (x_{i1}, \dots, x_{iD})^T \in \{1, 2, 3, \dots\}^D$.

Derive the M-step of the EM algorithm for the above mixture model, assuming that the responsibilities $\gamma_t(z_{ik})$ are given.

First, we write down the objective function

$$\mathbb{E}_{Z \sim Y_t(Z)} [\log p(\mathbf{x}, Z | \mu_1, \dots, \mu_K)] = \sum_{i=1}^N \sum_{k=1}^K \gamma_t(z_i=k) \log \left(\frac{1}{K} \prod_{d=1}^D \frac{(\mu_{kd} x_{id})^{(x_{id}-1)} \exp(-\mu_{kd} x_{id})}{x_{id}!} \right)$$

$$= C + \sum_{i=1}^N \sum_{k=1}^K \gamma_t(z_i=k) \sum_{d=1}^D ((x_{id}-1) \log \mu_{kd} - \mu_{kd} x_{id})$$

$$= C + \underbrace{\sum_{k=1}^K \sum_{d=1}^D \sum_{i=1}^N \gamma_t(z_i=k) ((x_{id}-1) \log \mu_{kd} - \mu_{kd} x_{id})}_{=: L_{kd}}$$

$$= C + \sum_{k=1}^K \sum_{d=1}^D L_{kd}$$

We see that the objective can be decomposed as a sum, where each parameter μ_{kd} is only encountered in a single term L_{kd}

Therefore, to find the optimal μ_{kd} , we only need to maximize the respective term L_{kd}

$$\begin{aligned} \frac{\partial L_{kd}}{\partial \mu_{kd}} &= \frac{\partial}{\partial \mu_{kd}} \left(\sum_{i=1}^N \gamma_t(z_i=k) ((x_{id}-1) \log \mu_{kd} - \mu_{kd} x_{id}) \right) \\ &= \sum_{i=1}^N \gamma_t(z_i=k) \left(\frac{x_{id}-1}{\mu_{kd}} - x_{id} \right) \stackrel{!}{=} 0 \end{aligned}$$

$$\mu_{kd} = \frac{\sum_{i=1}^N \gamma_t(z_i=k) (x_{id}-1)}{\sum_{i=1}^N \gamma_t(z_i=k) x_{id}}$$

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.

