



**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

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

**Date:** Saturday 11<sup>th</sup> July, 2020  
**Time:** 10:45 – 12:45

P 1	P 2	P 3	P 4	P 5	P 6	P 7	P 8	P 9	P 10	P 11	P 12
I											

## Working instructions

- This exam consists of **16 pages** with a total of **12 problems**.  
Please make sure now that you received a complete copy of the exam.
- The total amount of achievable credits in this exam is 55 credits.
- Allowed resources:
  - all materials that you will use on your own (lecture slides, calculator etc.)
  - **not allowed are any forms of collaboration between examinees and plagiarism**
- Only write on the provided sheets, **submitting your own additional sheets is not possible**.
- Last three pages can be used as scratch paper.
- All sheets (including scratch paper) have to be submitted to the upload queue. Missing pages will be considered empty.
- **Only use a black or blue color (no red or green)!**
- Write your answers only in the provided solution boxes or the scratch paper.
- **For problems that say "Justify your answer" you only get points if you provide a valid explanation.**
- **For problems that say "Prove" you only get points if you provide a valid mathematical proof.**
- If a problem does not say "Justify your answer" or "Prove" it's sufficient to only provide the correct answer.
- Exam duration - 120 minutes.

Left room from \_\_\_\_\_ to \_\_\_\_\_ / Early submission at \_\_\_\_\_

## Problem 1 KNN-Classification (4 credits)

0   
1   
2

- a) Assume you use a KNN-classifier on the following training data, that contains at least 100 samples of each class.

PS	acceleration	max. velocity [km/h]	cylinder capacity [cm <sup>3</sup> ]	weight [kg]	class
150	12.5	178	1968	2001	van
600	3.6	250	3996	2150	car
113	3.5	200	937	227	motorcycle
...	...	...	...	...	...

You observe that the obtained model performs bad on the test set. What might be the problem? Name at least two possible problems and explain how you would solve them.

- Underfitting → due to a small  $K$ , choosing a larger  $K$
- Overfitting → due to a large  $K$ , choosing a smaller  $K$   
problem: different range of features  
sol.: standardize the data:  $x_i = \frac{x_i - \mu_i}{\sigma_i}$
- problem: bad hyperparameter  $K$   
sol.: optimize hyperparameter  $K$  (grid-search)

0   
1   
2

- b) Would a decision tree have the same problems? Justify your answer.

Yes, we have to control the depth of the tree in order to not overfit. Also, it cannot be very small because then it can underfit.  
problem: different range of features → No  
Bad hyperparameters → no

## Problem 2 Overfitting (3 credits)

Explain overfitting. When does it occur? Why is overfitting unwanted? How can we spot overfitting? How can we avoid it?

0  
1  
2  
3

It occurs when the training error is decreasing and the test error after decreasing at the beginning starts increasing.

Is unwanted, because the model is not generalizing, the model is memorizing

By adding generalization ( $L_1, L_2$ )

Ocurs when we try to model the training data perfectly

### Problem 3 Probabilistic inference (7 credits)

0  
1  
2  
3  
4  
5  
6  
7

Consider the following probabilistic model

$$p(\lambda | a, b) = \text{Gamma}(\lambda | a, b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda)$$

$$p(x | \lambda) = \text{Poisson}(x | \lambda) = \frac{\lambda^x \exp(-\lambda)}{x!}$$

where  $a \in (1, \infty)$  and  $b \in (0, \infty)$ . We have observed a single data point  $x \in \mathbb{N}$ . Derive the maximum a posteriori (MAP) estimate of the parameter  $\lambda$  for the above probabilistic model. Show your work.

$$\begin{aligned}
 p(\lambda | x) &= p(x|\lambda) \cdot p(\lambda | a, b) \Rightarrow \\
 &= \log p(x|\lambda) + \log p(\lambda | a, b) \Rightarrow \\
 &= \log \left( \frac{\lambda^x \exp(-\lambda)}{x!} \right) + \log \left( \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda) \right) \\
 &= \log(\lambda^x) - \lambda + \log(\lambda^{a-1}) - b\lambda \Rightarrow \\
 &= x \log(\lambda) - \lambda + (a-1) \log(\lambda) - b\lambda \Rightarrow \\
 &= (x+a-1) \log(\lambda) - \lambda(1+b) \\
 \frac{\partial}{\partial \lambda} (x+a-1) \log(\lambda) - \lambda(1+b) &= \frac{x+a-1}{\lambda} - 1-b \\
 \frac{x+a-1}{\lambda} - 1-b &\stackrel{!}{=} 0 ; \quad \frac{x+a-1}{1-b} = \lambda \\
 \lambda_{\text{MAP}} &= \frac{x+a-1}{1+b}
 \end{aligned}$$

## Problem 4 Regression (5 credits)

- a) Assume you train a linear regression model on dataset  $D = \{\mathbf{x}_i, y_i\}_i$ ,  $\mathbf{x}_i \in \mathbb{R}^D$ ,  $y_i \in \mathbb{R}$  with the mean-square-error as loss function. After training is finished, you compute the MSE on individual data-points of the training-set. You notice that for three points you obtain a high MSE (1000 times higher than for the other points). Evaluation on the test-set shows that your regression model does not perform that well. What might be the reason for that? How would you improve performance of your model? Justify your answer.

0  
 1  
 2  
 3

Outliers in the dataset

$L_1 \rightarrow$  use different loss function

- b) You want to train another linear regression model and decide to use the log-cosh-loss:

0  
 1  
 2

$$E_{lc} = \sum_i \log \cosh(\mathbf{w}^T \mathbf{x}_i - y_i)$$

How do you learn the parameter  $\mathbf{w}$  of your model? Describe in one or two sentences.

Hint:  $\cosh(z) = 0.5(e^z + e^{-z})$

Gradient descent

Initialize  $\mathbf{w}$  randomly

Update until criterion is satisfied:  $\mathbf{w} = \mathbf{w} - t \nabla_{\mathbf{w}} f(\mathbf{w})$

## Problem 5 Classification (4 credits)

We would like to design a generative classification model for the following data.

Each data point is represented by a  $D$ -dimensional feature vector  $\mathbf{x} = (x_1, \dots, x_D)$ , where each entry  $x_j$  is a real number between 0 and 1, that is  $x_j \in [0, 1]$  for  $j = 1, \dots, D$ . Each data point belongs to one of  $K > 2$  possible classes, that is  $y \in \{1, \dots, K\}$ .

a) Which of the following distributions is the most reasonable choice for the class prior  $p(y)$ ?

- Bernoulli     Normal     Beta     Exponential     Categorical

b) We decide to model the class conditional distribution  $p(\mathbf{x}|y)$  as  $p(\mathbf{x}|y) = \prod_{j=1}^D p(x_j|y)$ , that is, the features  $x_j$  are conditionally independent given the class label  $y$ .

Which of the following distributions is the most reasonable choice for  $p(x_j|y)$ ?

- Categorical     Beta     Normal     Exponential     Bernoulli

0  c) What is the name of the posterior distribution  $p(y|\mathbf{x})$  for the model that you specified in subtasks (a) and (b)?  
1 Justify your answer.

2 Note that you need to provide the name of the distribution, not its probability density / mass functions. If the posterior distribution doesn't have a name, you should write "unknown distribution".

Since  $y$  can take  $K$  distinct values, the only possible choice is

Categorical

## Problem 6 Alternative characterization of vertices (4 credits)

Consider a non-empty convex set  $\mathcal{X} \subset \mathbb{R}^D$  and  $\mathbf{x} \in \mathcal{X}$ . Prove that if  $\mathbf{x}$  is a vertex of  $\mathcal{X}$  then  $\mathcal{X} \setminus \{\mathbf{x}\}$  is convex.

*Hint: additionally to the definition from the lecture you can use that  $\mathbf{x} \in \mathcal{X}$  is a vertex of  $\mathcal{X}$  if and only if for all  $\mathbf{x}_0, \mathbf{x}_1 \in \mathcal{X}$  with  $\mathbf{x}_0 \neq \mathbf{x}_1$  and all  $\lambda \in (0, 1)$  there holds that  $\mathbf{x} \neq \mathbf{x}_\lambda$ , where  $\mathbf{x}_\lambda = \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_0$  (i.e.  $\mathbf{x}$  does not lie between two different points from  $\mathcal{X}$ ).*

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

## Problem 7 Classification with Hinge loss and $L_\infty$ penalty (7 credits)

For  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^D$  and  $y_1, \dots, y_N \in \{-1, 1\}$  consider the following optimization problem with a fixed parameter  $\lambda > 0$  and  $\|\mathbf{w}\|_\infty = \max(|w_1|, \dots, |w_D|)$ .

$$\underset{\mathbf{w}, b}{\text{minimize}} \quad \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) + \lambda \|\mathbf{w}\|_\infty. \quad (1)$$

a) In this task you have to choose all correct options. Problem (1) as formulated above is

- |  |   |  |
|--|---|--|
| <input type="checkbox"/> concave.                            | <input type="checkbox"/> a quadratic problem.               | <input checked="" type="checkbox"/> unconstrained. |
| <input checked="" type="checkbox"/> not a quadratic problem. | <input checked="" type="checkbox"/> convex.                 | <input type="checkbox"/> constrained.              |
| <input type="checkbox"/> a linear problem.                   | <input checked="" type="checkbox"/> a minimization problem. | <input type="checkbox"/> non-convex.               |

0   
1   
2   
3   
4   
5  b) Reformulate problem (1) as an optimization problem with a linear objective and linear constraints. Justify your answer.

*Hint: you can introduce new variables to the problem.*

Non-linear Hinge loss can be removed by introducing new variables  $\epsilon \in \mathbb{R}^n$  and corresponding linear constraints

$$\underset{\mathbf{w}, b, \epsilon}{\text{min}} \quad \sum_{i=1}^n \epsilon_i + \lambda \|\mathbf{w}\|_\infty$$

$$\text{subject to } 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \leq \epsilon_i \text{ and } \epsilon_i \geq 0$$

## Problem 8 Deep learning (4 credits)

We are using a fully-connected neural network with 2 hidden layers for binary classification of points in  $\mathbb{R}^D$

$$f(\mathbf{x}, \mathbf{W}) = \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \sigma_0(\mathbf{W}_0 \mathbf{x}))).$$

where  $\mathbf{W} = \{\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2\}$  with  $\mathbf{W}_0 \in \mathbb{R}^{D_1 \times D}$ ,  $\mathbf{W}_1 \in \mathbb{R}^{D_2 \times D_1}$  and  $\mathbf{W}_2 \in \mathbb{R}^{1 \times D_2}$  are the weights of the neural network.

The neural network outputs probabilities of the positive class, i.e.  $p(y = 1 | \mathbf{x}, \mathbf{W}) = f(\mathbf{x}, \mathbf{W})$ , and is trained by minimizing the binary cross-entropy loss. We use the following activation functions:

$$\sigma_0(t) = t\sqrt{69}$$

$$\sigma_1(t) = -\frac{t}{54\pi}$$

$$\sigma_2(t) = \frac{1}{1 + \exp(-67t)}$$

The neural network achieves 100% classification accuracy on a dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ . Which of the following statements is true? Justify your answer.

1.  $\mathcal{D}$  is linearly separable.
2.  $\mathcal{D}$  is NOT linearly separable.
3. There is not enough information to determine if  $\mathcal{D}$  is linearly separable.

0
1
2
3
4

1)

Activation functions  $\sigma_0$  and  $\sigma_1$  are linear

$$f(\mathbf{x}, \mathbf{W}) = \sigma_2\left(W_2 \left(-\frac{1}{3\pi} I\right) W_1 (\sqrt{2} I) W_0 \mathbf{x}\right)$$

$$g(\mathbf{x}, \mathbf{W}) = \sigma_2(V\mathbf{x})$$

$$\text{where } V = W_2 \left(-\frac{1}{3} \pi I\right) W_1 (\sqrt{2} I) W_0 = -\frac{\sqrt{2}}{3} W_2 W_1 W_0$$

## Problem 9 Principal Component Analysis (4 credits)

Consider the data

$$\mathbf{X} = \begin{pmatrix} 0.37 & 0.95 & 0.73 & 0.60 \\ 0.16 & 0.16 & 0.06 & 0.87 \\ 0.60 & 0.71 & 0.02 & 0.97 \\ 0.83 & 0.21 & 0.18 & 0.18 \\ 0.30 & 0.52 & 0.43 & 0.29 \\ 0.61 & 0.14 & 0.29 & 0.37 \\ 0.46 & 0.79 & 0.20 & 0.51 \\ 0.59 & 0.05 & 0.61 & 0.17 \end{pmatrix},$$

where each row of  $\mathbf{X}$  represents a sample.

0      In each of the following PCA solutions the first row of  $\Gamma$  corresponds to the first principal component (associated with the first variance), the second row to the second, etc. Only one of these solutions is correct. Which one is it? For each wrong solution give a reason for why it is wrong!

1  
2  
3  
4

Variances	Principal component matrix $\Gamma$	Answer
$\begin{pmatrix} 0.16 \\ 0.10 \\ 0.05 \end{pmatrix}$	$\begin{pmatrix} 0.25 & -0.72 & 0.14 \\ 0.01 & 0.54 & 0.71 \\ -0.81 & -0.37 & 0.4 \\ -0.52 & 0.23 & -0.56 \end{pmatrix}$	Wrong. Principal components should be from $\mathbb{R}^4$
$\begin{pmatrix} 0.16 \\ 0.10 \\ 0.05 \end{pmatrix}$	$\begin{pmatrix} 0.25 & -0.72 & 0.14 & -0.63 \\ 0.01 & 0.54 & 0.71 & -0.46 \\ -0.81 & -0.37 & 0.41 & 0.18 \end{pmatrix}$	Correct
$\begin{pmatrix} 0.16 \\ -0.10 \\ 0.05 \\ 0.01 \end{pmatrix}$	$\begin{pmatrix} 0.25 & -0.72 & 0.14 & -0.63 \\ 0.01 & 0.54 & 0.71 & -0.46 \\ -0.81 & -0.37 & 0.41 & 0.18 \\ -0.52 & 0.23 & -0.56 & -0.60 \end{pmatrix}$	Wrong. Variances must be positive
$\begin{pmatrix} 0.16 \\ 0.05 \\ 0.10 \\ 0.01 \end{pmatrix}$	$\begin{pmatrix} 0.25 & -0.72 & 0.14 & -0.63 \\ -0.81 & -0.37 & 0.41 & 0.18 \\ 0.01 & 0.54 & 0.71 & -0.46 \\ -0.52 & 0.23 & -0.56 & -0.60 \end{pmatrix}$	Wrong. Variances are not monotonically decreasing
$\begin{pmatrix} 0.16 \\ 0.10 \\ 0.05 \end{pmatrix}$	$\begin{pmatrix} 0.25 & -0.72 & 0.14 & -0.63 \\ 0.01 & 0.54 & 0.71 & -0.46 \\ 0.50 & -0.72 & 0.14 & -0.63 \end{pmatrix}$	Wrong. Principal components are not orthogonal (1st and 3rd vectors)

## Problem 10 Mixture Models (1 credit)

Let  $z \sim \text{Cat}(\pi)$  be a random variable with categorical distribution on  $\{1, \dots, K\}$  with probabilities  $p(z = k) = \pi_k$  for  $k \in \{1, \dots, K\}$ . Furthermore, let  $x$  be a random variable dependent on  $z$  with an arbitrary likelihood, i.e.  $p(x | z)$  can be any probability distribution. Which of the following is the general form of  $p(z = k | x)$ ?

- $\pi_k \prod_{i=1}^K \log(p(x | z = i) \pi_i)$
- $p(x | z = k) \cdot \frac{\pi_k}{\sum_{i=1}^K \pi_i}$
- $p(z = k) \cdot \mathbb{E}_z [p(x | z)]$
- $p(x | z = k) \pi_k \left( \sum_{i=1}^K p(x | z = i) \pi_i \right)^{-1}$

## Problem 11 EM Algorithm (10 credits)

Consider a one-dimensional mixture of exponential distributions with  $K$  components and a uniform prior over components, i.e.

$$p(z_i = k) = \frac{1}{K} \quad p(x_i | \lambda_k, z_i = k) = \lambda_k \exp(-\lambda_k x_i) \quad \text{where } \lambda_k > 0.$$

We have observed  $N$  values  $x_i \in \mathbb{R}_{\geq 0}$  ( $i = 1 \dots N$ ) and want to fit this mixture model with the EM algorithm.

a) Derive the M-step, i.e. the responsibilities respectively the posterior  $\gamma(z_i = k) = p(z_i = k | x_i)$ .

The application of Bayes' theorem gave us

$$\begin{aligned} p(z_i = k | x_i, \lambda) &\propto p(x_i | \lambda_k, z_i = k) \cdot p(z_i = k) \\ &= \frac{1}{K} \cdot \lambda_k \exp(-\lambda_k x_i) \\ \gamma(z_i = k) &= p(z_i = k | x_i, \lambda) = \frac{\frac{1}{K} \cdot \lambda_k \exp(-\lambda_k x_i)}{\sum_{k=1}^K \frac{1}{K} \lambda_k \exp(-\lambda_k x_i)} \end{aligned}$$

0
1
2
3

0 b) Derive the E-step, i.e. find  $\arg \max_{\lambda} \mathbb{E}_{Z \sim \gamma} [\log p(Z, X | \lambda)]$ . Here  $Z$  represents all  $z_i$  and  $X$  all  $x_i$  ( $i = 1 \dots N$ ).

1  
2  
3  
4  
5

$$\begin{aligned} \mathbb{E}_{Z \sim \gamma} [\log p(Z, X | \lambda)] &= \sum_{i=1}^N \sum_{k=1}^K \gamma(z_i=k) \cdot \log P(z_i, x_i | \lambda) \\ &= \sum_{i=1}^N \sum_{k=1}^K \gamma(z_i=k) \log \frac{1}{\lambda_k} \exp(-\lambda_k x_i) \\ &= \sum_{i=1}^N \sum_{k=1}^K \gamma(z_i=k) (\log(\lambda_k) - \lambda_k x_i) + C \\ \underbrace{\sum_{z_i} \mathbb{E}_{Z \sim \gamma} [\log p(z_i, x_i | \lambda)]}_{\text{and find the root}} &= \sum_{i=1}^N \gamma(z_i=k) \cdot \left( \frac{1}{\lambda_k} - x_i \right) \end{aligned}$$

and find the root

$$\sum_{i=1}^N \gamma(z_i=k) \left( \frac{1}{\lambda_k} - x_i \right) = 0 \Leftrightarrow \lambda_k = \frac{\sum_{i=1}^N \gamma(z_i=k)}{\sum_{i=1}^N \gamma(z_i=k) \cdot x_i}$$

c) Is the EM algorithm guaranteed to converge to a global optimum in general? If yes, justify why. If no, how to avoid getting stuck in local optima or saddle points?

0  
1  
2

No

local optima can only be reached from very specific initial values. So this can for all intents and purposes be avoided with random initialization of the model parameters

### Problem 12 Differential Privacy (2 credits)

Let  $\mathcal{M}_f : \mathbb{R}^D \rightarrow \mathbb{R}^D$  be an  $\epsilon$  – DP mechanism with a privacy parameter  $\epsilon$  applied to the function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^D$ . Similarly, let  $\mathcal{N}_g : \mathbb{R}^D \rightarrow \mathbb{R}^D$  be a  $\sigma$  – DP mechanism with a privacy parameter  $\sigma$  applied to the function  $g$ .

Let  $h_1 : \mathbb{R}^D \rightarrow \mathbb{R}^D$  and  $h_2 : \mathbb{R}^D \rightarrow \mathbb{R}^D$  be arbitrary functions and  $\mathbf{X} \in \mathbb{R}^D$ . Can we provide differential privacy guarantees for the following mappings? If yes, what is their respective privacy parameter? If no, why not?

- a)  $\mathbf{X} \mapsto (\mathcal{M}_f(\mathbf{X}), \mathcal{N}_g(\mathbf{X}))$
- b)  $\mathbf{X} \mapsto h_1(\mathcal{N}_g(h_2(\mathbf{X})))$
- c)  $\mathbf{X} \mapsto h_2(\mathcal{M}_f(\mathbf{X}))$
- d)  $\mathbf{X} \mapsto (\mathcal{M}_f(h_1(\mathbf{X})), \mathcal{N}_g(\mathbf{X}))$

0  
1  
2

**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 students to write their solutions. The grid is composed of thin black lines forming small squares across the page.



