

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 for Graphs and Sequential Data

Exam: IN2323 / Endterm

Date: Tuesday 1st August, 2023

Examiner: Prof. Dr. Stephan Günnemann

Time: 11:00 – 12:15

	P 1	P 2	P 3	P 4	P 5	P 6	P 7
I							

Working instructions

- This exam consists of **12 pages** with a total of **7 problems**.
Please make sure now that you received a complete copy of the exam.
- The total amount of achievable credits in this exam is 35 credits.
- Detaching pages from the exam is prohibited.
- Allowed resources:
 - one A4 sheet of handwritten notes, two sides.
- **No other material (e.g. books, cell phones, calculators) is allowed!**
- Physically turn off all electronic devices, put them into your bag and close the bag.
- There is scratch paper at the end of the exam.
- Write your answers only in the provided solution boxes or the scratch paper.
- If you solve a task on the scratch paper, clearly reference it in the main solution box.
- All sheets (including scratch paper) have to be returned at the end.
- **Only use a black or a blue pen (no pencils, red or greens pens!)**
- **For problems that say “Justify your answer” you only get points if you provide a valid explanation.**
- **For problems that say “Derive” you only get points if you provide a valid mathematical derivation.**
- **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”, “Derive” or “Prove”, it is sufficient to only provide the correct answer.

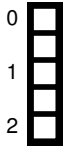
Left room from _____ to _____ / Early submission at _____

Problem 1 Hidden Markov Models (4 credits)

Consider a hidden Markov model with 3 states $\{1, 2, 3\}$ and 2 possible observations $\{a, b\}$. The initial distribution π , transition probabilities \mathbf{A} and emission probabilities \mathbf{B} are

$$\pi = \begin{matrix} & \begin{matrix} 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \begin{pmatrix} 1/2 \\ 1/2 \\ 0 \end{pmatrix} \end{matrix} \quad \mathbf{A} = \begin{matrix} & \begin{matrix} 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \begin{pmatrix} 1/2 & 0 & 1/2 \\ 0 & 1/2 & 1/2 \\ 1/3 & 1/3 & 1/3 \end{pmatrix} \end{matrix} \quad \mathbf{B} = \begin{matrix} & \begin{matrix} a & b \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \begin{pmatrix} 0 & 1 \\ 1/4 & 3/4 \\ 1 & 0 \end{pmatrix} \end{matrix},$$

where \mathbf{A}_{ij} specifies the probability of transitioning from state i to state j .



a) You have observed the sequence $X_{1:3} = [aba]$. Derive $\mathbb{P}(Z_3 \mid X_{1:3} = [aba])$ up to a normalizing constant.



b) What is the most likely state sequence $\arg \max_{Z_{1:3}} \mathbb{P}(Z_{1:3} \mid X_{1:3} = [aba])$? Justify your response.

Problem 2 Attention (8 credits)

Suppose we want to embed the sequence $S^{(X)} = [a, b, c, b, a, c]$ of length $N = 6$ over a vocabulary $\mathcal{V} = \{a, b, c\}$. We store the input sequence in a matrix $\mathbf{X} \in \{0, 1\}^{N \times |\mathcal{V}|}$, where the i -th row \mathbf{X}_i corresponds to a 1-hot representation of the i -th token. We use **masked-attention** to embed individual words which restricts the attention mechanism to a subset of other words in the sequence.

Masked attention is defined as:

$$\mathbf{Q} = \mathbf{X}\mathbf{W}^{(Q)} \quad \mathbf{K} = \mathbf{X}\mathbf{W}^{(K)} \quad \mathbf{V} = \mathbf{X}\mathbf{W}^{(V)}$$

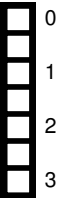
$$\mathbf{H} = \text{masked-softmax}_{\mathbf{M}}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$

When applying the softmax activation to obtain normalized attention scores, we ignore all entries where the mask $\mathbf{M} \in \{0, 1\}^{N \times N}$ has zero value:

$$\text{masked-softmax}_{\mathbf{M}}(\mathbf{A})_{i,j} = \begin{cases} \frac{\exp(\mathbf{A}_{i,j})}{\sum_{k: \mathbf{M}_{i,k} \neq 0} \exp(\mathbf{A}_{i,k})} & \text{if } \mathbf{M}_{i,j} \neq 0 \\ 0 & \text{otherwise} \end{cases}$$

Also, we use linear transformations $\mathbf{W}^{(Q)} \in \mathbb{R}^{|\mathcal{V}| \times d}$, $\mathbf{W}^{(K)} \in \mathbb{R}^{|\mathcal{V}| \times d}$, $\mathbf{W}^{(V)} \in \mathbb{R}^{|\mathcal{V}| \times d}$ to compute queries \mathbf{Q} , keys \mathbf{K} and values \mathbf{V} respectively.

a) For different realizations of attention masks, we now want to analyse which words are assigned identical embeddings. You are provided with different attention masks \mathbf{M} and your task is to find groups of input words that the corresponding masked-attention mechanism **can not** distinguish. That is, for each $\mathbf{M}^{(i)}$ list all groups of words of the sequence $S^{(X)}$ that will be assigned the same embedding no matter the choice of $\mathbf{W}^{(Q)}$, $\mathbf{W}^{(K)}$, $\mathbf{W}^{(V)}$. For example, if for $\mathbf{M}^{(i)}$ the first three tokens are assigned the same embedding and the last three tokens are assigned to the same embedding, your answer should be $\mathbf{M}^{(i)} : \{1, 2, 3\}, \{4, 5, 6\}$.



$$\mathbf{M}^{(1)} = \begin{pmatrix} 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 \end{pmatrix} \quad \mathbf{M}^{(2)} = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \quad \mathbf{M}^{(3)} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

b) Name and briefly explain a method that is employed in practice such that the words in any sequence can be distinguished from each other by the attention mechanism regardless of the choice of the attention mask \mathbf{M} (as long as each row $\mathbf{M}_{i,:}$ contains at least two non-zero elements).



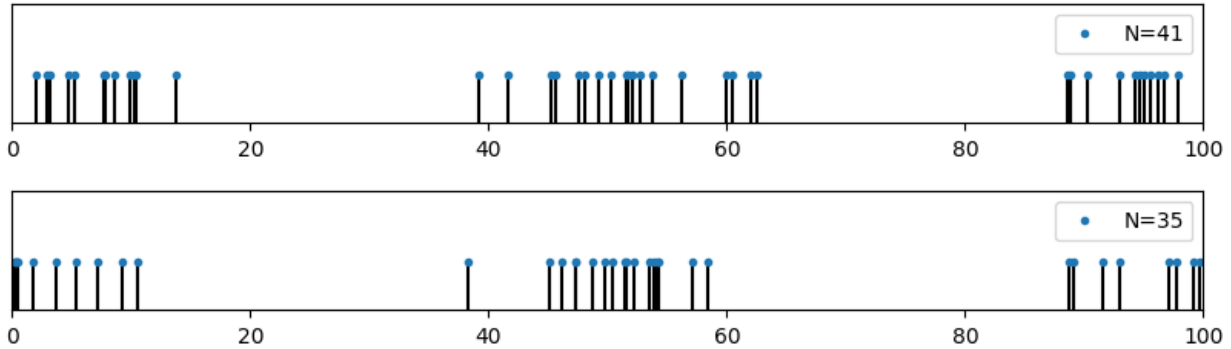
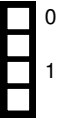
- 0 ☐ c) In practice, masked-attention is realized by setting unnormalized scores $\mathbf{A}_{i,j}$ to very low values wherever the attention mask $\mathbf{M}_{i,j} = 0$. Assume that the attention mask \mathbf{M} has ones everywhere except at position i, k where it is set to zero. Show that for all $j \in \{1, \dots, N\}$, we effectively recover the same behaviour as masking out the corresponding value $\mathbf{A}_{i,k}$ in the limit. That is, for all $j \in \{1, \dots, N\}$, show that $\lim_{\mathbf{A}_{i,k} \rightarrow -\infty} \text{softmax}(\mathbf{A})_{i,j} = \text{masked-softmax}_{\mathbf{M}}(\mathbf{A})_{i,j}$.
- 1 ☐
- 2 ☐
- 3 ☐

Problem 3 Temporal Point Processes (6 credits)

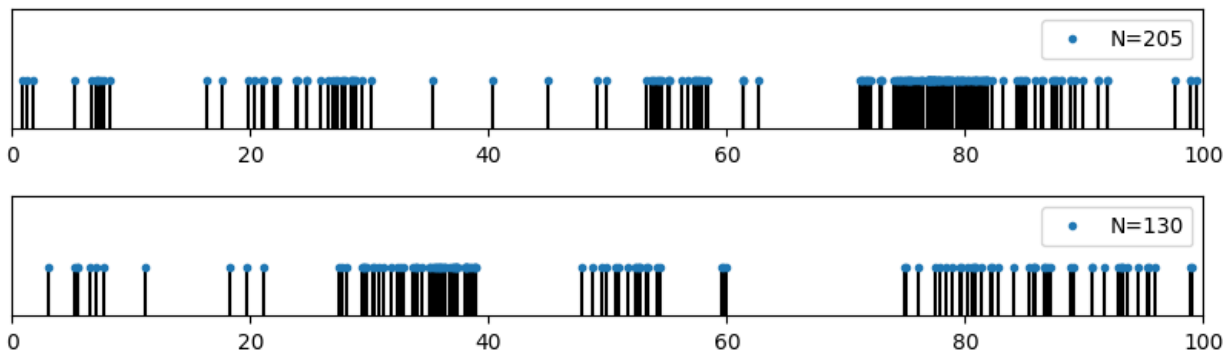
In the following, you are presented with five intensity functions of different temporal point processes. Your task is to subsequently match one of the intensity functions to the point process samples presented in the subtasks. Please justify your decision accordingly. Each intensity function can be only used once.

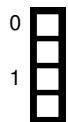
$$\begin{aligned}
 1) \lambda^*(t) &= 1.2 & 2) \lambda^*(t) &= \text{ReLU}(0.25 + \cos(\frac{1}{25}\pi t)) & 3) \lambda^*(t) &= 0.2 + 0.6 \sum_{t_i \in \mathcal{H}(t)} e^{-(t-t_i)} \\
 4) \lambda^*(t) &= 0.6 & 5) \lambda^*(t) &= 0.2 + 0.9 \sum_{t_i \in \mathcal{H}(t)} e^{-(t-t_i)}
 \end{aligned}$$

a) To which of the five intensity functions do the two samples most likely belong?

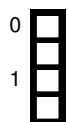
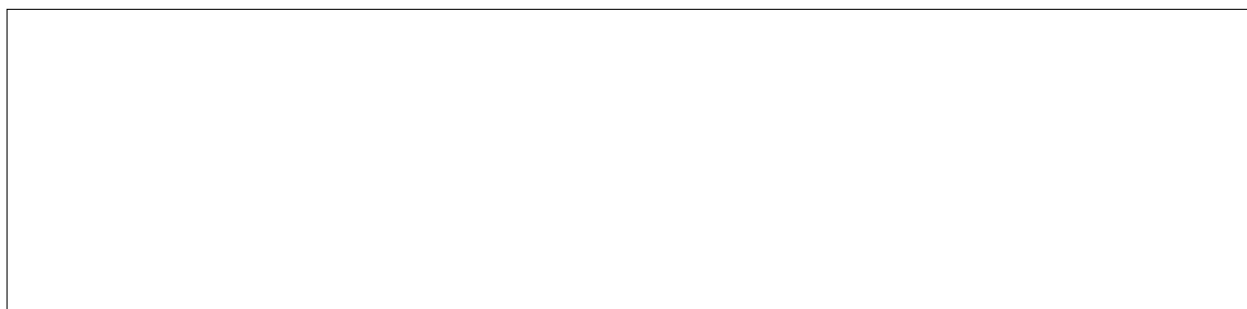
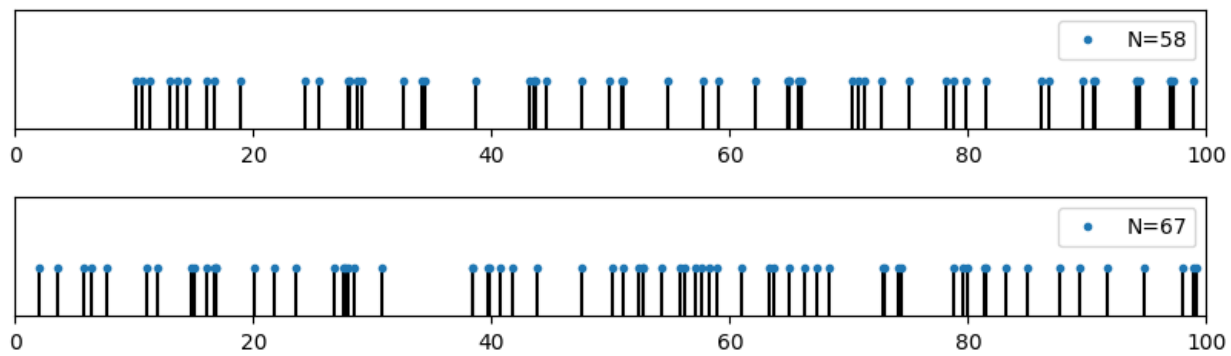


b) To which of the five intensity functions do the two samples most likely belong?

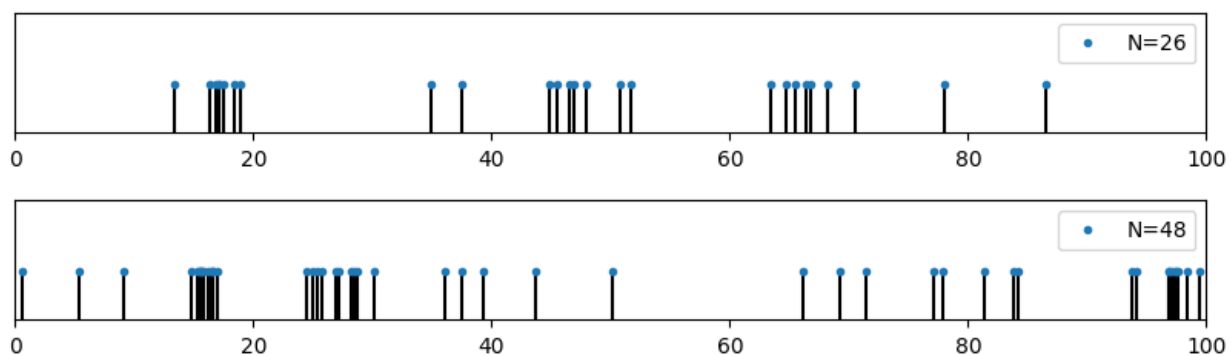




c) To which of the five intensity functions do the two samples most likely belong?



d) To which of the five intensity functions do the two samples most likely belong?



Problem 4 Spectral Clustering (5 credits)

a) Assume you are given a graph with 2 disconnected components of equal size and want to cluster the graph into 2 clusters. Provide an assignment vector \mathbf{f}_{C_1} for which it holds that

$$\mathbf{f}_{C_1}^T \mathbf{L} \mathbf{f}_{C_1} = 0 \quad \text{with } \|\mathbf{f}_{C_1}\|_2 > 0 \quad (4.1)$$

Justify.

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

b) Does the normalized cut minimize or maximize in-cluster associativity? Show your answer via the definition of the in-cluster associativity:

$$\text{assoc}(C_1, C_1) = \sum_{u \in C_1, v \in C_1} W_{uv} \quad (4.4)$$

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

Problem 5 Ranking (2 credits)

0 ☐ Assume you are given a function, which starts a random surfer without teleportation from some point v_i on a secret graph for some specified number of steps n . It returns the probability of finding the random surfer on each of the nodes after n steps. You observe that even for a sufficiently large number of steps n , you get significantly different outputs when you vary v_i . What can you say about the secret graph? How could you modify the random surfer function to solve this issue?

1 ☐

2 ☐

Problem 6 Equivariant Machine Learning on Graphs (5 credits)

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the feature matrix of an attributed graph with n nodes. We denote the feature vectors as $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$ (the columns of \mathbf{X}). Note that \mathbf{x}_i represents the i -th feature dimension of **all** nodes, $\mathbf{x}_i \in \mathbb{R}^n$.

A function f is called invariant to rescaling if

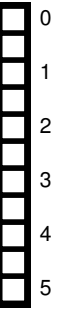
$$f(s_1 \mathbf{x}_1, \dots, s_d \mathbf{x}_d) = f(\mathbf{x}_1, \dots, \mathbf{x}_d) \quad s_i \in \mathbb{R} \setminus \{0\}$$

Consider the following function

$$g(\mathbf{x}_1, \dots, \mathbf{x}_d) = \psi \left(\bigotimes_{i=1}^d \left[\phi \left(\frac{\mathbf{x}_i}{\|\mathbf{x}_i\|} \right) + \phi \left(-\frac{\mathbf{x}_i}{\|\mathbf{x}_i\|} \right) \right] \right)$$

where ψ and ϕ are neural networks, \bigotimes an aggregation function such as sum, and $\|\cdot\|$ any ℓ_p -vector norm.

Show that g is invariant to rescaling.



Problem 7 Robustness - Discrete Randomized Smoothing (5 credits)

We want to certify our message-passing Graph Neural Network f_θ against edge perturbations using discrete randomized smoothing. We define the smoothed classifier for graphs $\mathbf{G} = (\mathbf{A}, \mathbf{X})$ as

$$g(\mathbf{G})_c = \Pr_{\phi}[f(\phi(\mathbf{A}), \mathbf{X}) = c]$$

where ϕ is the sparsity-aware smoothing distribution with edge deletion probability $p_d = \frac{1}{4}$ and edge addition probability $p_a = \frac{1}{2}$. Assume we know the adjacency matrix of the clean graph \mathbf{G} and the perturbed graph $\tilde{\mathbf{G}}$:

$$\mathbf{A} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \tilde{\mathbf{A}} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$$



a) Consider the two graphs \mathbf{G}_1 and \mathbf{G}_2 with the following adjacency matrices:

$$\mathbf{A}_1 = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \quad \mathbf{A}_2 = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$$

Compute the probabilities 1. $\Pr_{\phi}[\mathbf{A}_1 | \mathbf{A}]$, 2. $\Pr_{\phi}[\mathbf{A}_1 | \tilde{\mathbf{A}}]$, 3. $\Pr_{\phi}[\mathbf{A}_2 | \mathbf{A}]$, 4. $\Pr_{\phi}[\mathbf{A}_2 | \tilde{\mathbf{A}}]$.



b) To compute $\Pr_{\phi}[h(\phi(\tilde{\mathbf{A}}), \mathbf{X}) = c^*]$ under the worst-possible classifier h^* we have to select graphs that will be classified as c^* while ensuring $\Pr_{\phi}[h(\phi(\mathbf{A}), \mathbf{X}) = c^*] = g_{c^*}(\mathbf{G})$ for clean graph \mathbf{G} . Consider a classifier h with $h(\mathbf{G}_1) = c_{\text{other}}$ and $h(\mathbf{G}_2) = c^*$. Can h be a worst-case classifier? Why or why not? Justify your answer. *Hint:* If you did not solve the previous exercise you can use $\Pr_{\phi}[\mathbf{A}_1 | \mathbf{A}] = 0.2$, $\Pr_{\phi}[\mathbf{A}_1 | \tilde{\mathbf{A}}] = 0.1$, $\Pr_{\phi}[\mathbf{A}_2 | \mathbf{A}] = 0.1$, $\Pr_{\phi}[\mathbf{A}_2 | \tilde{\mathbf{A}}] = 0.2$.



c) Are the flipping probabilities $p_d = \frac{1}{4}$ and $p_a = \frac{1}{2}$ a good choice in practice? Why or why not?

[illegible]

