**TUM**

# Machine Learning for Graphs and Sequential Data

| **Exam:** | IN2323 / Endterm | **Date:** | Tuesday 1$^{st}$ 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 $\pi$, transition probabilities $A$ and emission probabilities $B$ are

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

where $A_{ij}$ specifies the probability of transitioning from state $i$ to state $j$.

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

> We can apply the forwards algorithm:
> $$\alpha_1 = \pi \odot B_{:,a} \propto \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$
> $$\alpha_2 = B_{:,b} \odot (A^T \alpha_1) \propto \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$
> $$\alpha_3 = B_{:,a} \odot (A^T \alpha_2) = \begin{pmatrix} 0 \\ 1/4 \\ 1 \end{pmatrix} \odot \begin{pmatrix} 0 \\ 1/2 \\ 1/2 \end{pmatrix} = \begin{pmatrix} 0 \\ 1/8 \\ 1/2 \end{pmatrix}$$
>
> We could also make the same argument for the first two calculations in prose: At $t = 1$, we are either in state 1 or 2. But $a$ can only be observed from state 2. At $t = 2$ we transition into either state 2 or 3. But $b$ can only be observed from state 2. At $t = 3$ we transition into state 2 or 3 with equal probability, but $a$ is more likely to be observed in state 3 by a factor of 4.

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

> From our calculations above, we know that $Z_{1:2} = [2, 2]$. The most likely choice for $Z_3$ is state 3.

# 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 $\boldsymbol{X} \in \{0, 1\}^{N \times |\mathcal{V}|}$, where the $i$-th row $\boldsymbol{X}_i$ 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:

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

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

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

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

Also, we use linear transformations $\boldsymbol{W}^{(Q)} \in \mathbb{R}^{|\mathcal{V}| \times d}$, $\boldsymbol{W}^{(K)} \in \mathbb{R}^{|\mathcal{V}| \times d}$, $\boldsymbol{W}^{(V)} \in \mathbb{R}^{|\mathcal{V}| \times d}$ to compute queries $\boldsymbol{Q}$, keys $\boldsymbol{K}$ and values $\boldsymbol{V}$ 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 $\boldsymbol{M}$ and your task is to find groups of input words that the corresponding masked-attention mechanism **can not** distinguish. That is, for each $\boldsymbol{M}^{(i)}$ list all groups of words of the sequence $S^{(X)}$ that will be assigned the same embedding no matter the choice of $\boldsymbol{W}^{(Q)}$, $\boldsymbol{W}^{(K)}$, $\boldsymbol{W}^{(V)}$. For example, if for $\boldsymbol{M}^{(i)}$ the first three tokens are assigned the same embedding and the last three tokens are assigned to the same embedding, your answer should be $\boldsymbol{M}^{(i)} : \{1, 2, 3\}, \{4, 5, 6\}$.

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

> $\boldsymbol{M}^{(1)} : \{1\}, \{2, 4\}, \{3\}, \{5\}, \{6\}$
> $\boldsymbol{M}^{(2)} : \{1, 5\}, \{2, 4\}, \{3, 6\}$
> $\boldsymbol{M}^{(3)} : \{1, 2, 3, 4, 5, 6\}$ (normalized attention score will always be 1.0)

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 $\boldsymbol{M}$ (as long as each row $\boldsymbol{M}_{i,:}$ contains at least two non-zero elements).

> Positional Encodings: Append / add a fixed embedding to each token representation that uniquely identifies its position in the sequence.

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

$$\text{softmax}(A)_{i,j} = \frac{\exp(A_{i,j})}{\sum_l \exp(A_{i,l})} = \frac{\exp(A_{i,j})}{\exp(A_{i,k}) + \sum_{l \neq k} \exp(A_{i,l})}$$

If $j = k$:

$$\lim_{A_{i,k} \to -\infty} \text{softmax}(A)_{i,k} = \frac{\lim_{A_{i,k} \to -\infty} \exp(A_{i,k})}{\lim_{A_{i,k} \to -\infty} \exp(A_{i,k}) + \sum_{l \neq k} \exp(A_{i,l})}$$

$$= \frac{0}{0 + \sum_{l \neq k} \exp(A_{i,l})} = 0$$

If $j \neq k$:

$$\lim_{A_{i,k} \to -\infty} \text{softmax}(A)_{i,j} = \frac{\lim_{A_{i,j} \to -\infty} \exp(A_{i,j})}{\lim_{A_{i,k} \to -\infty} \exp(A_{i,k}) + \sum_{l \neq k} \exp(A_{i,l})}$$

$$= \frac{\exp(A_{i,j})}{0 + \sum_{l \neq k} \exp(A_{i,l})}$$

$$= \frac{\exp(A_{i,j})}{\sum_{l : M_{i,l} \neq 0} \exp(A_{i,l})} = \text{masked-softmax}_M(A)_{i,j}$$

# 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.
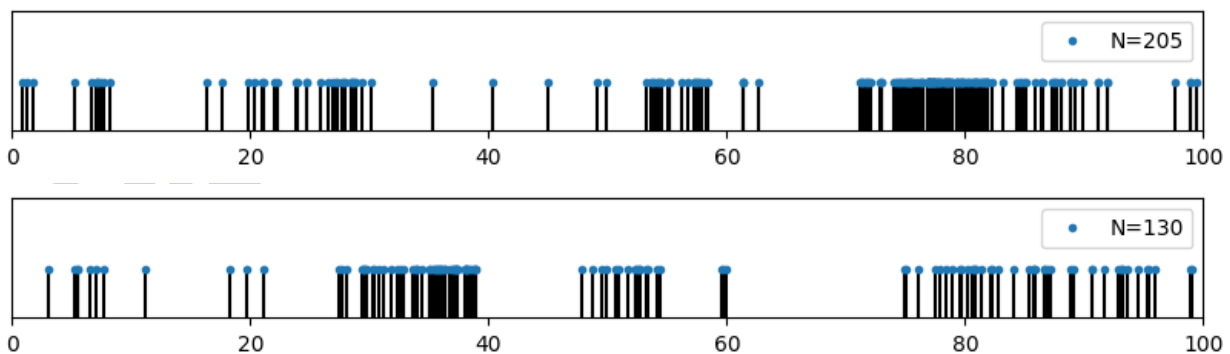
**1)** $\lambda^*(t) = 1.2$    **2)** $\lambda^*(t) = ReLU(0.25 + \cos(\frac{1}{25}\pi t))$    **3)** $\lambda^*(t) = 0.2 + 0.6 \sum_{t_i \in \mathcal{H}(t)} e^{-(t-t_i)}$

**4)** $\lambda^*(t) = 0.6$    **5)** $\lambda^*(t) = 0.2 + 0.9 \sum_{t_i \in \mathcal{H}(t)} e^{-(t-t_i)}$

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

The correct intensity is the cosine intensity 2).
The presented samples distinctly belong to a TPP that is not homogenous, as the intensity varies in time. Furthermore, the two samples show very similar patterns in time, indicating an inhomogeneous TPP. Lastly, we can see that the pattern matches the cosine function.
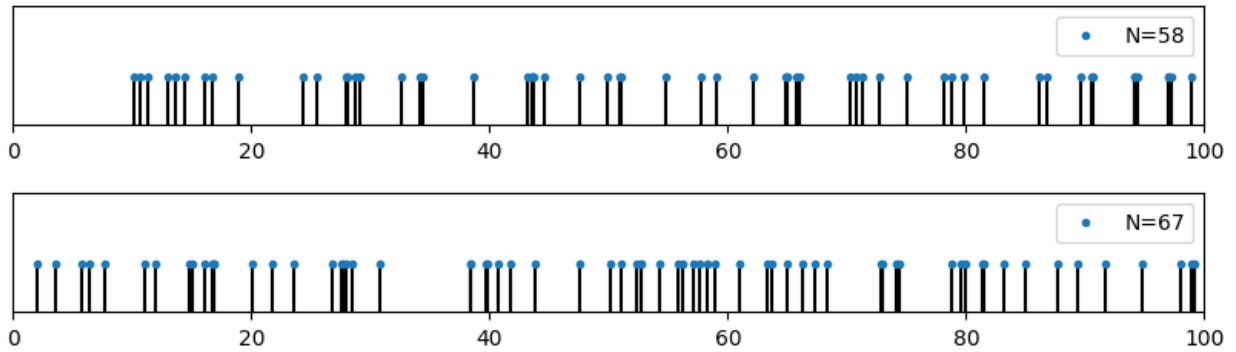
b) To which of the five intensity functions do the two samples most likely belong?

The correct intensity is the Hawkes intensity 5).
The two samples exhibit excitation (bursting of events) with periods of no events in between, which is a distinct property of the Hawkes process. As these two samples show stronger excitation when compared to d), they belong to the Hawkes process with the higher alpha value of 0.9.

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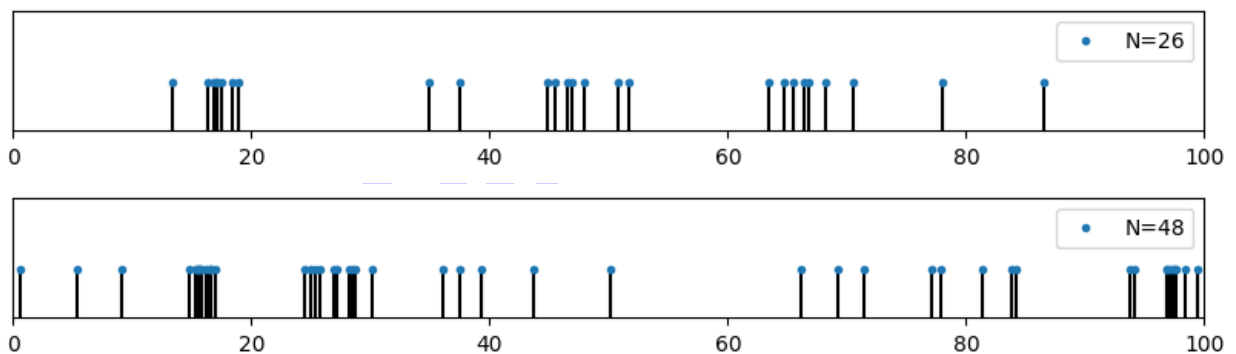




The correct intensity is the homogenous intensity 4).
The two samples exhibit a relatively homogenous event distribution, indicating an HPP. For an HPP, the expected number of events in an interval $[0, T]$ is given by $\lambda T$. Therefore, the intensity function 4) better matches the observed number of events than an HPP with intensity 1).

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





The correct intensity is the Hawkes intensity 3).
Again, the two samples exhibit excitation (bursting of events) with periods of no events in between, which is a distinct property of the Hawkes process. As these two samples show lower excitation than b), they belong to the Hawkes process with a lower alpha value of 0.6.

# 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

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

Justify.

☐ 0
☐ 1
☐ 2
☐ 3

Assuming we have the first connected component $A_1$ and the second connected component $A_2$. An assignment vector

$$f_{C_1,i} = \begin{cases} k, & \text{if } v_i \in A_1 \\ -k & \text{otherwise} \end{cases} \tag{4.2}$$

simply assigns $A_1 = C_1$ and $A_2 = C_2$. $k$ can be any constant. We can expand the given equation into

$$
\begin{aligned}
\mathbf{f}_{C_1}^T L \mathbf{f}_{C_1} = &= \frac{1}{2} \sum_{(u,v) \in E} W_{uv}(f_{C_1,v} - f_{C_1,v})^2 \\
&= \frac{1}{2} \sum_{\substack{(u,v) \in E \\ (u,v) \in C_1}} W_{uv}(f_{C_1,v} - f_{C_1,u})^2 + \frac{1}{2} \sum_{\substack{(u,v) \in E \\ (u,v) \in C_2}} W_{uv}(f_{C_1,v} - f_{C_1,u})^2 + \frac{1}{2} \sum_{\substack{(u,v) \in E \\ u \in C_1, v \in C_2}} W_{uv}(f_{C_1,v} - f_{C_1,u})^2 \\
&= \frac{1}{2} \sum_{\substack{(u,v) \in E \\ u \in C_1, v \in C_2}} W_{uv}(f_{C_1,v} - f_{C_1,u})^2 = 0
\end{aligned}
\tag{4.3}
$$

We can see that both sums over in-cluster edges result in 0, leaving us only with the sum over between-cluster edges. There are no between-cluster edges due to the fact that the two clusters, per assignment and definition of graph, share no edge.

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

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

☐ 0
☐ 1
☐ 2

The normalized cut maximizes in-cluster associativity!
We can rewrite the associativity as

$$\text{assoc}(C_1, C_1) = \sum_{u \in C_1, v \in C_1} W_{uv} = \sum_{u \in C_1, v \in V} W_{uv} - \sum_{u \in C_1, v \in \overline{C_1}} W_{uv} = \text{vol}(C_1) - \text{cut}(C_1, C_2) \tag{4.5}$$

We know that the normalized cut tries to minimize the cut while maximizing the volume of the clusters. Full points are also given for a solution via the graph laplacian.

## Problem 5   Ranking (2 credits)

0

1

2

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?

The graph could be:

- periodic

- disconnected

- more than one spider trap

This can be solved by Page Rank with Teleportation.

# Problem 6   Equivariant Machine Learning on Graphs (5 credits)

Let $X \in \mathbb{R}^{n \times d}$ be the feature matrix of an attributed graph with $n$ nodes. We denote the feature vectors as $x_1, \ldots, x_d \in \mathbb{R}^n$ (the columns of $X$). Note that $x_i$ represents the $i$-th feature dimension of **all** nodes, $x_i \in \mathbb{R}^n$.

A function $f$ is called invariant to rescaling if

$$f(s_1 x_1, \ldots, s_d x_d) = f(x_1, \ldots, x_d) \qquad s_i \in \mathbb{R} \setminus \{0\}$$

Consider the following function

$$g(x_1, \ldots, x_d) = \psi \left( \bigotimes_{i=1}^{D} \left[ \phi \left( \frac{x_i}{||x_i||} \right) + \phi \left( -\frac{x_i}{||x_i||} \right) \right] \right)$$
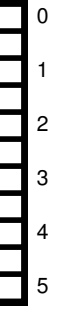
where $\psi$ and $\phi$ are neural networks, $\bigotimes$ an aggregation function such as sum, and $|| \cdot ||$ any $\ell_p$-vector norm.

Show that g is invariant to rescaling.

For any $s_1, \ldots, s_d \in \mathbb{R} \setminus \{0\}$ we have

$$
\begin{aligned}
g(s_1 x_1, \ldots, s_d x_d) &= \psi \left( \bigotimes_{i=1}^{D} \left[ \phi \left( \frac{s_i x_i}{||s_i x_i||} \right) + \phi \left( -\frac{s_i x_i}{||s_i x_i||} \right) \right] \right) \\
&= \psi \left( \bigotimes_{i=1}^{D} \left[ \phi \left( \frac{s_i x_i}{|s_i| \cdot ||x_i||} \right) + \phi \left( -\frac{s_i x_i}{|s_i| \cdot ||x_i||} \right) \right] \right) \\
&= \psi \left( \bigotimes_{i=1}^{D} \left[ \phi \left( \operatorname{sgn}(s_i) \frac{x_i}{||x_i||} \right) + \phi \left( -\operatorname{sgn}(s_i) \frac{x_i}{||x_i||} \right) \right] \right) \\
&\stackrel{(1)}{=} \psi \left( \bigotimes_{i=1}^{D} \left[ \phi \left( \frac{x_i}{||x_i||} \right) + \phi \left( -\frac{x_i}{||x_i||} \right) \right] \right) \\
&= g(x_1, \ldots, x_d)
\end{aligned}
$$

where (1) is due to sign invariance of the function $h(x) = \phi(x) + \phi(-x)$ for any vector $x \in \mathbb{R}^d$.
Since $g(s_1 x_1, \ldots, s_d x_d) = g(x_1, \ldots, x_d)$ we have that $g$ is invariant to rescaling. $\square$

## Problem 7  Robustness - Discrete Randomized Smoothing (5 credits)

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

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

where $\phi$ 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 $G$ and the perturbed graph $\tilde{G}$:

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

a) Consider the two graphs $G_1$ and $G_2$ with the following adjacency matrices:

$$A_1 = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \qquad A_2 = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$$

Compute the probabilities    1. $\Pr_\phi[A_1 \mid A]$,   2. $\Pr_\phi[A_1 \mid \tilde{A}]$,   3. $\Pr_\phi[A_2 \mid A]$,   4. $\Pr_\phi[A_2 \mid \tilde{A}]$.

1. $\Pr_\phi[A_1 \mid A] = (1 - p_d)^2 p_a^2 = (\frac{3}{4})^2(\frac{1}{2})^2 = \frac{9}{64}$

2. $\Pr_\phi[A_1 \mid \tilde{A}] = (1 - p_d)p_a^3 = \frac{3}{4}(\frac{1}{2})^3 = \frac{3}{32} = \frac{6}{64}$

3. $\Pr_\phi[A_2 \mid A] = p_d^2(1 - p_a)^2 = (\frac{1}{4})^2(\frac{1}{2})^2 = \frac{1}{64}$

4. $\Pr_\phi[A_2 \mid \tilde{A}] = p_d(1 - p_a)^3 = \frac{1}{4}(\frac{1}{2})^3 = \frac{2}{64}$

b) To compute $\Pr_\phi \left[ h(\phi(\tilde{A}), X) = c^* \right]$ under the worst-possible classifier $h^*$ we have to select graphs that will be classified as $c^*$ while ensuring $\Pr_\phi [h(\phi(A), X) = c^*] = g_{c^*}(G)$ for clean graph $G$. Consider a classifier $h$ with $h(G_1) = c_{\text{other}}$ and $h(G_2) = c^*$. 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 $\Pr_\phi[A_1 \mid A] = 0.2$, $\Pr_\phi[A_1 \mid \tilde{A}] = 0.1$, $\Pr_\phi[A_2 \mid A] = 0.1$, $\Pr_\phi[A_2 \mid \tilde{A}] = 0.2$.

No. In the worst-case either both graphs are classified the same or we would classify $G_1$ as $c^*$ and $G_2$ as $c_{\text{other}}$ since:

- The first graph $G_1$ is more likely to be sampled from $G$ than from $\tilde{G}$. The second graph $G_2$ is more likely to be sampled from $\tilde{G}$ than from $G$. Thus we would first classify $G_1$ as $c^*$.

- Alternative solution: The likelihood ratios are as follows: $\frac{\Pr_\phi[A_1|A]}{\Pr_\phi[A_1|\tilde{A}]} = \frac{9}{6} = \frac{3}{2} > \frac{1}{2} = \frac{\Pr_\phi[A_2|A]}{\Pr_\phi[A_2|\tilde{A}]}$. Higher likelihood ratio for $G_1$ thus we would first classify $G_1$ as $c^*$.

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?

In practice one would select a lower edge insertion probability $p_a$ to preserve the sparsity of the graph.

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