

**Esolution**

Place student sticker here

**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:** Friday 19<sup>th</sup> August, 2022

**Examiner:** Prof. Dr. Stephan Günnemann

**Time:** 08:15 – 09:30

P 1      P 2      P 3      P 4      P 5      P 6      P 7      P 8      P 9      P 10

--	--	--	--	--	--	--	--	--	--

### Working instructions

- This exam consists of **16 pages** with a total of **10 problems**.  
Please make sure now that you received a complete copy of the exam.
- The total amount of achievable credits in this exam is 86 credits.
- Detaching pages from the exam is prohibited.
- Allowed resources:
  - one A4 sheet of handwritten notes (two sides, not digitally written and printed).
- **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 (after problem 10).
- 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 Normalizing flows (8 credits)

You are given the task of density estimation on  $\mathbb{R}^2$  and plan on using normalizing flows. In the following we present some candidate transformations that will be used for **reverse parameterization**. For each of the transformations, state if it can be used to define a normalizing flow and justify your answers.

In all cases, the input is a vector  $\mathbf{x} = [x_1 \ x_2]^T$ . We denote the output of the transformation with  $\mathbf{z} \in \mathbb{R}^2$ .

0 ☐  
1 ☐  
2 ☐

a)

$$\mathbf{A} = \mathbf{W}^T \mathbf{W}$$

$$\mathbf{z} = \mathbf{A} \mathbf{x},$$

where  $\mathbf{W} \in \mathbb{R}^{2 \times 2}$ .

No.

The transformation may be non-invertible (consider  $m\mathbf{W} = 0$ ).

0 ☐  
1 ☐  
2 ☐

b)

$$\mathbf{z} = [x_1^2 \ x_2^2]^T$$

No.

The transformation is not a bijection.

0 ☐  
1 ☐  
2 ☐

c)

$$z_1 = \mathbf{V} \text{ReLU}(\mathbf{W} x_2 + \mathbf{b})$$

$$\mathbf{z} = [z_1 \ x_2]^T,$$

where  $\mathbf{W} \in \mathbb{R}^{h \times 1}$ ,  $\mathbf{V} \in \mathbb{R}^{1 \times h}$ ,  $\mathbf{b} \in \mathbb{R}^h$  and ReLU is applied elementwise.

No.

This is similar to a Real NVP architecture as the inverse of  $x_2$  is easy to obtain by assigning  $x_2 = z_2$ . However, the  $x_1 \mapsto z_1$  does not depend on  $x_1$  and all the points are mapped to a single point (which depends on  $x_2$ ). So the whole transformation is not a bijection.

d)

$$\mathbf{z} = \mathbf{a} \odot \mathbf{x} + \mathbf{b},$$

where  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^2$  and  $\odot$  is the elementwise product.

Yes if  $\mathbf{a} \neq \mathbf{0}$ .



Sample Solution

## Problem 2 Variational inference (10 credits)

Suppose we are given a latent variable model for a sequence of observations  $x_1, \dots, x_N \in \{0, 1\}$  and latent variables  $z_1, \dots, z_N \in [0, 1]$  with

$$p(z_1, \dots, z_N) = \prod_{n=1}^N \text{Beta}(z_n \mid \alpha, \beta) = \prod_{n=1}^N \frac{1}{B(\alpha, \beta)} z_n^{\alpha-1} (1 - z_n)^{\beta-1}$$

$$p(x_1, \dots, x_N \mid z_1, \dots, z_N) = \prod_{n=1}^N \text{Bern}(x_n \mid z_n) = \prod_{n=1}^N z_n^{x_n} (1 - z_n)^{1-x_n}$$

with parameters  $\alpha, \beta > 0$  and normalizing constant  $B(\alpha, \beta)$ . We define the variational distribution

$$q(z_1, \dots, z_N) = \prod_{n=1}^N \text{Beta}(z_n \mid \gamma, \delta) = \prod_{n=1}^N \frac{1}{B(\gamma, \delta)} z_n^{\gamma-1} (1 - z_n)^{\delta-1}$$

with parameters  $\gamma, \delta > 0$ .

Assume that  $\alpha, \beta$  are known and fixed. Prove or disprove the following statement:

There **exist** observations  $x_1, \dots, x_N \in \{0, 1\}$  and values of  $\gamma, \delta > 0$  such that the ELBO is tight, i.e.  $\exists x_1, \dots, x_N, \exists \gamma, \delta : \log(p(x_1, \dots, x_N)) = \mathcal{L}((\alpha, \beta), (\gamma, \delta))$ .

In the following, we use bold  $\mathbf{x}$  and  $\mathbf{z}$  for the two sequences.

The ELBO is tight if  $\text{KL}(p(\mathbf{z} \mid \mathbf{x}) \parallel q(\mathbf{z})) = 0$ , i.e.  $p(\mathbf{z} \mid \mathbf{x}) = q(\mathbf{z})$ .

We begin by determining  $p(\mathbf{z} \mid \mathbf{x})$ .

Applying Bayes formula shows that

$$p(\mathbf{z} \mid \mathbf{x}) \propto p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z}) \quad (2.1)$$

$$\propto \left( \prod_{n=1}^N z_n^{x_n} (1 - z_n)^{1-x_n} \right) \left( \prod_{n=1}^N \frac{1}{B(\alpha, \beta)} z_n^{\alpha-1} (1 - z_n)^{\beta-1} \right) \quad (2.2)$$

$$= \prod_{n=1}^N z_n^{\alpha-1+x_n} (1 - z_n)^{\beta-1+(1-x_n)} \quad (2.3)$$

$$= \prod_{n=1}^N z_n^{\alpha-1+x_n} (1 - z_n)^{\beta-x_n} \quad (2.4)$$

$$\propto \prod_{n=1}^N \text{Beta}(z_n \mid \alpha + x_n, \beta - x_n + 1) \quad (2.5)$$

Alternatively, we can use the fact that the Beta distribution is the conjugate prior of the Bernoulli distribution to immediately arrive at the above result.

If all observations  $x_n$  have the same value, i.e.  $\forall n : x_n = c$  with  $c \in \{0, 1\}$ , we can ensure that  $p(\mathbf{z} \mid \mathbf{x}) = q(\mathbf{z})$  by setting  $\gamma = \alpha + c$  and  $\delta = \beta - c + 1$ .

### Problem 3 Robustness (9 credits)

In the lecture, we have derived a convex relaxation for the ReLU activation function. Now, we want to generalize this result to the flexible ReLU (FReLU) activation function

$$FReLU(x) = \begin{cases} x + b & \text{if } x > 0 \\ b & \text{if } x \leq 0 \end{cases}$$

with variable input  $x \in \mathbb{R}$  and **constant parameter**  $b \in \mathbb{R}$ .

Let  $y \in \mathbb{R}$  be the variable we use to model the function's output. Now, given input bounds  $l, u \in \mathbb{R}$  with  $l \leq x \leq u$ , provide a set of **linear constraints** corresponding to the convex hull of  $\{[x \ FReLU(x)]^T \mid l \leq x \leq u\}$ .

**Hint:** You will have to make a case distinction to account for different ranges of  $l$  and  $u$ .

We have to distinguish three cases:

Case I:  $l, u \leq 0$

A single constraint suffices:  $y = b$

Case II:  $l, u > 0$

A single constraint suffices:  $y = x + b$

Case III:  $l < 0, u > 0$

Now, we need three linear constraints:

$$y \geq b$$

$$y \geq x + b$$

$$y \leq \frac{u}{u-l}(x-l) + b$$

	0
	1
	2
	3
	4
	5
	6
	7
	8
	9

## Problem 4 Autoregressive models (8 credits)

0	
1	
2	
3	
4	
5	

a)

An autoregressive process of order  $p$ ,  $AR(p)$ , is defined as:

$$X_t = c + \sum_{i=1}^p \varphi_i X_{t-i} + \varepsilon_t,$$

with independently distributed noise variables  $\varepsilon_t \sim \mathcal{N}(0, \sigma^2)$ .

Provided that the  $AR(p)$  is stationary, derive its first moment  $\mathbb{E}[X_t]$  as a function of  $c$  and  $\varphi_i$ .

$$\begin{aligned} \mathbb{E}[X_t] &= \mathbb{E}\left[c + \sum_{i=1}^p \varphi_i X_{t-i} + \varepsilon_t\right] && \text{Linearity} \\ \mathbb{E}[X_t] &= c + \sum_{i=1}^p \varphi_i \mathbb{E}[X_{t-i}] + \mathbb{E}[\varepsilon_t] && \mathbb{E}[\varepsilon_t] = 0 \forall t \text{ by definition} \\ \mathbb{E}[X_t] &= c + \sum_{i=1}^p \varphi_i \mathbb{E}[X_{t-i}] + 0 && \mathbb{E}[X_t] = \mu \forall t \text{ by stationarity} \\ \mu &= c + \sum_{i=1}^p \varphi_i \mu \\ \mu &= c + \mu \sum_{i=1}^p \varphi_i \\ \mu &= \frac{c}{1 - \sum_{i=1}^p \varphi_i} \end{aligned}$$

0	
1	
2	
3	

b) Let us define a process  $X_t$  as

$$X_t = \sin^2\left(-\frac{\pi}{2}t\right) + \frac{2}{3}X_{t-1} + \varepsilon_t$$

with independently distributed noise variables  $\varepsilon_t \sim \mathcal{N}(0, \sigma^2)$ .

Decide if the process  $X_t$  is stationary. Justify your answer.

We have

$$\mathbb{E}[X_t] = \mathbb{E}\left[\sin^2\left(-\frac{\pi}{2}t\right)\right] + \mathbb{E}\left[\frac{2}{3}X_{t-1}\right] + \mathbb{E}[\varepsilon_t] = \sin^2\left(-\frac{\pi}{2}t\right) + \mathbb{E}\left[\frac{2}{3}X_{t-1}\right]$$

Assume that the process was stationary, i.e.  $\mathbb{E}[X_t] = \mu \forall t$  with constant  $\mu$ . Then, we would have

$$\begin{aligned} \mu &= \sin^2\left(-\frac{\pi}{2}t\right) + \frac{2}{3}\mu \\ \Leftrightarrow \frac{1}{3}\mu &= \sin^2\left(-\frac{\pi}{2}t\right) \end{aligned}$$

This contradicts our assumption that  $\mu$  is a constant.

## Problem 5 Hidden Markov Models (9 credits)

Consider a hidden Markov model with 2 states  $\{1, 2\}$  and 4 possible observations  $\{c, e, i, n\}$ . The initial distribution  $\pi$ , transition probabilities  $\mathbf{A}$  and emission probabilities  $\mathbf{B}$  are

$$\pi = \begin{matrix} 1 \\ 2 \end{matrix} \begin{pmatrix} 2/5 \\ 3/5 \end{pmatrix} \quad \mathbf{A} = \begin{matrix} & 1 & 2 \\ 1 & \begin{pmatrix} 1/3 & 2/3 \end{pmatrix} \\ 2 & \begin{pmatrix} 3/5 & 2/5 \end{pmatrix} \end{matrix} \quad \mathbf{B} = \begin{matrix} & c & e & i & n \\ 1 & \begin{pmatrix} 2/5 & 1/5 & 1/5 & 1/5 \end{pmatrix} \\ 2 & \begin{pmatrix} 0 & 1/5 & 3/5 & 1/5 \end{pmatrix} \end{matrix},$$

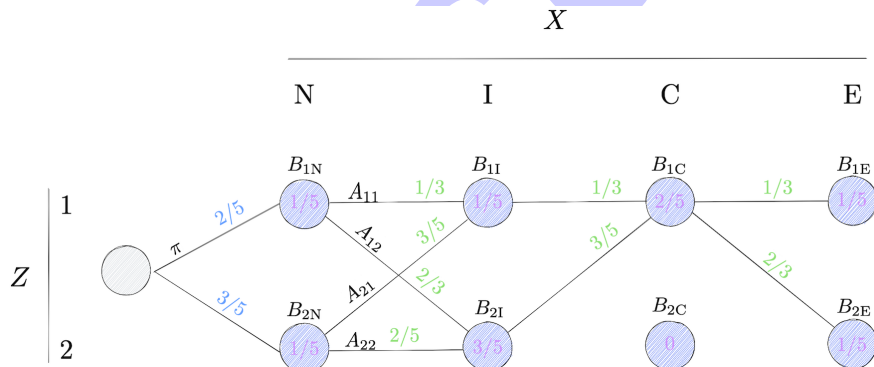
where  $\mathbf{A}_{ij}$  specifies the probability of transitioning from state  $i$  to state  $j$ .

a) We have observed the sequence  $X_{1:3} = [\text{nic}]$ . What is the most likely latent state  $Z_3$  given these observations? Justify your answer. What is this type of inference called?

This type of inference is called *filtering / smoothing*. To compute  $P(Z_3|X_{1:3})$  we can either use the forward algorithm or identify from the emission matrix  $\mathbf{B}$  that “c” can only be observed from latent state  $Z = 1$ . Hence, we get:

$$P(Z_3|X_{1:3}) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

b) The full observed sequence is  $X_{1:4} = [\text{nice}]$ . What is the most likely latent state sequence  $Z_{1:4}$  given these observations? Justify your answer.



- 1 We note that “c” can only be observed from latent state  $Z = 1$ . Hence, all paths going through  $B_{2C}$  can be ignored.
- 2 Since all paths have to go through  $B_{1C}$  and we have  $B_{1E} = B_{2E}$ , the maximising path has to also go along  $B_{2E}$  as the edge weight  $A_{12} = 2/3$  is larger than  $A_{11} = 1/3$ . Hence,  $Z_{3:4} = [12]$ .
- 3 We note that the probabilities for  $B_N$  are identical. Hence, the paths to  $B_{1I}$  only depend on the edge weights and the path along  $Z_1 = 2$  is favoured, its weight is  $\propto 9/25$ .
- 4 We note that the probabilities for  $B_N$  are identical. Hence, the paths to  $B_{2I}$  only depend on the edge weights and the path along  $Z_1 = 1$  is favoured, its weight is  $\propto 4/15$ .
- 5 Comparing the paths to  $B_{1C}$ , we multiply the paths weights for  $B_{1I}$  and  $B_{2I}$  with  $1/5 \times 1/3$  and  $3/5 \times 3/5$ , respectively. We get  $9/375$  for  $Z_{1:3} = [211]$  and  $36/375$  for  $Z_{1:3} = [121]$ . Hence,  $Z_{1:3} = [121]$  is more likely.

We conclude that the most likely latent state sequence  $\max_Z P(Z|X_{1:4} = [\text{nice}]) = [1212]$ .

## Problem 6 Temporal Point Processes (10 credits)

Assume that we use a Hawkes process to model a discrete event sequence  $\{t_1, \dots, t_N\}$  with  $t_i \in [0, T]$ . Further assume that (like in the lecture) we use an exponential triggering kernel, i.e.  $k_\omega(t - t_i) = \exp(-\omega(t - t_i))$ . Prove that the log-likelihood-function of the process is

$$\log p_\theta(\{t_1, \dots, t_N\}) = \sum_{i=1}^N \log \left( \mu + \alpha \sum_{j < i} \exp(-\omega(t_i - t_j)) \right) - \mu T + \frac{\alpha}{\omega} \sum_{i=1}^N (\exp(-\omega(T - t_i)) - 1)$$

The Hawkes process uses the conditional density

$$\lambda^*(t) = \mu + \alpha \sum_{t_j \in \mathcal{H}(t)} k_\omega(t - t_j), \quad (6.1)$$

with  $k_\omega(t - t_j) = \exp(-\omega(t - t_j))$ . The log-likelihood of a TPP is

$$\log p_\theta(\{t_1, \dots, t_N\}) = \sum_{i=1}^N \log \lambda^*(t_i) - \int_0^T \lambda^*(u) du. \quad (6.2)$$

We first insert the conditional density into the sum in Eq. 6.2:

$$\sum_{i=1}^N \log \lambda^*(t_i) = \sum_{i=1}^N \log \left( \mu + \alpha \sum_{j < i} \exp(-\omega(t_i - t_j)) \right). \quad (6.3)$$

The result is identical to the first summand of the log-likelihood specified in the problem statement. Next, we evaluate the integral. Since  $\lambda^*(u)$  changes after each event, we have to split the integral into a sum of segments. Doing so, we obtain

$$\begin{aligned} \int_0^T \lambda^*(u) du &= \int_0^{t_1} \mu du + \sum_{i=1}^N \int_{t_i}^{t_{i+1}} \mu + \alpha \sum_{j \leq i} \exp(-\omega(t - t_j)) du \\ &= \mu T + \alpha \sum_{i=1}^N \sum_{j \leq i} \int_{t_i}^{t_{i+1}} \exp(-\omega(t - t_j)) du \\ &= \mu T - \frac{\alpha}{\omega} \sum_{i=1}^N \sum_{j \leq i} (\exp(-\omega(t_{i+1} - t_j)) - \exp(-\omega(t_i - t_j))) \\ &\stackrel{(1)}{=} \mu T - \frac{\alpha}{\omega} \sum_{i=1}^N (\exp(-\omega(T - t_i)) - \exp(-\omega(t_i - t_i))) \\ &= \mu T - \frac{\alpha}{\omega} \sum_{i=1}^N (\exp(-\omega(T - t_i)) - 1) \end{aligned} \quad (6.4)$$

where we denote  $t_{N+1} = T$  for convenience and in (1) we used that the double sum cancels out. Overall, we arrive at

$$\log p_\theta(\{t_1, \dots, t_N\}) = \sum_{i=1}^N \log \left( \mu + \alpha \sum_{j < i} \exp(-\omega(t_i - t_j)) \right) - \mu T + \frac{\alpha}{\omega} \sum_{i=1}^N (\exp(-\omega(T - t_i)) - 1). \quad (6.5)$$



## Problem 7 Graphs – Generative Models (8 credits)

Let  $\mathbf{A} \in \{0, 1\}^{N \times N}$  be the adjacency matrix of a graph generated by a stochastic block model with  $\pi = \begin{bmatrix} a & 1-a \end{bmatrix}^T$ ,  $\eta = \begin{bmatrix} p & q \\ q & p \end{bmatrix}$  and parameters  $a, p, q \in [0, 1]$ . Let  $\deg(n) = \sum_{j=1}^N A_{n,j}$  be the degree of node  $n$ .

Derive the expected degree  $\mathbb{E}[\deg(n)]$  of an arbitrary node  $n$ .

In the following, we assume that we have no self-loops, i.e.  $A_{n,n} = 0$  (but allowing self-loops is also fine). Without loss of generality assume that  $n = N$ . Due to linearity of expectation, we have

$$\mathbb{E}[\deg] = \sum_{j=1}^{N-1} \mathbb{E}[A_{N,j}] = \sum_{j=1}^{N-1} \Pr[A_{N,j} = 1],$$

where the last equality follows from the fact that  $A_{N,j}$  is a binary random variable.

We can distinguish four possible community assignments for nodes  $N$  and  $j$  and apply the law of total probability:

$$\begin{aligned} \Pr[A_{N,j} = 1] &= \Pr[A_{N,j} = 1 \mid z_N = 1 = z_j] \Pr[z_N = 1 = z_j] \\ &\quad + \Pr[A_{N,j} = 1 \mid z_N = 2 = z_j] \Pr[z_N = 2 = z_j] \\ &\quad + \Pr[A_{N,j} = 1 \mid z_N = 1 \wedge z_j = 2] \Pr[z_N = 1 \wedge z_j = 2] \\ &\quad + \Pr[A_{N,j} = 1 \mid z_N = 2 \wedge z_j = 1] \Pr[z_N = 2 \wedge z_j = 1] \\ &= p \cdot (a^2 + (1-a)^2) + q \cdot (2 \cdot a(1-a)). \end{aligned}$$

and thus

$$\mathbb{E}[\deg] = (N-1) (p \cdot (a^2 + (1-a)^2) + q \cdot (2 \cdot a(1-a))).$$

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3
<input type="checkbox"/>	4
<input type="checkbox"/>	5
<input type="checkbox"/>	6
<input type="checkbox"/>	7
<input type="checkbox"/>	8

## Problem 8 Graphs – Clustering (10 credits)

Let  $\mathbf{A} \in \{0, 1\}^{N \times N}$  be the adjacency matrix of an undirected graph (i.e. symmetric adjacency matrix) generated by a stochastic block model with  $\pi = [a \quad 1 - a]^T$ ,  $\eta = \begin{bmatrix} p & q \\ q & p \end{bmatrix}$  and parameters  $a, p, q \in [0, 1]$ .

- 0 ☐ a) Assume that  $p, q, a \in [0, 1]$  are known and fixed. Can  $\Pr(\mathbf{z} \mid \mathbf{A}, \eta, \pi)$ , the probability mass function of  
1 ☐ community assignments  $\mathbf{z}$  given  $\mathbf{A}$ , be evaluated in polynomial time? That is, can it be evaluated in  $\mathcal{O}(N^c)$ ,  
2 ☐ where  $N$  is the number of nodes and  $c \in \mathbb{R}_+$ ? Justify your answer.

No. We have

$$\Pr(\mathbf{z} \mid \mathbf{A}, \eta, \pi) = \frac{\Pr(\mathbf{A} \mid \mathbf{z}, \eta) \Pr(\mathbf{z} \mid \pi)}{\Pr(\mathbf{A} \mid \eta, \pi)} = \frac{\Pr(\mathbf{A} \mid \mathbf{z}, \eta) \Pr(\mathbf{z} \mid \pi)}{\sum_{\mathbf{z}' \in \{0, 1\}^N} \Pr(\mathbf{A} \mid \mathbf{z}', \eta) \Pr(\mathbf{z}' \mid \pi)}$$

and evaluating the denominator requires exponential time.

- 0 ☐ b) Now assume that  $p = 0$ . Further assume that  $\mathbf{A}$  is a connected graph (i.e. each pair of nodes  $(i, j)$  is  
1 ☐ connected by a path). Propose a procedure for finding the most likely community assignment, i.e.  
2 ☐  
3 ☐  
4 ☐  
5 ☐  
6 ☐  
7 ☐  
8 ☐

$$\max_{\mathbf{z} \in \{0, 1\}^N} \Pr(\mathbf{z} \mid \mathbf{A}, \eta, \pi)$$

in polynomial time  $\mathcal{O}(N^c)$ . Justify your answers.

Because  $p = 0$ ,  $\mathbf{A}$  must be a bipartite graph, i.e.  $\mathbb{V} = \mathbb{H} \cup \mathbb{J}$ .

Because  $\mathbf{A}$  is connected, the partitions can be found in polynomial time, e.g. by starting depth-first-search at an arbitrary node, and adding nodes at even depth to  $\mathbb{H}$  and at odd depth to  $\mathbb{J}$ .

Because  $\mathbf{A}$  is connected and bipartite, there are only two possible community assignments: Either  $\mathbf{z}_n = 1 \iff n \in \mathbb{H}$  or  $\mathbf{z}_n = 1 \iff n \in \mathbb{J}$ .

Define  $\hat{\mathbf{z}} \in \{0, 1\}^N$  with  $\hat{\mathbf{z}}_n = 1 \iff n \in \mathbb{H}$ .

Then

$$\Pr(\hat{\mathbf{z}} \mid \mathbf{A}, \eta, \pi) = \frac{\Pr(\mathbf{A} \mid \hat{\mathbf{z}}, \eta) \Pr(\hat{\mathbf{z}} \mid \pi)}{\Pr(\mathbf{A} \mid \hat{\mathbf{z}}, \eta) \Pr(\hat{\mathbf{z}} \mid \pi) + \Pr(\mathbf{A} \mid 1 - \hat{\mathbf{z}}, \eta) \Pr(1 - \hat{\mathbf{z}} \mid \pi)}$$

If the probability is  $\geq 0.5$ , we have found our most likely community assignment. Otherwise,  $1 - \hat{\mathbf{z}}$  is the most likely community assignment.

The above probability can be computed in polynomial time.

## Problem 9 Limitations of Graph Neural Networks (6 credits)

a) What is oversmoothing and what causes it?

Oversmoothing refers to the effect that in the limit of infinite layers (under some assumptions), the embedding vectors of each node (in a connected component) converge to the same constant vector and hence, the nodes will get indistinguishable w.r.t. their attributes.

This is caused by the fact that in this limit, the influence of all nodes on an arbitrary node follows the PageRank solution and hence, the final representation is only dependent on the global graph structure and not on the local neighbourhood anymore.

0  
1  
2  
3

b) The Personalized Propagation of Neural Predictions (PPNP) architecture is designed to overcome the problem of oversmoothing. Briefly explain its two key building blocks.

- First, it separates the feature transformation from the propagation along the graph structure. It does so by applying a learnable function on the features of each node, e.g. a neural network before the propagation step.
- Second, for the propagation step, it uses personalized pagerank applied on the transformed features of each node. This means for each node a personalized teleportation probability is introduced and only the particular node is in the teleportation set.

0  
1  
2  
3

## Problem 10 Page Rank (8 credits)

Recall the spam farm discussed in our exercise. It consists of the spammer's own pages  $S_{\text{own}}$  with target page  $t$  and  $k$  supporting pages, as well as links from the accessible pages  $S_{\text{acc}}$  to the target page. **Different from the exercise, every page within  $S_{\text{own}}$  has a link to every other page within  $S_{\text{own}}$**  (see Fig. 10.1).

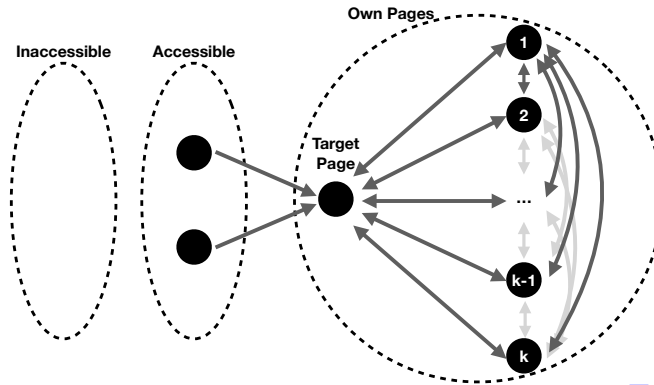


Figure 10.1

Let  $n$  be the total number of pages on the web,  $E$  the set of all edges,  $r_p$  the PageRank score of a page  $p$  and  $d_p$  the degree of a page  $p$ . Let  $x_{\text{acc}} = \sum_{p \in S_{\text{acc}}, (p,t) \in E} \frac{r_p}{d_p}$  be the amount of PageRank contributed by the accessible pages. We are using PageRank with teleports, where  $(1 - \beta)$  is the teleport probability.

a) Derive the PageRank  $r_s$  of a supporting page  $r_s$  as a function of  $\beta$ ,  $r_t$ ,  $k$ ,  $n$ .

The PageRank score of a supporting page consists of the PageRank coming from the target page  $r_t$  and all supporting pages:

$$r_s = \beta \frac{r_t}{k} + \beta * \frac{1}{k} * \sum_{i=1}^{k-1} r_i + \frac{1 - \beta}{n} \quad (10.1)$$

We can recognize that the supporting pages are all symmetric, i.e. have the same edges. Thus, we can get rid of the sum with:

$$r_s = \beta \frac{r_t}{k} + \beta * \frac{k-1}{k} * r_s + \frac{1 - \beta}{n} \quad (10.2)$$

Lastly, we solve for  $r_s$ :

$$r_s = r_t \frac{\beta}{k - \beta k + \beta} + \frac{k - \beta k}{n * (k - \beta k + \beta)} \quad (10.3)$$

b) Derive the PageRank  $r_t$  of the target page as a function of  $x_{\text{acc}}$ ,  $k$ ,  $\beta$ ,  $n$ . You do not have to simplify.

The PageRank of a target page consists of the accessible pages as well as the PageRank from all supporting pages:

$$r_t = \beta x_{\text{acc}} + \beta \sum_{i=1}^k \frac{r_i}{k} + \frac{1 - \beta}{n} \quad (10.4)$$

We can again simplify due to the fact that all supporting pages have the same PageRank and insert the solution from a) for  $r_s$ :

$$r_t = \beta x_{\text{acc}} + \frac{\beta^2}{k - \beta k + \beta} r_t + \beta \frac{k - \beta k}{n(k - \beta k + \beta)} + \frac{1 - \beta}{n} \quad (10.5)$$

Lastly, we can rewrite w.r.t.  $r_t$ :

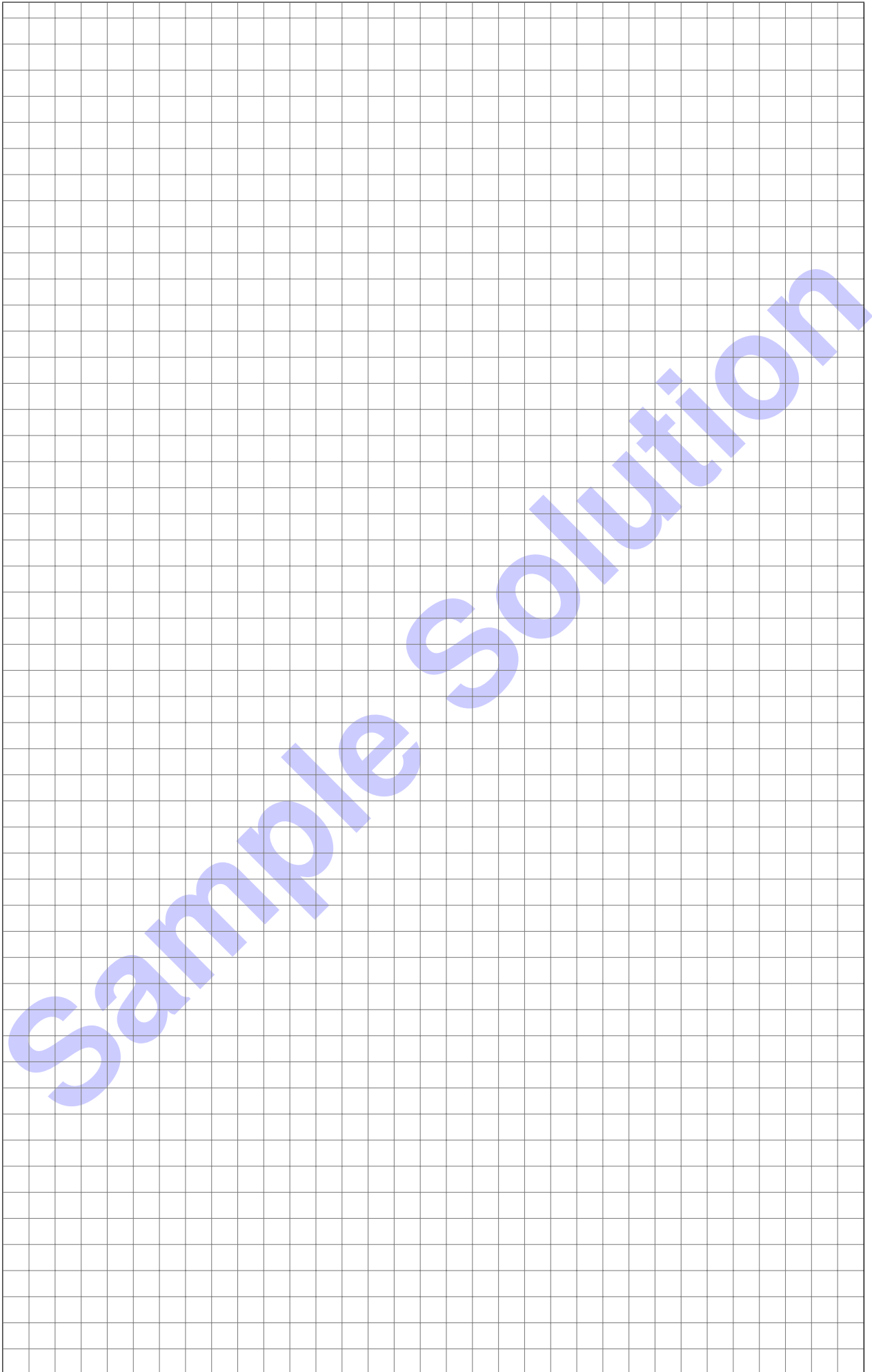
$$r_t = \left( \beta x_{\text{acc}} + \beta \frac{k - \beta k}{n(k - \beta k + \beta)} + \frac{1 - \beta}{n} \right) * \left( 1 - \frac{\beta^2}{k - \beta k + \beta} \right)^{-1} \quad (10.6)$$

c) How can the spammer modify the edges of the  $k$  supporting pages to increase the PageRank score  $r_t$  of the target page? Justify your answer.

Since currently, the is PageRank score is shared by the spam farm due to connections between the supporting pages, the spammer could increase the PageRank of the target page by decreasing the number of edges between supporting pages or completely removing edges between supporting pages.

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

A large grid of graph paper for solutions, with a diagonal watermark reading "Sample Solution". The grid is composed of small squares, and the watermark is written in a large, light blue font.



Sample Solution