

**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 (Problem sheet)

**Graded Exercise:** IN2323 / Retake

**Date:** Thursday 14<sup>th</sup> October, 2021

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

**Time:** 14:15 – 15:30

### Working instructions

- **DO NOT SUBMIT THIS SHEET! ONLY SUBMIT YOUR PERSONALIZED ANSWER SHEET THAT IS DISTRIBUTED THROUGH TUMEXAM!**
- Make sure that you solve the version of the problem stated on your personalized answer sheet (e.g., Problem 1 (Version B), Problem 2 (Version A), etc.)

## Problem 1: Normalizing Flows (Version A) (6 credits)

|   |  |
|---|--|
| 0 |  |
| 1 |  |
| 2 |  |
| 3 |  |
| 4 |  |
| 5 |  |
| 6 |  |

We consider two transformations  $f_1, f_2 : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  for use in normalizing flows.

Let

$$f_1(\mathbf{z}) = \begin{bmatrix} (1 + \max(0, z_2)) \cdot z_1 - \min(0, z_2) \cdot z_3 \\ (z_2)^3 \\ z_1 - z_3 \end{bmatrix} \quad \text{and} \quad f_2(\mathbf{z}) = \begin{bmatrix} (z_1)^3 \\ (z_3)^3 \cdot \exp(z_2) \\ z_1 \cdot |z_3| \end{bmatrix}.$$

Prove or disprove whether  $f_1$  and/or  $f_2$  are invertible.

Transformation  $f_1$  is invertible. Let  $\mathbf{x} = f(\mathbf{z})$  for some unknown  $\mathbf{z} \in \mathbb{R}^3$ . We can first solve for  $z_2$  to find  $z_2 = \sqrt[3]{x_2}$ . Then, we can make a case distinction. If  $z_2 \geq 0$  our problem can be restated as

$$x_1 = (1 + \sqrt[3]{x_2}) \cdot z_1 \quad (1.1)$$

$$x_3 = z_1 - z_3, \quad (1.2)$$

which can be uniquely solved because both equations are linearly independent (in the first equation,  $z_1$  always has a non-zero coefficient). If  $z_2 < 0$ , our problem can be restated as

$$x_1 = z_1 - \sqrt[3]{x_2} \cdot z_3 \quad (1.3)$$

$$x_3 = z_1 - z_3, \quad (1.4)$$

for which we can again find a unique solution (in the first equation,  $z_3$  has a positive coefficient, in the second equation,  $z_3$  has a negative coefficient).

Transformation  $f_2$  is not invertible. All  $\mathbf{z}$  with  $z_1 = z_3 = 0$  get mapped to  $\begin{bmatrix} 0 & 0 & 0 \end{bmatrix}^T$

## Problem 1: Normalizing Flows (Version B) (6 credits)

We consider two transformations  $f_1, f_2 : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  for use in normalizing flows.

Let

$$f_1(\mathbf{z}) = \begin{bmatrix} (z_2)^3 \\ (1 + \max(0, z_2)) \cdot z_1 - \min(0, z_2) \cdot z_3 \\ z_1 - z_3 \end{bmatrix} \quad \text{and} \quad f_2(\mathbf{z}) = \begin{bmatrix} (z_1)^3 \\ \ln(1 + |z_3|) \cdot \exp(z_2) \\ z_1 \cdot z_3 \end{bmatrix}.$$

Prove or disprove whether  $f_1$  and/or  $f_2$  are invertible.

|                          |   |
|--------------------------|---|
| <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 |

Transformation  $f_1$  is invertible. Let  $\mathbf{x} = f(\mathbf{z})$  for some unknown  $\mathbf{z} \in \mathbb{R}^3$ . We can first solve for  $z_2$  to find  $z_2 = \sqrt[3]{x_1}$ . Then, we can make a case distinction. If  $z_2 \geq 0$  our problem can be restated as

$$x_2 = (1 + \sqrt[3]{x_1}) \cdot z_1 \quad (2.1)$$

$$x_3 = z_1 - z_3, \quad (2.2)$$

which can be uniquely solved because both equations are linearly independent (in the first equation,  $z_1$  always has a non-zero coefficient). If  $z_2 < 0$ , our problem can be restated as

$$x_2 = z_1 - \sqrt[3]{x_1} \cdot z_3 \quad (2.3)$$

$$x_3 = z_1 - z_3, \quad (2.4)$$

for which we can again find a unique solution (in the first equation,  $z_3$  has a positive coefficient, in the second equation,  $z_3$  has a negative coefficient).

Transformation  $f_2$  is not invertible. All  $\mathbf{z}$  with  $z_1 = z_3 = 0$  get mapped to  $[0 \ 0 \ 0]^T$ .

## Problem 1: Normalizing Flows (Version C) (6 credits)

|   |  |
|---|--|
| 0 |  |
| 1 |  |
| 2 |  |
| 3 |  |
| 4 |  |
| 5 |  |
| 6 |  |

We consider two transformations  $f_1, f_2 : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  for use in normalizing flows.

Let

$$f_1(\mathbf{z}) = \begin{bmatrix} \min(0, z_2) \cdot z_3 + (1 + \max(0, z_2)) \cdot z_1 \\ (z_2)^3 \\ z_1 + 2 \cdot z_3 \end{bmatrix} \quad \text{and} \quad f_2(\mathbf{z}) = \begin{bmatrix} \ln(1 + |z_1|) \\ (z_3)^3 \cdot \exp(z_2) \\ z_1 \cdot z_3 \end{bmatrix}.$$

Prove or disprove whether  $f_1$  and/or  $f_2$  are invertible.

Transformation  $f_1$  is invertible. Let  $\mathbf{x} = f(\mathbf{z})$  for some unknown  $\mathbf{z} \in \mathbb{R}^3$ . We can first solve for  $z_2$  to find  $z_2 = \sqrt[3]{x_2}$ . Then, we can make a case distinction. If  $z_2 \geq 0$  our problem can be restated as

$$x_1 = (1 + \sqrt[3]{x_2}) \cdot z_1 \quad (3.1)$$

$$x_3 = z_1 + 2 \cdot z_3, \quad (3.2)$$

which can be uniquely solved because both equations are linearly independent (in the first equation,  $z_1$  always has a non-zero coefficient). If  $z_2 < 0$ , our problem can be restated as

$$x_1 = z_1 + \sqrt[3]{x_2} \cdot z_3 \quad (3.3)$$

$$x_3 = z_1 + 2 \cdot z_3, \quad (3.4)$$

for which we can again find a unique solution (in the first equation,  $z_3$  has a negative coefficient, in the second equation,  $z_3$  has a positive coefficient).

Transformation  $f_2$  is not invertible. All  $\mathbf{z}$  with  $z_1 = z_3 = 0$  get mapped to  $\begin{bmatrix} 0 & 0 & 0 \end{bmatrix}^T$

## Problem 1: Normalizing Flows (Version D) (6 credits)

We consider two transformations  $f_1, f_2 : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  for use in normalizing flows.

Let

$$f_1(\mathbf{z}) = \begin{bmatrix} (z_2)^3 \\ \min(0, z_2) \cdot z_3 + (1 + \max(0, z_2)) \cdot z_1 \\ z_1 + 2 \cdot z_3 \end{bmatrix} \quad \text{and} \quad f_2(\mathbf{z}) = \begin{bmatrix} (z_2)^3 \\ z_2 \cdot |z_1| \\ (z_1)^3 \cdot \exp(z_3) \end{bmatrix}.$$

Prove or disprove whether  $f_1$  and/or  $f_2$  are invertible.

|                          |   |
|--------------------------|---|
| <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 |

Transformation  $f_1$  is invertible. Let  $\mathbf{x} = f(\mathbf{z})$  for some unknown  $\mathbf{z} \in \mathbb{R}^3$ . We can first solve for  $z_2$  to find  $z_2 = \sqrt[3]{x_1}$ . Then, we can make a case distinction. If  $z_2 \geq 0$  our problem can be restated as

$$x_2 = (1 + \sqrt[3]{x_1}) \cdot z_1 \quad (4.1)$$

$$x_3 = z_1 + 2 \cdot z_3, \quad (4.2)$$

which can be uniquely solved because both equations are linearly independent (in the first equation,  $z_1$  always has a non-zero coefficient). If  $z_2 < 0$ , our problem can be restated as

$$x_2 = z_1 + \sqrt[3]{x_1} \cdot z_3 \quad (4.3)$$

$$x_3 = z_1 + 2 \cdot z_3, \quad (4.4)$$

for which we can again find a unique solution (in the first equation,  $z_3$  has a negative coefficient, in the second equation,  $z_3$  has a positive coefficient).

Transformation  $f_2$  is not invertible. All  $\mathbf{z}$  with  $z_2 = z_1 = 0$  get mapped to  $\begin{bmatrix} 0 & 0 & 0 \end{bmatrix}^T$

## Problem 2: Variational Inference (Version A) (7 credits)

Suppose we are given a latent variable model

$$p(z) = \mathcal{N}(z; 0, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)$$

$$p_\theta(x|z) = \mathcal{N}(x; z + 5, \theta^2) = \frac{1}{\theta\sqrt{2\pi}} \exp\left(-\frac{(x - z - 5)^2}{2\theta^2}\right)$$

where  $x, z \in \mathbb{R}$ . We parametrize the variational distribution  $q_\phi(z)$  as:

$$q_\phi(z) = \mathcal{N}(z; \phi, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(z - \phi)^2}{2}\right)$$

a) Derive the evidence lower bound (ELBO) for this particular parametrization. Simplify the parts depending on  $\phi$  as far as possible.

*Reminder:* The ELBO for parameters  $\theta$  and variational distribution  $q_\phi$  is defined as

$$\mathcal{L}(\theta, q_\phi) = \mathbb{E}_{\mathbf{z} \sim q_\phi} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z})].$$

*Hint:* Given a random variable  $X$ , the variance decomposition  $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$  can be rewritten as

$$\mathbb{E}[X^2] = \text{Var}(X) + \mathbb{E}[X]^2.$$

$$\begin{aligned} \mathcal{L}(\theta, q_\phi) &= \mathbb{E}_{\mathbf{z} \sim q_\phi} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z})] \\ &= \mathbb{E}_{\mathbf{z} \sim q_\phi} [\log p_\theta(\mathbf{x} | \mathbf{z}) + \log p(\mathbf{z}) - \log q_\phi(\mathbf{z})] \\ &= \mathbb{E}_{\mathbf{z} \sim q_\phi} \left[ -\frac{(x - z - 5)^2}{2\theta^2} - \frac{z^2}{2} + \frac{(z - \phi)^2}{2} \right] + \underbrace{\log \frac{1}{\theta\sqrt{2\pi}} + \log \frac{1}{\sqrt{2\pi}} - \log \frac{1}{\sqrt{2\pi}}}_{\text{constant}} \\ &= \mathbb{E}_{\mathbf{z} \sim q_\phi} \left[ \underbrace{-\frac{(x - 5)^2}{2\theta^2}}_{\text{constant w.r.t. } z, \phi} + 2\frac{(x - 5)z}{2\theta^2} - \frac{z^2}{2\theta^2} - \frac{z^2}{2} + \frac{z^2}{2} - 2\frac{z\phi}{2} + \frac{\phi^2}{2} \right] + C \\ &= \frac{x - 5}{\theta^2} \mathbb{E}_{\mathbf{z} \sim q_\phi} [z] - \frac{1}{2\theta^2} \mathbb{E}_{\mathbf{z} \sim q_\phi} [z^2] - \phi \mathbb{E}_{\mathbf{z} \sim q_\phi} [z] + \frac{\phi^2}{2} + C \\ &= \frac{x - 5}{\theta^2} \phi - \frac{\phi^2 + 1}{2\theta^2} - \phi^2 + \frac{\phi^2}{2} + C \\ &= \frac{x - 5}{\theta^2} \phi - \left( \frac{1}{2\theta^2} + \frac{1}{2} \right) \phi^2 + C \end{aligned}$$

b) Suppose  $\theta$  is fixed. Derive the value of  $\phi$  that maximizes the ELBO.

$$\begin{aligned}\frac{d\mathcal{L}(\theta, q_\phi)}{d\phi} &= \frac{x-5}{\theta^2} - 2\left(\frac{1}{2\theta^2} + \frac{1}{2}\right)\phi \\ &= \frac{x-5}{\theta^2} - \frac{1+\theta^2}{\theta^2}\phi\end{aligned}$$

Set  $\frac{d\mathcal{L}(\theta, q_\phi)}{d\phi} = 0$  and solve for  $\phi$ :

$$\begin{aligned}\frac{x-5}{\theta^2} - \frac{1+\theta^2}{\theta^2}\phi &= 0 \\ \phi &= \frac{x-5}{\theta^2} \frac{\cancel{\theta^2}}{1+\theta^2} \\ &= \frac{x-5}{1+\theta^2}\end{aligned}$$

Sample Solution

## Problem 2: Variational Inference (Version B) (7 credits)

Suppose we are given a latent variable model

$$p(z) = \mathcal{N}(z; 0, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)$$

$$p_\theta(x|z) = \mathcal{N}(x; 2z + 4, \theta^2) = \frac{1}{\theta\sqrt{2\pi}} \exp\left(-\frac{(x - 2z - 4)^2}{2\theta^2}\right)$$

where  $x, z \in \mathbb{R}$ . We parametrize the variational distribution  $q_\mu(z)$  as:

$$q_\mu(z) = \mathcal{N}(z; \mu, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(z - \mu)^2}{2}\right)$$

a) Derive the evidence lower bound (ELBO) for this particular parametrization. Simplify the parts depending on  $\mu$  as far as possible.

*Reminder:* The ELBO for parameters  $\theta$  and variational distribution  $q_\mu$  is defined as

$$\mathcal{L}(\theta, q_\mu) = \mathbb{E}_{\mathbf{z} \sim q_\mu} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\mu(\mathbf{z})].$$

*Hint:* Given a random variable  $X$ , the variance decomposition  $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$  can be rewritten as

$$\mathbb{E}[X^2] = \text{Var}(X) + \mathbb{E}[X]^2.$$

$$\begin{aligned} \mathcal{L}(\theta, q_\mu) &= \mathbb{E}_{\mathbf{z} \sim q_\mu} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})] \\ &= \mathbb{E}_{\mathbf{z} \sim q_\mu} [\log p_\theta(\mathbf{x} | \mathbf{z}) + \log p(\mathbf{z}) - \log q(\mathbf{z})] \\ &= \mathbb{E}_{\mathbf{z} \sim q_\mu} \left[ -\frac{(x - 2z - 4)^2}{2\theta^2} - \frac{z^2}{2} + \frac{(z - \mu)^2}{2} \right] + \underbrace{\log \frac{1}{\theta\sqrt{2\pi}} + \log \frac{1}{\sqrt{2\pi}} - \log \frac{1}{\sqrt{2\pi}}}_{\text{constant}} \\ &= \mathbb{E}_{\mathbf{z} \sim q_\mu} \left[ \underbrace{-\frac{(x - 4)^2}{2\theta^2}}_{\text{constant w.r.t. } z, \mu} + 2\frac{(x - 4)z}{2\theta^2} - \frac{4z^2}{2\theta^2} - \frac{z^2}{2} + \frac{z^2}{2} - 2\frac{z\mu}{2} + \frac{\mu^2}{2} \right] + C \\ &= \frac{x - 4}{\theta^2} \mathbb{E}_{\mathbf{z} \sim q_\mu} [z] - \frac{2}{\theta^2} \mathbb{E}_{\mathbf{z} \sim q_\mu} [z^2] - \mu \mathbb{E}_{\mathbf{z} \sim q_\mu} [z] + \frac{\mu^2}{2} + C \\ &= \frac{x - 4}{\theta^2} \mu - \frac{2\mu^2 + 2}{\theta^2} - \mu^2 + \frac{\mu^2}{2} + C \\ &= \frac{x - 4}{\theta^2} \mu - \left( \frac{2}{\theta^2} + \frac{1}{2} \right) \mu^2 + C \end{aligned}$$

b) Suppose  $\theta$  is fixed. Derive the value of  $\mu$  that maximizes the ELBO.



$$\begin{aligned}\frac{d\mathcal{L}(\theta, q_\mu)}{d\mu} &= \frac{x-4}{\theta^2} - 2\left(\frac{2}{\theta^2} + \frac{1}{2}\right)\mu \\ &= \frac{x-4}{\theta^2} - \frac{4+\theta^2}{\theta^2}\mu\end{aligned}$$

Set  $\frac{d\mathcal{L}(\theta, q_\mu)}{d\mu} = 0$  and solve for  $\mu$ :

$$\begin{aligned}\frac{x-4}{\theta^2} - \frac{4+\theta^2}{\theta^2}\mu &= 0 \\ \mu &= \frac{x-4}{\theta^2} \cancel{\theta^2} \frac{1}{4+\theta^2} \\ &= \frac{x-4}{4+\theta^2}\end{aligned}$$

## Problem 2: Variational Inference (Version C) (7 credits)

Suppose we are given a latent variable model

$$p(z) = \mathcal{N}(z; 0, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)$$

$$p_\theta(x|z) = \mathcal{N}(x; z + 3, \theta^2) = \frac{1}{\theta\sqrt{2\pi}} \exp\left(-\frac{(x - z - 3)^2}{2\theta^2}\right)$$

where  $x, z \in \mathbb{R}$ . We parametrize the variational distribution  $q_\phi(z)$  as:

$$q_\phi(z) = \mathcal{N}(z; \phi, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(z - \phi)^2}{2}\right)$$

a) Derive the evidence lower bound (ELBO) for this particular parametrization. Simplify the parts depending on  $\phi$  as far as possible.

*Reminder:* The ELBO for parameters  $\theta$  and variational distribution  $q_\phi$  is defined as

$$\mathcal{L}(\theta, q_\phi) = \mathbb{E}_{\mathbf{z} \sim q_\phi} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z})].$$

*Hint:* Given a random variable  $X$ , the variance decomposition  $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$  can be rewritten as

$$\mathbb{E}[X^2] = \text{Var}(X) + \mathbb{E}[X]^2.$$

$$\begin{aligned} \mathcal{L}(\theta, q_\phi) &= \mathbb{E}_{\mathbf{z} \sim q_\phi} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})] \\ &= \mathbb{E}_{\mathbf{z} \sim q_\phi} [\log p_\theta(\mathbf{x} | \mathbf{z}) + \log p(\mathbf{z}) - \log q(\mathbf{z})] \\ &= \mathbb{E}_{\mathbf{z} \sim q_\phi} \left[ -\frac{(x - z - 3)^2}{2\theta^2} - \frac{z^2}{2} + \frac{(z - \phi)^2}{2} \right] + \underbrace{\log \frac{1}{\theta\sqrt{2\pi}} + \log \frac{1}{\sqrt{2\pi}} - \log \frac{1}{\sqrt{2\pi}}}_{\text{constant}} \\ &= \mathbb{E}_{\mathbf{z} \sim q_\phi} \left[ \underbrace{-\frac{(x - 3)^2}{2\theta^2}}_{\text{constant w.r.t. } z, \phi} + 2\frac{(x - 3)z}{2\theta^2} - \frac{z^2}{2\theta^2} - \frac{z^2}{2} + \frac{z^2}{2} - 2\frac{z\phi}{2} + \frac{\phi^2}{2} \right] + C \\ &= \frac{x - 3}{\theta^2} \mathbb{E}_{\mathbf{z} \sim q_\phi} [z] - \frac{1}{2\theta^2} \mathbb{E}_{\mathbf{z} \sim q_\phi} [z^2] - \phi \mathbb{E}_{\mathbf{z} \sim q_\phi} [z] + \frac{\phi^2}{2} + C \\ &= \frac{x - 3}{\theta^2} \phi - \frac{\phi^2 + 1}{2\theta^2} - \phi^2 + \frac{\phi^2}{2} + C \\ &= \frac{x - 3}{\theta^2} \phi - \left( \frac{1}{2\theta^2} + \frac{1}{2} \right) \phi^2 + C \end{aligned}$$

b) Suppose  $\theta$  is fixed. Derive the value of  $\phi$  that maximizes the ELBO.

$$\begin{aligned}\frac{d\mathcal{L}(\theta, q_\phi)}{d\phi} &= \frac{x-3}{\theta^2} - 2\left(\frac{1}{2\theta^2} + \frac{1}{2}\right)\phi \\ &= \frac{x-3}{\theta^2} - \frac{1+\theta^2}{\theta^2}\phi\end{aligned}$$

Set  $\frac{d\mathcal{L}(\theta, q_\phi)}{d\phi} = 0$  and solve for  $\phi$ :

$$\begin{aligned}\frac{x-3}{\theta^2} - \frac{1+\theta^2}{\theta^2}\phi &= 0 \\ \phi &= \frac{x-3}{\theta^2} \frac{\cancel{\theta^2}}{1+\theta^2} \\ &= \frac{x-3}{1+\theta^2}\end{aligned}$$

## Problem 2: Variational Inference (Version D) (7 credits)

Suppose we are given a latent variable model

$$p(z) = \mathcal{N}(z; 0, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)$$

$$p_\theta(x|z) = \mathcal{N}(x; z + 7, \theta^2) = \frac{1}{\theta\sqrt{2\pi}} \exp\left(-\frac{(x - z - 7)^2}{2\theta^2}\right)$$

where  $x, z \in \mathbb{R}$ . We parametrize the variational distribution  $q_\mu(z)$  as:

$$q_\mu(z) = \mathcal{N}(z; \mu, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(z - \mu)^2}{2}\right)$$

a) Derive the evidence lower bound (ELBO) for this particular parametrization. Simplify the parts depending on  $\mu$  as far as possible.

*Reminder:* The ELBO for parameters  $\theta$  and variational distribution  $q_\mu$  is defined as

$$\mathcal{L}(\theta, q_\mu) = \mathbb{E}_{\mathbf{z} \sim q_\mu} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\mu(\mathbf{z})].$$

*Hint:* Given a random variable  $X$ , the variance decomposition  $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$  can be rewritten as

$$\mathbb{E}[X^2] = \text{Var}(X) + \mathbb{E}[X]^2.$$

$$\begin{aligned} \mathcal{L}(\theta, q_\mu) &= \mathbb{E}_{\mathbf{z} \sim q_\mu} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})] \\ &= \mathbb{E}_{\mathbf{z} \sim q_\mu} [\log p_\theta(\mathbf{x} | \mathbf{z}) + \log p(\mathbf{z}) - \log q(\mathbf{z})] \\ &= \mathbb{E}_{\mathbf{z} \sim q_\mu} \left[ -\frac{(x - z - 7)^2}{2\theta^2} - \frac{z^2}{2} + \frac{(z - \mu)^2}{2} \right] + \underbrace{\log \frac{1}{\theta\sqrt{2\pi}} + \log \frac{1}{\sqrt{2\pi}} - \log \frac{1}{\sqrt{2\pi}}}_{\text{constant}} \\ &= \mathbb{E}_{\mathbf{z} \sim q_\mu} \left[ \underbrace{-\frac{(x - 7)^2}{2\theta^2}}_{\text{constant w.r.t. } z, \mu} + 2 \frac{(x - 7)z}{2\theta^2} - \frac{z^2}{2\theta^2} - \frac{z^2}{2} + \frac{z^2}{2} - 2 \frac{z\mu}{2} + \frac{\mu^2}{2} \right] + C \\ &= \frac{x - 7}{\theta^2} \mathbb{E}_{\mathbf{z} \sim q_\mu} [z] - \frac{1}{2\theta^2} \mathbb{E}_{\mathbf{z} \sim q_\mu} [z^2] - \mu \mathbb{E}_{\mathbf{z} \sim q_\mu} [z] + \frac{\mu^2}{2} + C \\ &= \frac{x - 7}{\theta^2} \mu - \frac{\mu^2 + 1}{2\theta^2} - \mu^2 + \frac{\mu^2}{2} + C \\ &= \frac{x - 7}{\theta^2} \mu - \left( \frac{1}{2\theta^2} + \frac{1}{2} \right) \mu^2 + C \end{aligned}$$

b) Suppose  $\theta$  is fixed. Derive the value of  $\mu$  that maximizes the ELBO.

$$\begin{aligned}\frac{d\mathcal{L}(\theta, q_\mu)}{d\mu} &= \frac{x-7}{\theta^2} - 2\left(\frac{1}{2\theta^2} + \frac{1}{2}\right)\mu \\ &= \frac{x-7}{\theta^2} - \frac{1+\theta^2}{\theta^2}\mu\end{aligned}$$

Set  $\frac{d\mathcal{L}(\theta, q_\mu)}{d\mu} = 0$  and solve for  $\mu$ :

$$\begin{aligned}\frac{x-7}{\theta^2} - \frac{1+\theta^2}{\theta^2}\mu &= 0 \\ \mu &= \frac{x-7}{\theta^2} \frac{\cancel{\theta^2}}{1+\theta^2} \\ &= \frac{x-7}{1+\theta^2}\end{aligned}$$

### Problem 3: Variational Autoencoder (Version A) (2 credits)



We would like to define a variational autoencoder model for black-and-white images. Each image is represented as a binary vector  $\mathbf{x} \in \{0, 1\}^N$ . We define the conditional distribution  $p_\theta(\mathbf{x}|\mathbf{z})$  as follows.

1. We obtain the distribution parameters as

$$\lambda = \exp(f_\theta(\mathbf{z})),$$

where  $\mathbf{z} \in \mathbb{R}^L$  is the latent variable and  $f_\theta : \mathbb{R}^L \rightarrow \mathbb{R}^N$  is the decoder neural network.

2. We obtain the conditional distribution as

$$p_\theta(\mathbf{x}|\mathbf{z}) = \prod_{i=1}^N \text{Exponential}(x_i|\lambda_i)$$

where  $\text{Exponential}(x|\lambda)$  is the exponential distribution with probability density function

$$\text{Exponential}(x_i|\lambda_i) = \begin{cases} \lambda_i e^{-\lambda_i x_i} & \text{if } x_i \geq 0, \\ 0 & \text{else.} \end{cases}$$

What is the main problem with the above definition of  $p_\theta(\mathbf{x}|\mathbf{z})$ ? Explain how we can modify the above definition to fix this problem. Justify your answer.

The exponential distribution is not suitable for modeling binary data. We have to

1. Replace the exponential distribution with a Bernoulli distribution.
2. Pass the decoder output through a sigmoid function to obtain the distribution parameters (instead of the exponential function).

### Problem 3: Variational Autoencoder (Version B) (2 credits)

We would like to define a variational autoencoder model for black-and-white images. Each image is represented as a binary vector  $\mathbf{x} \in \{0, 1\}^N$ . We define the conditional distribution  $p_\theta(\mathbf{x}|\mathbf{z})$  as follows.



1. We obtain the distribution parameters as

$$\lambda = \exp(f_\theta(\mathbf{z})),$$

where  $\mathbf{z} \in \mathbb{R}^L$  is the latent variable and  $f_\theta : \mathbb{R}^L \rightarrow \mathbb{R}^N$  is the decoder neural network.

2. We obtain the conditional distribution as

$$p_\theta(\mathbf{x}|\mathbf{z}) = \prod_{i=1}^N \text{Exponential}(x_i|\lambda_i)$$

where  $\text{Exponential}(x|\lambda)$  is the exponential distribution with probability density function

$$\text{Exponential}(x_i|\lambda_i) = \begin{cases} \lambda_i e^{-\lambda_i x_i} & \text{if } x_i \geq 0, \\ 0 & \text{else.} \end{cases}$$

What is the main problem with the above definition of  $p_\theta(\mathbf{x}|\mathbf{z})$ ? Explain how we can modify the above definition to fix this problem. Justify your answer.

The exponential distribution is not suitable for modeling binary data. We have to

1. Replace the exponential distribution with a Bernoulli distribution.
2. Pass the decoder output through a sigmoid function to obtain the distribution parameters (instead of the exponential function).

### Problem 3: Variational Autoencoder (Version C) (2 credits)



We would like to define a variational autoencoder model for black-and-white images. Each image is represented as a binary vector  $\mathbf{x} \in \{0, 1\}^N$ . We define the conditional distribution  $p_\theta(\mathbf{x}|\mathbf{z})$  as follows.

1. We obtain the distribution parameters as

$$\lambda = \exp(f_\theta(\mathbf{z})),$$

where  $\mathbf{z} \in \mathbb{R}^L$  is the latent variable and  $f_\theta : \mathbb{R}^L \rightarrow \mathbb{R}^N$  is the decoder neural network.

2. We obtain the conditional distribution as

$$p_\theta(\mathbf{x}|\mathbf{z}) = \prod_{i=1}^N \text{Exponential}(x_i|\lambda_i)$$

where  $\text{Exponential}(x|\lambda)$  is the exponential distribution with probability density function

$$\text{Exponential}(x_i|\lambda_i) = \begin{cases} \lambda_i e^{-\lambda_i x_i} & \text{if } x_i \geq 0, \\ 0 & \text{else.} \end{cases}$$

What is the main problem with the above definition of  $p_\theta(\mathbf{x}|\mathbf{z})$ ? Explain how we can modify the above definition to fix this problem. Justify your answer.

The exponential distribution is not suitable for modeling binary data. We have to

1. Replace the exponential distribution with a Bernoulli distribution.
2. Pass the decoder output through a sigmoid function to obtain the distribution parameters (instead of the exponential function).



### Problem 3: Variational Autoencoder (Version D) (2 credits)

We would like to define a variational autoencoder model for black-and-white images. Each image is represented as a binary vector  $\mathbf{x} \in \{0, 1\}^N$ . We define the conditional distribution  $p_\theta(\mathbf{x}|\mathbf{z})$  as follows.



1. We obtain the distribution parameters as

$$\lambda = \exp(f_\theta(\mathbf{z})),$$

where  $\mathbf{z} \in \mathbb{R}^L$  is the latent variable and  $f_\theta : \mathbb{R}^L \rightarrow \mathbb{R}^N$  is the decoder neural network.

2. We obtain the conditional distribution as

$$p_\theta(\mathbf{x}|\mathbf{z}) = \prod_{i=1}^N \text{Exponential}(x_i|\lambda_i)$$

where  $\text{Exponential}(x|\lambda)$  is the exponential distribution with probability density function

$$\text{Exponential}(x_i|\lambda_i) = \begin{cases} \lambda_i e^{-\lambda_i x_i} & \text{if } x_i \geq 0, \\ 0 & \text{else.} \end{cases}$$

What is the main problem with the above definition of  $p_\theta(\mathbf{x}|\mathbf{z})$ ? Explain how we can modify the above definition to fix this problem. Justify your answer.

The exponential distribution is not suitable for modeling binary data. We have to

1. Replace the exponential distribution with a Bernoulli distribution.
2. Pass the decoder output through a sigmoid function to obtain the distribution parameters (instead of the exponential function).

## Problem 4: Robustness - Convex Relaxation (Version A) (7 credits)

In the lecture, we have derived a tight convex relaxation for the ReLU activation function. Now we want to generalize this result to the LeakyReLU activation function

$$\text{LeakyReLU}(x) = \begin{cases} x & \text{for } x \geq 0 \\ \alpha x & \text{for } x < 0 \end{cases}$$

with  $\alpha \in (0, 1)$ .

Let  $x, y \in \mathbb{R}$  be the variables we use to model the function's input and output, respectively. Assume we know that  $l \leq x \leq u$  with  $l, u \in \mathbb{R}$ . Specify a set of **linear constraints** on  $[x \ y]^T$  that model the **convex hull** of  $[x \ \text{LeakyReLU}(x)]^T$ , i.e. whose feasible region is

$$\left\{ \lambda \begin{bmatrix} x_1 \\ \text{LeakyReLU}(x_1) \end{bmatrix} + (1 - \lambda) \begin{bmatrix} x_2 \\ \text{LeakyReLU}(x_2) \end{bmatrix} \mid x_1, x_2 \in [l, u] \wedge \lambda \in [0, 1] \right\}.$$

*Reminder:* A linear constraint is an inequality or equality relation between terms that are linear in  $x$  and  $y$ .

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

As for the normal ReLU nonlinearity, we can distinguish three cases:

Case 1:

$l, u < 0$ .

In this case, we can use the single constraint  $y = \alpha x$ .

Case 2:

$l, u \geq 0$ .

In this case, we can use the single constraint  $y = x$ .

Case 3:

$l < 0 \leq u$ .

In this case, we need three linear constraints:

$$y \geq \alpha x,$$

$$y \geq x,$$

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

## Problem 4: Robustness - Convex Relaxation (Version B) (7 credits)

In the lecture, we have derived a tight convex relaxation for the ReLU activation function. Now we want to generalize this result to the LeakyReLU activation function

$$\text{LeakyReLU}(x) = \begin{cases} x & \text{for } x \geq 0 \\ \alpha x & \text{for } x < 0 \end{cases}$$

with  $\alpha \in (0, 1)$ .

Let  $x, y \in \mathbb{R}$  be the variables we use to model the function's input and output, respectively. Assume we know that  $l \leq x \leq u$  with  $l, u \in \mathbb{R}$ . Specify a set of **linear constraints** on  $[x \ y]^T$  that model the **convex hull** of  $[x \ \text{LeakyReLU}(x)]^T$ , i.e. whose feasible region is

$$\left\{ \lambda \begin{bmatrix} x_1 \\ \text{LeakyReLU}(x_1) \end{bmatrix} + (1 - \lambda) \begin{bmatrix} x_2 \\ \text{LeakyReLU}(x_2) \end{bmatrix} \mid x_1, x_2 \in [l, u] \wedge \lambda \in [0, 1] \right\}.$$

*Reminder:* A linear constraint is an inequality or equality relation between terms that are linear in  $x$  and  $y$ .

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

As for the normal ReLU nonlinearity, we can distinguish three cases:

Case 1:

$l, u < 0$ .

In this case, we can use the single constraint  $y = \alpha x$ .

Case 2:

$l, u \geq 0$ .

In this case, we can use the single constraint  $y = x$ .

Case 3:

$l < 0 \leq u$ .

In this case, we need three linear constraints:

$$y \geq \alpha x,$$

$$y \geq x,$$

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

## Problem 4: Robustness - Convex Relaxation (Version C) (7 credits)

In the lecture, we have derived a tight convex relaxation for the ReLU activation function. Now we want to generalize this result to the LeakyReLU activation function

$$\text{LeakyReLU}(x) = \begin{cases} x & \text{for } x \geq 0 \\ \alpha x & \text{for } x < 0 \end{cases}$$

with  $\alpha \in (0, 1)$ .

Let  $x, y \in \mathbb{R}$  be the variables we use to model the function's input and output, respectively. Assume we know that  $l \leq x \leq u$  with  $l, u \in \mathbb{R}$ . Specify a set of **linear constraints** on  $[x \ y]^T$  that model the **convex hull** of  $[x \ \text{LeakyReLU}(x)]^T$ , i.e. whose feasible region is

$$\left\{ \lambda \begin{bmatrix} x_1 \\ \text{LeakyReLU}(x_1) \end{bmatrix} + (1 - \lambda) \begin{bmatrix} x_2 \\ \text{LeakyReLU}(x_2) \end{bmatrix} \mid x_1, x_2 \in [l, u] \wedge \lambda \in [0, 1] \right\}.$$

*Reminder:* A linear constraint is an inequality or equality relation between terms that are linear in  $x$  and  $y$ .

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

As for the normal ReLU nonlinearity, we can distinguish three cases:

Case 1:

$l, u < 0$ .

In this case, we can use the single constraint  $y = \alpha x$ .

Case 2:

$l, u \geq 0$ .

In this case, we can use the single constraint  $y = x$ .

Case 3:

$l < 0 \leq u$ .

In this case, we need three linear constraints:

$$y \geq \alpha x,$$

$$y \geq x,$$

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

## Problem 4: Robustness - Convex Relaxation (Version D) (7 credits)

In the lecture, we have derived a tight convex relaxation for the ReLU activation function. Now we want to generalize this result to the LeakyReLU activation function

$$\text{LeakyReLU}(x) = \begin{cases} x & \text{for } x \geq 0 \\ \alpha x & \text{for } x < 0 \end{cases}$$

with  $\alpha \in (0, 1)$ .

Let  $x, y \in \mathbb{R}$  be the variables we use to model the function's input and output, respectively. Assume we know that  $l \leq x \leq u$  with  $l, u \in \mathbb{R}$ . Specify a set of **linear constraints** on  $[x \ y]^T$  that model the **convex hull** of  $[x \ \text{LeakyReLU}(x)]^T$ , i.e. whose feasible region is

$$\left\{ \lambda \begin{bmatrix} x_1 \\ \text{LeakyReLU}(x_1) \end{bmatrix} + (1 - \lambda) \begin{bmatrix} x_2 \\ \text{LeakyReLU}(x_2) \end{bmatrix} \mid x_1, x_2 \in [l, u] \wedge \lambda \in [0, 1] \right\}.$$

*Reminder:* A linear constraint is an inequality or equality relation between terms that are linear in  $x$  and  $y$ .

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

As for the normal ReLU nonlinearity, we can distinguish three cases:

Case 1:

$l, u < 0$ .

In this case, we can use the single constraint  $y = \alpha x$ .

Case 2:

$l, u \geq 0$ .

In this case, we can use the single constraint  $y = x$ .

Case 3:

$l < 0 \leq u$ .

In this case, we need three linear constraints:

$$y \geq \alpha x,$$

$$y \geq x,$$

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

## Problem 5: Markov Chain Language Model (Version A) (7 credits)

We want to use a Markov chain to model a very simple language consisting of the 4 words I, orange, like, eat. While the words are borrowed from the English language, our simple language is not bound to its grammatical rules. The words map to the Markov chain parameters as follows.

$$\pi = \begin{matrix} & \begin{matrix} I \\ \text{orange} \\ \text{like} \\ \text{eat} \end{matrix} \end{matrix} \begin{pmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \\ \pi_4 \end{pmatrix} \quad \mathbf{A} = \begin{matrix} & \begin{matrix} I & \text{orange} & \text{like} & \text{eat} \end{matrix} \\ \begin{matrix} I \\ \text{orange} \\ \text{like} \\ \text{eat} \end{matrix} & \begin{pmatrix} A_{11} & & & \\ \vdots & \ddots & & \\ & & \ddots & \\ A_{41} & & & A_{44} \end{pmatrix} \end{matrix}$$

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

a) Fit the Markov chain to the following dataset of example sentences by computing the most likely parameters.

- I like orange
- I eat orange
- orange eat orange
- I like I

$$\pi = \begin{matrix} & \begin{matrix} I \\ \text{orange} \\ \text{like} \\ \text{eat} \end{matrix} \end{matrix} \begin{pmatrix} 3/4 \\ 1/4 \\ 0 \\ 0 \end{pmatrix} \quad \mathbf{A} = \begin{matrix} & \begin{matrix} I & \text{orange} & \text{like} & \text{eat} \end{matrix} \\ \begin{matrix} I \\ \text{orange} \\ \text{like} \\ \text{eat} \end{matrix} & \begin{pmatrix} 0 & 0 & 2/3 & 1/3 \\ 0 & 0 & 0 & 1 \\ 1/2 & 1/2 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \end{matrix}$$

For the remaining problems assume that you are given the following Markov chain parameters that were fit to a larger dataset.

$$\pi = \begin{matrix} & \begin{matrix} I \\ \text{orange} \\ \text{like} \\ \text{eat} \end{matrix} \end{matrix} \begin{pmatrix} 4/6 \\ 2/6 \\ 0 \\ 0 \end{pmatrix} \quad \mathbf{A} = \begin{matrix} & \begin{matrix} I & \text{orange} & \text{like} & \text{eat} \end{matrix} \\ \begin{matrix} I \\ \text{orange} \\ \text{like} \\ \text{eat} \end{matrix} & \begin{pmatrix} 0 & 1/6 & 3/6 & 2/6 \\ 0 & 0 & 2/6 & 4/6 \\ 1/6 & 3/6 & 1/6 & 1/6 \\ 1/6 & 5/6 & 0 & 0 \end{pmatrix} \end{matrix}$$

b) Which of the following two sentences is more likely according to the model? Justify your answer.

- 1) I like orange
- 2) orange eat I

Sentence 1 has a likelihood of

$$\frac{4}{6} \cdot \frac{3}{6} \cdot \frac{3}{6} = \frac{36}{216}$$

while sentence 2 has a likelihood of

$$\frac{2}{6} \cdot \frac{4}{6} \cdot \frac{1}{6} = \frac{8}{216}$$

Therefore sentence 1 is more likely according to the model.

c) Given that the 3rd word  $X_3$  of a sentence is orange, compute the (unnormalized) probability distribution over the previous word  $X_2$ . Justify your answer.

$$\begin{aligned} \Pr(X_2 \mid X_3 = \text{orange}) &= \frac{\Pr(X_3 = \text{orange} \mid X_2) \cdot \Pr(X_2)}{\Pr(X_3 = \text{orange})} \\ &\propto \Pr(X_3 = \text{orange} \mid X_2) \cdot \sum_{j=1}^3 \Pr(X_2 \mid X_1) \cdot \Pr(X_1) \\ &= \mathbf{A}_{:,2} \odot \mathbf{A}^T \boldsymbol{\pi} \\ &\propto \begin{pmatrix} 1 \\ 0 \\ 3 \\ 5 \end{pmatrix} \odot \mathbf{A}^T \begin{pmatrix} 4 \\ 2 \\ 0 \\ 0 \end{pmatrix} \\ &\propto \begin{pmatrix} 1 \\ 0 \\ 3 \\ 5 \end{pmatrix} \odot \begin{pmatrix} 0 \\ 4 \\ 16 \\ 16 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 48 \\ 80 \end{pmatrix} \end{aligned}$$



## Problem 5: Markov Chain Language Model (Version B) (7 credits)

We want to use a Markov chain to model a very simple language consisting of the 4 words I, orange, see, like. While the words are borrowed from the English language, our simple language is not bound to its grammatical rules. The words map to the Markov chain parameters as follows.

$$\pi = \begin{matrix} & \begin{matrix} \text{I} \\ \text{orange} \\ \text{see} \\ \text{like} \end{matrix} \end{matrix} \begin{pmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \\ \pi_4 \end{pmatrix} \quad \mathbf{A} = \begin{matrix} & \begin{matrix} \text{I} & \text{orange} & \text{see} & \text{like} \end{matrix} \\ \begin{matrix} \text{I} \\ \text{orange} \\ \text{see} \\ \text{like} \end{matrix} & \begin{pmatrix} A_{11} & & & \\ \vdots & \ddots & & \\ & & \ddots & \\ A_{41} & & & A_{44} \end{pmatrix} \end{matrix}$$

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

a) Fit the Markov chain to the following dataset of example sentences by computing the most likely parameters.

- I see I
- I like orange
- orange like orange
- I see orange

$$\pi = \begin{matrix} & \begin{matrix} \text{I} \\ \text{orange} \\ \text{see} \\ \text{like} \end{matrix} \end{matrix} \begin{pmatrix} 3/4 \\ 1/4 \\ 0 \\ 0 \end{pmatrix} \quad \mathbf{A} = \begin{matrix} & \begin{matrix} \text{I} & \text{orange} & \text{see} & \text{like} \end{matrix} \\ \begin{matrix} \text{I} \\ \text{orange} \\ \text{see} \\ \text{like} \end{matrix} & \begin{pmatrix} 0 & 0 & 2/3 & 1/3 \\ 0 & 0 & 0 & 1 \\ 1/2 & 1/2 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \end{matrix}$$

For the remaining problems assume that you are given the following Markov chain parameters that were fit to a larger dataset.

$$\pi = \begin{matrix} & \begin{matrix} \text{I} \\ \text{orange} \\ \text{see} \\ \text{like} \end{matrix} \end{matrix} \begin{pmatrix} 4/6 \\ 2/6 \\ 0 \\ 0 \end{pmatrix} \quad \mathbf{A} = \begin{matrix} & \begin{matrix} \text{I} & \text{orange} & \text{see} & \text{like} \end{matrix} \\ \begin{matrix} \text{I} \\ \text{orange} \\ \text{see} \\ \text{like} \end{matrix} & \begin{pmatrix} 0 & 1/6 & 3/6 & 2/6 \\ 0 & 0 & 2/6 & 4/6 \\ 1/6 & 3/6 & 1/6 & 1/6 \\ 1/6 & 5/6 & 0 & 0 \end{pmatrix} \end{matrix}$$

b) Which of the following two sentences is more likely according to the model? Justify your answer.

- 1) I see orange
- 2) orange like I



Sentence 1 has a likelihood of

$$\frac{4}{6} \cdot \frac{3}{6} \cdot \frac{3}{6} = \frac{36}{216}$$

while sentence 2 has a likelihood of

$$\frac{2}{6} \cdot \frac{4}{6} \cdot \frac{1}{6} = \frac{8}{216}$$

Therefore sentence 1 is more likely according to the model.

c) Given that the 3rd word  $X_3$  of a sentence is orange, compute the (unnormalized) probability distribution over the previous word  $X_2$ . Justify your answer.

$$\begin{aligned} \Pr(X_2 \mid X_3 = \text{orange}) &= \frac{\Pr(X_3 = \text{orange} \mid X_2) \cdot \Pr(X_2)}{\Pr(X_3 = \text{orange})} \\ &\propto \Pr(X_3 = \text{orange} \mid X_2) \cdot \sum_{j=1}^3 \Pr(X_2 \mid X_1) \cdot \Pr(X_1) \\ &= \mathbf{A}_{:,2} \odot \mathbf{A}^T \boldsymbol{\pi} \\ &\propto \begin{pmatrix} 1 \\ 0 \\ 3 \\ 5 \end{pmatrix} \odot \mathbf{A}^T \begin{pmatrix} 4 \\ 2 \\ 0 \\ 0 \end{pmatrix} \\ &\propto \begin{pmatrix} 1 \\ 0 \\ 3 \\ 5 \end{pmatrix} \odot \begin{pmatrix} 0 \\ 4 \\ 16 \\ 16 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 48 \\ 80 \end{pmatrix} \end{aligned}$$



## Problem 5: Markov Chain Language Model (Version C) (7 credits)

We want to use a Markov chain to model a very simple language consisting of the 4 words you, apple, see, like. While the words are borrowed from the English language, our simple language is not bound to its grammatical rules. The words map to the Markov chain parameters as follows.

$$\pi = \begin{matrix} & \text{you} & \text{apple} & \text{see} & \text{like} \\ \begin{matrix} \text{you} \\ \text{apple} \\ \text{see} \\ \text{like} \end{matrix} & \begin{pmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \\ \pi_4 \end{pmatrix} \end{matrix} \quad \mathbf{A} = \begin{matrix} & \text{you} & \text{apple} & \text{see} & \text{like} \\ \begin{matrix} \text{you} \\ \text{apple} \\ \text{see} \\ \text{like} \end{matrix} & \begin{pmatrix} A_{11} & & & \\ \vdots & \ddots & & \\ & & \ddots & \\ A_{41} & & & A_{44} \end{pmatrix} \end{matrix}$$

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

a) Fit the Markov chain to the following dataset of example sentences by computing the most likely parameters.

- you see apple
- you like apple
- apple like apple
- you see you

$$\pi = \begin{matrix} & \text{you} & \text{apple} & \text{see} & \text{like} \\ \begin{matrix} \text{you} \\ \text{apple} \\ \text{see} \\ \text{like} \end{matrix} & \begin{pmatrix} 3/4 \\ 1/4 \\ 0 \\ 0 \end{pmatrix} \end{matrix} \quad \mathbf{A} = \begin{matrix} & \text{you} & \text{apple} & \text{see} & \text{like} \\ \begin{matrix} \text{you} \\ \text{apple} \\ \text{see} \\ \text{like} \end{matrix} & \begin{pmatrix} 0 & 0 & 2/3 & 1/3 \\ 0 & 0 & 0 & 1 \\ 1/2 & 1/2 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \end{matrix}$$

For the remaining problems assume that you are given the following Markov chain parameters that were fit to a larger dataset.

$$\pi = \begin{matrix} & \text{you} & \text{apple} & \text{see} & \text{like} \\ \begin{matrix} \text{you} \\ \text{apple} \\ \text{see} \\ \text{like} \end{matrix} & \begin{pmatrix} 4/6 \\ 2/6 \\ 0 \\ 0 \end{pmatrix} \end{matrix} \quad \mathbf{A} = \begin{matrix} & \text{you} & \text{apple} & \text{see} & \text{like} \\ \begin{matrix} \text{you} \\ \text{apple} \\ \text{see} \\ \text{like} \end{matrix} & \begin{pmatrix} 0 & 1/6 & 3/6 & 2/6 \\ 0 & 0 & 2/6 & 4/6 \\ 1/6 & 3/6 & 1/6 & 1/6 \\ 1/6 & 5/6 & 0 & 0 \end{pmatrix} \end{matrix}$$

b) Which of the following two sentences is more likely according to the model? Justify your answer.

- 1) you see apple
- 2) apple like you

Sentence 1 has a likelihood of

$$\frac{4}{6} \cdot \frac{3}{6} \cdot \frac{3}{6} = \frac{36}{216}$$

while sentence 2 has a likelihood of

$$\frac{2}{6} \cdot \frac{4}{6} \cdot \frac{1}{6} = \frac{8}{216}$$

Therefore sentence 1 is more likely according to the model.

c) Given that the 3rd word  $X_3$  of a sentence is apple, compute the (unnormalized) probability distribution over the previous word  $X_2$ . Justify your answer.

$$\begin{aligned} \Pr(X_2 \mid X_3 = \text{apple}) &= \frac{\Pr(X_3 = \text{apple} \mid X_2) \cdot \Pr(X_2)}{\Pr(X_3 = \text{apple})} \\ &\propto \Pr(X_3 = \text{apple} \mid X_2) \cdot \sum_{j=1}^3 \Pr(X_2 \mid X_1) \cdot \Pr(X_1) \\ &= \mathbf{A}_{:,2} \odot \mathbf{A}^T \boldsymbol{\pi} \\ &\propto \begin{pmatrix} 1 \\ 0 \\ 3 \\ 5 \end{pmatrix} \odot \mathbf{A}^T \begin{pmatrix} 4 \\ 2 \\ 0 \\ 0 \end{pmatrix} \\ &\propto \begin{pmatrix} 1 \\ 0 \\ 3 \\ 5 \end{pmatrix} \odot \begin{pmatrix} 0 \\ 4 \\ 16 \\ 16 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 48 \\ 80 \end{pmatrix} \end{aligned}$$



## Problem 5: Markov Chain Language Model (Version D) (7 credits)

We want to use a Markov chain to model a very simple language consisting of the 4 words they, apple, like, eat. While the words are borrowed from the English language, our simple language is not bound to its grammatical rules. The words map to the Markov chain parameters as follows.

$$\pi = \begin{matrix} & \begin{matrix} \text{they} \\ \text{apple} \\ \text{like} \\ \text{eat} \end{matrix} \end{matrix} \begin{pmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \\ \pi_4 \end{pmatrix} \quad \mathbf{A} = \begin{matrix} & \begin{matrix} \text{they} & \text{apple} & \text{like} & \text{eat} \end{matrix} \\ \begin{matrix} \text{they} \\ \text{apple} \\ \text{like} \\ \text{eat} \end{matrix} & \begin{pmatrix} A_{11} & & \dots & A_{14} \\ \vdots & \ddots & & \vdots \\ A_{41} & & \dots & A_{44} \end{pmatrix} \end{matrix}$$

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

a) Fit the Markov chain to the following dataset of example sentences by computing the most likely parameters.

- they like they
- they eat apple
- apple eat apple
- they like apple

$$\pi = \begin{matrix} & \begin{matrix} \text{they} \\ \text{apple} \\ \text{like} \\ \text{eat} \end{matrix} \end{matrix} \begin{pmatrix} 3/4 \\ 1/4 \\ 0 \\ 0 \end{pmatrix} \quad \mathbf{A} = \begin{matrix} & \begin{matrix} \text{they} & \text{apple} & \text{like} & \text{eat} \end{matrix} \\ \begin{matrix} \text{they} \\ \text{apple} \\ \text{like} \\ \text{eat} \end{matrix} & \begin{pmatrix} 0 & 0 & 2/3 & 1/3 \\ 0 & 0 & 0 & 1 \\ 1/2 & 1/2 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \end{matrix}$$

For the remaining problems assume that you are given the following Markov chain parameters that were fit to a larger dataset.

$$\pi = \begin{matrix} & \begin{matrix} \text{they} \\ \text{apple} \\ \text{like} \\ \text{eat} \end{matrix} \end{matrix} \begin{pmatrix} 4/6 \\ 2/6 \\ 0 \\ 0 \end{pmatrix} \quad \mathbf{A} = \begin{matrix} & \begin{matrix} \text{they} & \text{apple} & \text{like} & \text{eat} \end{matrix} \\ \begin{matrix} \text{they} \\ \text{apple} \\ \text{like} \\ \text{eat} \end{matrix} & \begin{pmatrix} 0 & 1/6 & 3/6 & 2/6 \\ 0 & 0 & 2/6 & 4/6 \\ 1/6 & 3/6 & 1/6 & 1/6 \\ 1/6 & 5/6 & 0 & 0 \end{pmatrix} \end{matrix}$$

b) Which of the following two sentences is more likely according to the model? Justify your answer.

- 1) they like apple
- 2) apple eat they

Sentence 1 has a likelihood of

$$\frac{4}{6} \cdot \frac{3}{6} \cdot \frac{3}{6} = \frac{36}{216}$$

while sentence 2 has a likelihood of

$$\frac{2}{6} \cdot \frac{4}{6} \cdot \frac{1}{6} = \frac{8}{216}$$

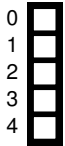
Therefore sentence 1 is more likely according to the model.

c) Given that the 3rd word  $X_3$  of a sentence is apple, compute the (unnormalized) probability distribution over the previous word  $X_2$ . Justify your answer.

$$\begin{aligned} \Pr(X_2 \mid X_3 = \text{apple}) &= \frac{\Pr(X_3 = \text{apple} \mid X_2) \cdot \Pr(X_2)}{\Pr(X_3 = \text{apple})} \\ &\propto \Pr(X_3 = \text{apple} \mid X_2) \cdot \sum_{j=1}^3 \Pr(X_2 \mid X_1) \cdot \Pr(X_1) \\ &= \mathbf{A}_{:,2} \odot \mathbf{A}^T \boldsymbol{\pi} \\ &\propto \begin{pmatrix} 1 \\ 0 \\ 3 \\ 5 \end{pmatrix} \odot \mathbf{A}^T \begin{pmatrix} 4 \\ 2 \\ 0 \\ 0 \end{pmatrix} \\ &\propto \begin{pmatrix} 1 \\ 0 \\ 3 \\ 5 \end{pmatrix} \odot \begin{pmatrix} 0 \\ 4 \\ 16 \\ 16 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 48 \\ 80 \end{pmatrix} \end{aligned}$$



## Problem 6: Neural Sequence Models (Version A) (4 credits)



We want to find out the limitations of our neural models for sequential data. To do that, we construct a dataset where the inputs are multiple sequences of  $n > 10$  numbers  $[x_1, x_2, \dots, x_n]$ ,  $x_i \in \mathbb{R}$ , where the corresponding target for each sequence is  $y = x_1 + x_n$ . We use four different encoders:

1. RNN with positional encoding
2. Transformer with positional encoding
3. Transformer without positional encoding
4. Dilated causal convolution with 2 hidden layers. We set dilation size to 2.

After processing the sequence with the above described encoders, we have access to hidden states  $\mathbf{h}_i \in \mathbb{R}^h$ , corresponding to the  $i$ -th place in the sequence. We use the last hidden state  $\mathbf{h}_n$  to make the prediction. For each of the four encoders, write down if they can learn the given task *in theory*. Justify your answer. For those encoders that can learn it, what issues might you encounter in practice?

Models that cannot learn:

3. because the model is equivariant
4. because the receptive field is not wide enough

Practical issues:

1. Long term-dependency leads to gradient issues
2. None, or taking the last state might be an issue instead of aggregating/having a dedicated state, or positional encoding might not capture the desired behavior

## Problem 6: Neural Sequence Models (Version B) (4 credits)

We want to find out the limitations of our neural models for sequential data. To do that, we construct a dataset where the inputs are multiple sequences of  $n > 10$  numbers  $[x_1, x_2, \dots, x_n]$ ,  $x_i \in \mathbb{R}$ , where the corresponding target for each sequence is  $y = x_1 + x_n$ . We use four different encoders:

1. Recurrent neural network
2. Transformer without positional encoding
3. Transformer with positional encoding
4. Sliding window neural network that takes  $[x_{i-k}, \dots, x_{i-1}, x_i]$  and outputs  $\mathbf{h}_i \in \mathbb{R}^h$ , for each  $i$ . We set  $k = 5$

After processing the sequence with the above described encoders, we have access to hidden states  $\mathbf{h}_i \in \mathbb{R}^h$ , corresponding to the  $i$ -th place in the sequence. We use the last hidden state  $\mathbf{h}_n$  to make the prediction. For each of the four encoders, write down if they can learn the given task *in theory*. Justify your answer. For those encoders that can learn it, what issues might you encounter in practice?

Models that cannot learn:

- 2. because the model is equivariant
- 4. because the receptive field is not wide enough

Practical issues:

- 1. Long term-dependency leads to gradient issues
- 3. None, or taking the last state might be an issue instead of aggregating/having a dedicated state, or positional encoding might not capture the desired behavior

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

## Problem 6: Neural Sequence Models (Version C) (4 credits)

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

We want to find out the limitations of our neural models for sequential data. To do that, we construct a dataset where the inputs are multiple sequences of  $n > 10$  numbers  $[x_1, x_2, \dots, x_n]$ ,  $x_i \in \mathbb{R}$ , where the corresponding target for each sequence is  $y = x_1 + x_n$ . We use four different encoders:

1. Transformer with positional encoding
2. Transformer without positional encoding
3. Multilayer neural network that takes vector in  $\mathbb{R}^n$  as input (all numbers concatenated) and outputs  $\mathbb{R}^{n \times h}$
4. Recurrent neural network

After processing the sequence with the above described encoders, we have access to hidden states  $\mathbf{h}_i \in \mathbb{R}^h$ , corresponding to the  $i$ -th place in the sequence. We use the last hidden state  $\mathbf{h}_n$  to make the prediction. For each of the four encoders, write down if they can learn the given task *in theory*. Justify your answer. For those encoders that can learn it, what issues might you encounter in practice?

Models that cannot learn:

2. because the model is equivariant

Practical issues:

4. Long term-dependency leads to gradient issues

1., 3. None, or taking the last state might be an issue instead of aggregating/having a dedicated state, or positional encoding might not capture the desired behavior etc.



## Problem 6: Neural Sequence Models (Version D) (4 credits)

We want to find out the limitations of our neural models for sequential data. To do that, we construct a dataset where the inputs are multiple sequences of  $n > 10$  numbers  $[x_1, x_2, \dots, x_n]$ ,  $x_i \in \mathbb{R}$ , where the corresponding target for each sequence is  $y = x_1 + x_n$ . We use four different encoders:

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

1. Recurrent neural network
2. Dilated causal convolution with 2 hidden layers. We set dilation size to 2.
3. Transformer with positional encoding
4. Transformer without positional encoding

After processing the sequence with the above described encoders, we have access to hidden states  $\mathbf{h}_i \in \mathbb{R}^h$ , corresponding to the  $i$ -th place in the sequence. We use the last hidden state  $\mathbf{h}_n$  to make the prediction. For each of the four encoders, write down if they can learn the given task *in theory*. Justify your answer. For those encoders that can learn it, what issues might you encounter in practice?

Models that cannot learn:

2. because the receptive field is not wide enough
4. because the model is equivariant

Practical issues:

1. Long term-dependency leads to gradient issues
3. None

## Problem 7: Temporal Point Process (Version A) (6 credits)

We fit a homogeneous Poisson process with intensity parameter  $\mu$  to model event occurrences in a time interval  $[0, T]$ . We have observed a single sequence that contains  $n$  points  $\{t_1, t_2, \dots, t_n\}$ ,  $t_i \in [0, T]$ .

a) Derive the maximum likelihood estimate of the parameter  $\mu$ .

$$\begin{aligned}\log p(\{t_1, t_2, \dots, t_n\}) &= \sum_i^n \log \mu - \int_0^T \mu dt \\ &= n \log \mu - \mu T\end{aligned}$$

Take the derivative:

$$\frac{d \log p}{d\mu} = \frac{n}{\mu} - T$$

Set to zero and solve for  $\mu$ :

$$\frac{n}{\mu} - T = 0 \Rightarrow \mu = \frac{n}{T}$$

b) Suppose we install a sensor next to a busy road that records the times when cars drive by. We model the times as described above, using the events from the whole day as one sequence. We estimate  $\mu$  using data we collected in one year. Our task is to find the least busy 2 hour interval in each day to close down the road for maintenance. Can we use the homogeneous Poisson process to achieve this? If not, can you suggest an alternative model? Justify your answer.

No. Alternative: inhomogeneous, but not Hawkes (unless good justification)

## Problem 7: Temporal Point Process (Version B) (6 credits)

We fit a homogeneous Poisson process with intensity parameter  $\mu$  to model event occurrences in a time interval  $[0, 5]$ . We have observed a single sequence  $\{0.7, 0.8, 1.5, 2.3, 4.7\}$ .

a) Derive the maximum likelihood estimate of the parameter  $\mu$ .

0  
1  
2  
3  
4

$$\begin{aligned}\log p(\{t_1, t_2, \dots, t_n\}) &= \sum_i^n \log \mu - \int_0^T \mu dt \\ &= n \log \mu - \mu T\end{aligned}$$

Take the derivative:

$$\frac{d \log p}{d \mu} = \frac{n}{\mu} - T$$

Set to zero and solve for  $\mu$ :

$$\frac{n}{\mu} - T = 0 \Rightarrow \mu = \frac{n}{T}$$

$$\mu = 1$$

b) Suppose we install a sensor next to a busy road that records the times when cars drive by. We model the times as described above, using the events from the whole day as one sequence. For each day of the week we estimate the parameter  $\mu$  using data we collected in one year. That means we have  $\mu_{\text{Mon}}, \mu_{\text{Tue}}, \dots, \mu_{\text{Sun}}$ , each  $\mu$  corresponding to one day of the week. Our task is to find the least busy day of the week to close down the road for maintenance. Can we use the homogeneous Poisson process to achieve this? If not, can you suggest an alternative model? Justify your answer.

0  
1  
2

Yes. Any satisfying explanation, e.g., we don't need anything more than homogeneous for such data.

## Problem 7: Temporal Point Process (Version C) (6 credits)

We fit a homogeneous Poisson process with intensity parameter  $\mu$  to model event occurrences in a time interval  $[0, 2]$ . We have observed a single sequence  $\{0.1, 0.8, 1.3, 1.5, 1.7, 1.9\}$ .

a) Derive the maximum likelihood estimate of the parameter  $\mu$ .

$$\begin{aligned}\log p(\{t_1, t_2, \dots, t_n\}) &= \sum_i^n \log \mu - \int_0^T \mu dt \\ &= n \log \mu - \mu T\end{aligned}$$

Take the derivative:

$$\frac{d \log p}{d \mu} = \frac{n}{\mu} - T$$

Set to zero and solve for  $\mu$ :

$$\frac{n}{\mu} - T = 0 \Rightarrow \mu = \frac{n}{T}$$

$$\mu = 3$$

b) Suppose we install a sensor next to a busy road that records the times when cars drive by. We model the times as described above, using the events from the whole day as one sequence. Using our model, we want to estimate the probability that less than 100 cars will pass our sensor in a day. Can we use the homogeneous Poisson process to achieve this? If not, can you suggest an alternative model? Justify your answer.

Yes. Any satisfying explanation, e.g., we don't need anything more than homogeneous for such data.

## Problem 7: Temporal Point Process (Version D) (6 credits)

We fit a homogeneous Poisson process with intensity parameter  $\mu$  to model event occurrences in a time interval  $[3, 13]$ . We have observed a single sequence  $\{3.5, 4.3, 4.5, 7.1, 8.3\}$ .

a) Derive the maximum likelihood estimate of the parameter  $\mu$ .

0  
1  
2  
3  
4

$$\begin{aligned}\log p(\{t_1, t_2, \dots, t_n\}) &= \sum_i^n \log \mu - \int_0^T \mu dt \\ &= n \log \mu - \mu T\end{aligned}$$

Take the derivative:

$$\frac{d \log p}{d\mu} = \frac{n}{\mu} - T$$

Set to zero and solve for  $\mu$ :

$$\frac{n}{\mu} - T = 0 \Rightarrow \mu = \frac{n}{T}$$

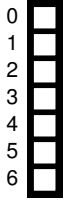
$$\mu = 1/2$$

b) Suppose we install a sensor next to a busy road that records the times when cars drive by. We model the times as described above, using the events from the whole day as one sequence. Using our model, we want to answer whether fast vehicles get stuck behind slower vehicles. That is, we want to see if observing one vehicle leads to a few more following behind it. Can we use the homogeneous Poisson process to achieve this? If not, can you suggest an alternative model? Justify your answer.

0  
1  
2

No. Alternative: Hawkes, but not inhomogeneous

## Problem 8: Clustering (Version A) (6 credits)



We consider the graph  $G = (E, V)$  with adjacency matrix  $\mathbf{A}$  where the nodes are separated into two clusters,  $C$  and  $\bar{C}$ . We define the associated random walk  $\Pr(X_{t+1} = j | X_t = i) = \frac{A_{ij}}{d_i}$  where  $d_i = \sum_j A_{ij}$  is the degree of node  $i$  and  $\Pr(X_0 = i) = \frac{d_i}{\text{vol}(V)}$  is the starting distribution where  $\text{vol}(V) = \sum_{i \in V} d_i$  is the volume of the set of nodes  $V$ . We define the probability to transition from cluster  $C$  to cluster  $\bar{C}$  in the first random walk step as  $\Pr(\bar{C} | C) = \Pr(X_1 \in \bar{C} | X_0 \in C)$  and vice versa. Show that the normalized cut satisfies the equation

$$\text{Ncut}(C, \bar{C}) = \Pr(\bar{C} | C) + \Pr(C | \bar{C}).$$

*Reminder:* The normalized cut of an undirected graph is defined as

$$\text{Ncut}(C, \bar{C}) = \frac{\text{cut}(C, \bar{C})}{\text{vol}(C)} + \frac{\text{cut}(C, \bar{C})}{\text{vol}(\bar{C})}.$$

We use the definition of the conditional probability:

$$\Pr(\bar{C} | C) = \frac{\Pr(X_1 \in \bar{C}, X_0 \in C)}{\Pr(X_0 \in C)}$$

First we compute the numerator:

$$\begin{aligned} \Pr(X_1 \in \bar{C}, X_0 \in C) &= \sum_{i \in \bar{C}, j \in C} \Pr(X_1 = j, X_0 = i) \\ &= \sum_{i \in \bar{C}, j \in C} \Pr(X_1 = j | X_0 = i) \Pr(X_0 = i) \\ &= \sum_{i \in \bar{C}, j \in C} \frac{d_i}{\text{vol}(V)} \frac{A_{ij}}{d_i} \\ &= \frac{1}{\text{vol}(V)} \sum_{i \in \bar{C}, j \in C} A_{ij} \end{aligned}$$

Second, we compute the denominator:

$$\Pr(X_0 \in C) = \sum_i \frac{d_i}{\text{vol}(V)} = \frac{\text{vol}(C)}{\text{vol}(V)}$$

We combine the two previous expressions and obtain:

$$\Pr(\bar{C} | C) + \Pr(C | \bar{C}) = \frac{\sum_{i \in \bar{C}, j \in C} A_{ij}}{\text{vol}(C)} + \frac{\sum_{i \in \bar{C}, j \in C} A_{ij}}{\text{vol}(\bar{C})} = \text{Ncut}$$

## Problem 8: Clustering (Version B) (6 credits)

We consider the graph  $G = (E, V)$  with adjacency matrix  $\mathbf{A}$  where the nodes are separated into two clusters,  $C$  and  $\bar{C}$ . We define the associated random walk  $\Pr(X_{t+1} = j | X_t = i) = \frac{A_{ij}}{d_i}$  where  $d_i = \sum_j A_{ij}$  is the degree of node  $i$  and  $\Pr(X_0 = i) = \frac{d_i}{\text{vol}(V)}$  is the starting distribution where  $\text{vol}(V) = \sum_{i \in V} d_i$  is the volume of the set of nodes  $V$ . We define the probability to transition from cluster  $C$  to cluster  $\bar{C}$  in the first random walk step as  $\Pr(\bar{C} | C) = \Pr(X_1 \in \bar{C} | X_0 \in C)$  and vice versa. Show that the normalized cut satisfies the equation

$$\text{Ncut}(C, \bar{C}) = \Pr(\bar{C} | C) + \Pr(C | \bar{C}).$$

*Reminder:* The normalized cut of an undirected graph is defined as

$$\text{Ncut}(C, \bar{C}) = \frac{\text{cut}(C, \bar{C})}{\text{vol}(C)} + \frac{\text{cut}(C, \bar{C})}{\text{vol}(\bar{C})}.$$

We use the definition of the conditional probability:

$$\Pr(\bar{C} | C) = \frac{\Pr(X_1 \in \bar{C}, X_0 \in C)}{\Pr(X_0 \in C)}$$

First we compute the numerator:

$$\begin{aligned} \Pr(X_1 \in \bar{C}, X_0 \in C) &= \sum_{i \in \bar{C}, j \in C} \Pr(X_1 = j, X_0 = i) \\ &= \sum_{i \in \bar{C}, j \in C} \Pr(X_1 = j | X_0 = i) \Pr(X_0 = i) \\ &= \sum_{i \in \bar{C}, j \in C} \frac{d_i}{\text{vol}(V)} \frac{A_{ij}}{d_i} \\ &= \frac{1}{\text{vol}(V)} \sum_{i \in \bar{C}, j \in C} A_{ij} \end{aligned}$$

Second, we compute the denominator:

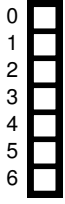
$$\Pr(X_0 \in C) = \sum_i \frac{d_i}{\text{vol}(V)} = \frac{\text{vol}(C)}{\text{vol}(V)}$$

We combine the two previous expressions and obtain:

$$\Pr(\bar{C} | C) + \Pr(C | \bar{C}) = \frac{\sum_{i \in \bar{C}, j \in C} A_{ij}}{\text{vol}(C)} + \frac{\sum_{i \in \bar{C}, j \in C} A_{ij}}{\text{vol}(\bar{C})} = \text{Ncut}$$

|                          |   |
|--------------------------|---|
| <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 |

## Problem 8: Clustering (Version C) (6 credits)



We consider the graph  $G = (E, V)$  with adjacency matrix  $\mathbf{A}$  where the nodes are separated into two clusters,  $C$  and  $\bar{C}$ . We define the associated random walk  $\Pr(X_{t+1} = j | X_t = i) = \frac{A_{ij}}{d_i}$  where  $d_i = \sum_j A_{ij}$  is the degree of node  $i$  and  $\Pr(X_0 = i) = \frac{d_i}{\text{vol}(V)}$  is the starting distribution where  $\text{vol}(V) = \sum_{i \in V} d_i$  is the volume of the set of nodes  $V$ . We define the probability to transition from cluster  $C$  to cluster  $\bar{C}$  in the first random walk step as  $\Pr(\bar{C} | C) = \Pr(X_1 \in \bar{C} | X_0 \in C)$  and vice versa. Show that the normalized cut satisfies the equation

$$\text{Ncut}(C, \bar{C}) = \Pr(\bar{C} | C) + \Pr(C | \bar{C}).$$

*Reminder:* The normalized cut of an undirected graph is defined as

$$\text{Ncut}(C, \bar{C}) = \frac{\text{cut}(C, \bar{C})}{\text{vol}(C)} + \frac{\text{cut}(C, \bar{C})}{\text{vol}(\bar{C})}.$$

We use the definition of the conditional probability:

$$\Pr(\bar{C} | C) = \frac{\Pr(X_1 \in \bar{C}, X_0 \in C)}{\Pr(X_0 \in C)}$$

First we compute the numerator:

$$\begin{aligned} \Pr(X_1 \in \bar{C}, X_0 \in C) &= \sum_{i \in \bar{C}, j \in C} \Pr(X_1 = j, X_0 = i) \\ &= \sum_{i \in \bar{C}, j \in C} \Pr(X_1 = j | X_0 = i) \Pr(X_0 = i) \\ &= \sum_{i \in \bar{C}, j \in C} \frac{d_i}{\text{vol}(V)} \frac{A_{ij}}{d_i} \\ &= \frac{1}{\text{vol}(V)} \sum_{i \in \bar{C}, j \in C} A_{ij} \end{aligned}$$

Second, we compute the denominator:

$$\Pr(X_0 \in C) = \sum_i \frac{d_i}{\text{vol}(V)} = \frac{\text{vol}(C)}{\text{vol}(V)}$$

We combine the two previous expressions and obtain:

$$\Pr(\bar{C} | C) + \Pr(C | \bar{C}) = \frac{\sum_{i \in \bar{C}, j \in C} A_{ij}}{\text{vol}(C)} + \frac{\sum_{i \in \bar{C}, j \in C} A_{ij}}{\text{vol}(\bar{C})} = \text{Ncut}$$



## Problem 8: Clustering (Version D) (6 credits)

We consider the graph  $G = (E, V)$  with adjacency matrix  $\mathbf{A}$  where the nodes are separated into two clusters,  $C$  and  $\bar{C}$ . We define the associated random walk  $\Pr(X_{t+1} = j | X_t = i) = \frac{A_{ij}}{d_i}$  where  $d_i = \sum_j A_{ij}$  is the degree of node  $i$  and  $\Pr(X_0 = i) = \frac{d_i}{\text{vol}(V)}$  is the starting distribution where  $\text{vol}(V) = \sum_{i \in V} d_i$  is the volume of the set of nodes  $V$ . We define the probability to transition from cluster  $C$  to cluster  $\bar{C}$  in the first random walk step as  $\Pr(\bar{C} | C) = \Pr(X_1 \in \bar{C} | X_0 \in C)$  and vice versa. Show that the normalized cut satisfies the equation

$$\text{Ncut}(C, \bar{C}) = \Pr(\bar{C} | C) + \Pr(C | \bar{C}).$$

*Reminder:* The normalized cut of an undirected graph is defined as

$$\text{Ncut}(C, \bar{C}) = \frac{\text{cut}(C, \bar{C})}{\text{vol}(C)} + \frac{\text{cut}(C, \bar{C})}{\text{vol}(\bar{C})}.$$

We use the definition of the conditional probability:

$$\Pr(\bar{C} | C) = \frac{\Pr(X_1 \in \bar{C}, X_0 \in C)}{\Pr(X_0 \in C)}$$

First we compute the numerator:

$$\begin{aligned} \Pr(X_1 \in \bar{C}, X_0 \in C) &= \sum_{i \in \bar{C}, j \in C} \Pr(X_1 = j, X_0 = i) \\ &= \sum_{i \in \bar{C}, j \in C} \Pr(X_1 = j | X_0 = i) \Pr(X_0 = i) \\ &= \sum_{i \in \bar{C}, j \in C} \frac{d_i}{\text{vol}(V)} \frac{A_{ij}}{d_i} \\ &= \frac{1}{\text{vol}(V)} \sum_{i \in \bar{C}, j \in C} A_{ij} \end{aligned}$$

Second, we compute the denominator:

$$\Pr(X_0 \in C) = \sum_i \frac{d_i}{\text{vol}(V)} = \frac{\text{vol}(C)}{\text{vol}(V)}$$

We combine the two previous expressions and obtain:

$$\Pr(\bar{C} | C) + \Pr(C | \bar{C}) = \frac{\sum_{i \in \bar{C}, j \in C} A_{ij}}{\text{vol}(C)} + \frac{\sum_{i \in \bar{C}, j \in C} A_{ij}}{\text{vol}(\bar{C})} = \text{Ncut}$$

|                          |   |
|--------------------------|---|
| <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 |

## Problem 9: Embeddings & Ranking (Version A) (6 credits)

We consider a graph  $G = (V, E)$  with adjacency matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  where  $A_{ij} = \mathbb{1}_{(i,j) \in E}$  indicates if an edge exist between node  $i$  and node  $j$  in the graph  $G$ . The node features are represented by the matrix  $\mathbf{X} \in \mathbb{R}^{n \times D}$ . We consider the three following models  $M_k, k \in \{1, 2, 3\}$  which produce node embeddings  $\mathbf{E}_k = M_k(G, \mathbf{X}) \in \mathbb{R}^{n \times D'}$ . The vector  $\mathbf{E}_k[i, :] \in \mathbb{R}^{D'}$  denotes the embedding of node  $i$  for model  $M_k$ :

- $M_1$ : Node2Vec.
- $M_2$ : Mean of Graph2Gauss i.e.  $\mathbf{E}_2[i, :] = \boldsymbol{\mu}_i$  where the Graph2Gauss mapping transforms node  $i$  into the Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i))$ .
- $M_3$ : Spectral embedding with  $k$  smallest eigenvectors.

- 0 ☐ a) We modify the attributed graph such that all nodes have the same features, i.e. the adjacency matrix is  $\mathbf{A}' = \mathbf{A}$  and  
 1 ☐ the new node attributes are  $\mathbf{X}'$  such that  $\mathbf{X}'[i, :] = \mathbf{X}'[j, :]$  for all  $(i, j)$ . For which model will the new node embeddings  
 2 ☐  $\mathbf{E}'_k = M_k(G', \mathbf{X}')$  be different from the embeddings obtained with the original attributed graphs  $\mathbf{E}_k = M_k(G, \mathbf{X})$ ? Justify  
 3 ☐ your answer.

- $\mathbf{E}_1 = \mathbf{E}'_1$ : Node2Vec does not account for node features.
- $\mathbf{E}'_2$  such that  $\mathbf{E}'_2[i, :] = \mathbf{E}'_2[j, :]$ : Graph2Gauss encodes a node based on its features only.
- $\mathbf{E}_3 = \mathbf{E}'_3$ : Spectral embedding does not account for node features.

- 0 ☐ b) We modify the attributed graph such that the graph is a clique, i.e. the new adjacency matrix is  $\mathbf{A}' = \mathbf{1} - \mathbf{I}$   
 1 ☐ where  $\mathbf{1}$  is the all-ones matrix, and the node attributes are  $\mathbf{X}' = \mathbf{X}$ . For which model will the new node embeddings  
 2 ☐  $\mathbf{E}'_k = M_k(G', \mathbf{X}')$  be different from the embeddings obtained with the original attributed graph  $\mathbf{E}_k = M_k(G, \mathbf{X})$ ? Justify  
 3 ☐ your answer.

- $E'_1$  such that  $E'_1[i, :] = E'_1[j, :]$ : All nodes have the same neighborhood.
- $E'_2 \neq E_2$ : The loss used for Graph2Gauss training depends on the graph edges.
- $E'_3$  such that  $E'_3[i, :] = E'_3[j, :]$ : All nodes have the same neighborhood.

Sample Solution

## Problem 9: Embeddings & Ranking (Version B) (6 credits)

We consider a graph  $G = (V, E)$  with adjacency matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  where  $A_{ij} = \mathbb{1}_{(i,j) \in E}$  indicates if an edge exist between node  $i$  and node  $j$  in the graph  $G$ . The node features are represented by the matrix  $\mathbf{X} \in \mathbb{R}^{n \times D}$ . We consider the three following models  $M_k, k \in \{1, 2, 3\}$  which produce node embeddings  $\mathbf{E}_k = M_k(G, \mathbf{X}) \in \mathbb{R}^{n \times D'}$ . The vector  $\mathbf{E}_k[i, :] \in \mathbb{R}^{D'}$  denotes the embedding of node  $i$  for model  $M_k$ :

- $M_1$ : Node2Vec.
- $M_2$ : Mean of Graph2Gauss i.e.  $\mathbf{E}_2[i, :] = \boldsymbol{\mu}_i$  where the Graph2Gauss mapping transforms node  $i$  into the Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i))$ .
- $M_3$ : Spectral embedding with  $k$  smallest eigenvectors.

- 0 ☐ a) We modify the attributed graph such that all nodes have the same features, i.e. the adjacency matrix is  $\mathbf{A}' = \mathbf{A}$  and  
1 ☐ the new node attributes are  $\mathbf{X}'$  such that  $\mathbf{X}'[i, :] = \mathbf{X}'[j, :]$  for all  $(i, j)$ . For which model will the new node embeddings  
2 ☐  $\mathbf{E}'_k = M_k(G', \mathbf{X}')$  be different from the embeddings obtained with the original attributed graphs  $\mathbf{E}_k = M_k(G, \mathbf{X})$ ? Justify  
3 ☐ your answer.

- $\mathbf{E}_1 = \mathbf{E}'_1$ : Node2Vec does not account for node features.
- $\mathbf{E}'_2$  such that  $\mathbf{E}'_2[i, :] = \mathbf{E}'_2[j, :]$ : Graph2Gauss encodes a node based on its features only.
- $\mathbf{E}_3 = \mathbf{E}'_3$ : Spectral embedding does not account for node features.

- 0 ☐ b) We modify the attributed graph such that the graph is a clique, i.e. the new adjacency matrix is  $\mathbf{A}' = \mathbf{1} - \mathbf{I}$   
1 ☐ where  $\mathbf{1}$  is the all-ones matrix, and the node attributes are  $\mathbf{X}' = \mathbf{X}$ . For which model will the new node embeddings  
2 ☐  $\mathbf{E}'_k = M_k(G', \mathbf{X}')$  be different from the embeddings obtained with the original attributed graph  $\mathbf{E}_k = M_k(G, \mathbf{X})$ ? Justify  
3 ☐ your answer.

- $E'_1$  such that  $E'_1[i, :] = E'_1[j, :]$ : All nodes have the same neighborhood.
- $E'_2 \neq E_2$ : The loss used for Graph2Gauss training depends on the graph edges.
- $E'_3$  such that  $E'_3[i, :] = E'_3[j, :]$ : All nodes have the same neighborhood.

Sample Solution

## Problem 9: Embeddings & Ranking (Version C) (6 credits)

We consider a graph  $G = (V, E)$  with adjacency matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  where  $A_{ij} = \mathbb{1}_{(i,j) \in E}$  indicates if an edge exist between node  $i$  and node  $j$  in the graph  $G$ . The node features are represented by the matrix  $\mathbf{X} \in \mathbb{R}^{n \times D}$ . We consider the three following models  $M_k, k \in \{1, 2, 3\}$  which produce node embeddings  $\mathbf{E}_k = M_k(G, \mathbf{X}) \in \mathbb{R}^{n \times D'}$ . The vector  $\mathbf{E}_k[i, :] \in \mathbb{R}^{D'}$  denotes the embedding of node  $i$  for model  $M_k$ :

- $M_1$ : Node2Vec.
- $M_2$ : Mean of Graph2Gauss i.e.  $\mathbf{E}_2[i, :] = \boldsymbol{\mu}_i$  where the Graph2Gauss mapping transforms node  $i$  into the Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i))$ .
- $M_3$ : Spectral embedding with  $k$  largest eigenvectors.

- 0 ☐ a) We modify the attributed graph such that all nodes have the same features, i.e. the adjacency matrix is  $\mathbf{A}' = \mathbf{A}$  and  
1 ☐ the new node attributes are  $\mathbf{X}'$  such that  $\mathbf{X}'[i, :] = \mathbf{X}'[j, :]$  for all  $(i, j)$ . For which model will the new node embeddings  
2 ☐  $\mathbf{E}'_k = M_k(G', \mathbf{X}')$  be different from the embeddings obtained with the original attributed graphs  $\mathbf{E}_k = M_k(G, \mathbf{X})$ ? Justify  
3 ☐ your answer.

- $\mathbf{E}_1 = \mathbf{E}'_1$ : Node2Vec does not account for node features.
- $\mathbf{E}'_2$  such that  $\mathbf{E}'_2[i, :] = \mathbf{E}'_2[j, :]$ : Graph2Gauss encodes a node based on its features only.
- $\mathbf{E}_3 = \mathbf{E}'_3$ : Spectral embedding does not account for node features.

- 0 ☐ b) We modify the attributed graph such that the graph is a clique, i.e. the new adjacency matrix is  $\mathbf{A}' = \mathbf{1} - \mathbf{I}$   
1 ☐ where  $\mathbf{1}$  is the all-ones matrix, and the node attributes are  $\mathbf{X}' = \mathbf{X}$ . For which model will the new node embeddings  
2 ☐  $\mathbf{E}'_k = M_k(G', \mathbf{X}')$  be different from the embeddings obtained with the original attributed graph  $\mathbf{E}_k = M_k(G, \mathbf{X})$ ? Justify  
3 ☐ your answer.

- $E'_1$  such that  $E'_1[i, :] = E'_1[j, :]$ : All nodes have the same neighborhood.
- $E'_2 \neq E_2$ : The loss used for Graph2Gauss training depends on the graph edges.
- $E'_3$  such that  $E'_3[i, :] = E'_3[j, :]$ : All nodes have the same neighborhood.

Sample Solution

## Problem 9: Embeddings & Ranking (Version D) (6 credits)

We consider a graph  $G = (V, E)$  with adjacency matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  where  $A_{ij} = \mathbb{1}_{(i,j) \in E}$  indicates if an edge exist between node  $i$  and node  $j$  in the graph  $G$ . The node features are represented by the matrix  $\mathbf{X} \in \mathbb{R}^{n \times D}$ . We consider the three following models  $M_k, k \in \{1, 2, 3\}$  which produce node embeddings  $\mathbf{E}_k = M_k(G, \mathbf{X}) \in \mathbb{R}^{n \times D'}$ . The vector  $\mathbf{E}_k[i, :] \in \mathbb{R}^{D'}$  denotes the embedding of node  $i$  for model  $M_k$ :

- $M_1$ : Node2Vec.
- $M_2$ : Mean of Graph2Gauss i.e.  $\mathbf{E}_2[i, :] = \boldsymbol{\mu}_i$  where the Graph2Gauss mapping transforms node  $i$  into the Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i))$ .
- $M_3$ : Spectral embedding with  $k$  largest eigenvectors.

- 0 ☐ a) We modify the attributed graph such that all nodes have the same features, i.e. the adjacency matrix is  $\mathbf{A}' = \mathbf{A}$  and  
 1 ☐ the new node attributes are  $\mathbf{X}'$  such that  $\mathbf{X}'[i, :] = \mathbf{X}'[j, :]$  for all  $(i, j)$ . For which model will the new node embeddings  
 2 ☐  $\mathbf{E}'_k = M_k(G', \mathbf{X}')$  be different from the embeddings obtained with the original attributed graphs  $\mathbf{E}_k = M_k(G, \mathbf{X})$ ? Justify  
 3 ☐ your answer.

- $\mathbf{E}_1 = \mathbf{E}'_1$ : Node2Vec does not account for node features.
- $\mathbf{E}'_2$  such that  $\mathbf{E}'_2[i, :] = \mathbf{E}'_2[j, :]$ : Graph2Gauss encodes a node based on its features only.
- $\mathbf{E}_3 = \mathbf{E}'_3$ : Spectral embedding does not account for node features.

- 0 ☐ b) We modify the attributed graph such that the graph is a clique, i.e. the new adjacency matrix is  $\mathbf{A}' = \mathbf{1} - \mathbf{I}$   
 1 ☐ where  $\mathbf{1}$  is the all-ones matrix, and the node attributes are  $\mathbf{X}' = \mathbf{X}$ . For which model will the new node embeddings  
 2 ☐  $\mathbf{E}'_k = M_k(G', \mathbf{X}')$  be different from the embeddings obtained with the original attributed graph  $\mathbf{E}_k = M_k(G, \mathbf{X})$ ? Justify  
 3 ☐ your answer.



- $E'_1$  such that  $E'_1[i, :] = E'_1[j, :]$ : All nodes have the same neighborhood.
- $E'_2 \neq E_2$ : The loss used for Graph2Gauss training depends on the graph edges.
- $E'_3$  such that  $E'_3[i, :] = E'_3[j, :]$ : All nodes have the same neighborhood.

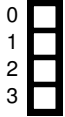
Sample Solution

## Problem 10: Semi-Supervised Learning (Version A) (6 credits)

In this problem, we consider a Stochastic Block Model with two ground-truth communities  $C_1$  and  $C_2$ . The SBM has community proportions  $\pi$  and edge probability  $\nu$  given as

$$\pi = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} \quad \text{and} \quad \nu = \begin{bmatrix} 0.2 & 0.9 \\ 0.9 & 0.2 \end{bmatrix}.$$

We consider a sampled graph  $G$  with  $n$  nodes from the SBM defined as above where the node labels are defined as the ground-truth communities of the SBM. The task is now to predict the labels of all nodes of the graphs where only a fraction of the node labels is available for training.

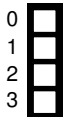


a) Do you expect label propagation with the optimization problem

$$\min \sum_{i,j} w_{ij} (\mathbf{y}_i - \mathbf{y}_j)^T (\mathbf{y}_i - \mathbf{y}_j)$$

to work well for this task? If not, propose a modification of the optimization problem which would solve the problem. Justify your answer.

No, this optimization problem assumes a homophilic graph, i.e. connected nodes are likely to have same labels. One could, for example, introduce  $\min \sum_{i,j} w_{ij} (\mathbf{y}_i - \mathbf{y}_j)^T \mathbf{H} (\mathbf{y}_i - \mathbf{y}_j)$  with  $\mathbf{H} = \begin{bmatrix} .2 & .9 \\ .9 & .2 \end{bmatrix}$  which accounts for nodes from different classes to be more likely to connect.



b) The nodes are now assigned node features sampled as

$$\mathbf{h}_v^{(0)} \sim \mathcal{N} \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \text{ for } v \in C_1 \quad \text{and} \quad \mathbf{h}_v^{(0)} \sim \mathcal{N} \left( \begin{bmatrix} -1 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \text{ for } v \in C_2.$$

We define  $N(v)$  as the 1-hop neighborhood of node  $v$ . Do you expect a one-layer GNN with the message passing step  $\mathbf{m}_v^{(1)}(\mathbf{h}_1^{(0)}, \dots, \mathbf{h}_n^{(0)}) = \frac{1}{|N(v)|} \sum_{u \in N(v)} (\mathbf{W} \mathbf{h}_u^{(0)} + \mathbf{b})$  and the update step  $\mathbf{h}_v^{(1)} = \text{ReLU}(\mathbf{Q} \mathbf{h}_v^{(0)} + \mathbf{p} + \mathbf{m}_v^{(1)})$  to work well for this task? If not, propose a modification to the message passing and/or update step that would solve the problem. Justify your answer.

No, the aggregation scheme assumes that two neighboring nodes have similar embeddings. One could e.g.  $\mathbf{m}_v(\mathbf{h}_1^{(0)}, \dots, \mathbf{h}_n^{(0)}) = \frac{1}{|N_2(v)|} \sum_{u \in N_2(v)} \mathbf{W} \mathbf{h}_u^{(0)} + \mathbf{b}$  where  $N_2(v)$  is the 2-hop neighborhood.

## Problem 10: Semi-Supervised Learning (Version B) (6 credits)

In this problem, we consider a Stochastic Block Model with two ground-truth communities  $C_1$  and  $C_2$ . The SBM has community proportions  $\pi$  and edge probability  $\nu$  given as

$$\pi = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} \quad \text{and} \quad \nu = \begin{bmatrix} 0.2 & 0.9 \\ 0.9 & 0.2 \end{bmatrix}.$$

We consider a sampled graph  $G$  with  $n$  nodes from the SBM defined as above where the node labels are defined as the ground-truth communities of the SBM. The task is now to predict the labels of all nodes of the graphs where only a fraction of the node labels is available for training.

a) Do you expect label propagation with the optimization problem

$$\min \sum_{i,j} w_{ij} (\mathbf{y}_i - \mathbf{y}_j)^T (\mathbf{y}_i - \mathbf{y}_j)$$

to work well for this task? If not, propose a modification of the optimization problem which would solve the problem. Justify your answer.

0  
1  
2  
3

No, this optimization problem assumes a homophilic graph, i.e. connected nodes are likely to have same labels. One could, for example, introduce  $\min \sum_{i,j} w_{ij} (\mathbf{y}_i - \mathbf{y}_j)^T \mathbf{H} (\mathbf{y}_i - \mathbf{y}_j)$  with  $\mathbf{H} = \begin{bmatrix} .2 & .9 \\ .9 & .2 \end{bmatrix}$  which accounts for nodes from different classes to be more likely to connect.

b) The nodes are now assigned node features sampled as

$$\mathbf{h}_v^{(0)} \sim \mathcal{N} \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \text{ for } v \in C_1 \quad \text{and} \quad \mathbf{h}_v^{(0)} \sim \mathcal{N} \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \text{ for } v \in C_2.$$

0  
1  
2  
3

We define  $N(v)$  as the 1-hop neighborhood of node  $v$ . Do you expect a one-layer GNN with the message passing step  $\mathbf{m}_v^{(1)}(\mathbf{h}_1^{(0)}, \dots, \mathbf{h}_n^{(0)}) = \frac{1}{|N(v)|} \sum_{u \in N(v)} (\mathbf{W} \mathbf{h}_u^{(0)} + \mathbf{b})$  and the update step  $\mathbf{h}_v^{(1)} = \text{ReLU}(\mathbf{Q} \mathbf{h}_v^{(0)} + \mathbf{p} + \mathbf{m}_v^{(1)})$  to work well for this task? If not, propose a modification to the message passing and/or update step that would solve the problem. Justify your answer.

No, the aggregation scheme assumes that two neighboring nodes have similar embeddings. One could e.g.  $\mathbf{m}_v(\mathbf{h}_1^{(0)}, \dots, \mathbf{h}_n^{(0)}) = \frac{1}{|N_2(v)|} \sum_{u \in N_2(v)} \mathbf{W} \mathbf{h}_u^{(0)} + \mathbf{b}$  where  $N_2(v)$  is the 2-hop neighborhood.

## Problem 10: Semi-Supervised Learning (Version C) (6 credits)

In this problem, we consider a Stochastic Block Model with two ground-truth communities  $C_1$  and  $C_2$ . The SBM has community proportions  $\pi$  and edge probability  $\nu$  given as

$$\pi = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} \quad \text{and} \quad \nu = \begin{bmatrix} 0.1 & 0.8 \\ 0.8 & 0.1 \end{bmatrix}.$$

We consider a sampled graph  $G$  with  $n$  nodes from the SBM defined as above where the node labels are defined as the ground-truth communities of the SBM. The task is now to predict the labels of all nodes of the graphs where only a fraction of the node labels is available for training.

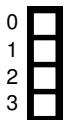


a) Do you expect label propagation with the optimization problem

$$\min \sum_{i,j} w_{ij} (\mathbf{y}_i - \mathbf{y}_j)^T (\mathbf{y}_i - \mathbf{y}_j)$$

to work well for this task? If not, propose a modification of the optimization problem which would solve the problem. Justify your answer.

No, this optimization problem assumes a homophilic graph, i.e. connected nodes are likely to have same labels. One could, for example, introduce  $\min \sum_{i,j} w_{ij} (\mathbf{y}_i - \mathbf{y}_j)^T \mathbf{H} (\mathbf{y}_i - \mathbf{y}_j)$  with  $\mathbf{H} = \begin{bmatrix} .2 & .9 \\ .9 & .2 \end{bmatrix}$  which accounts for nodes from different classes to be more likely to connect.



b) The nodes are now assigned node features sampled as

$$\mathbf{h}_v^{(0)} \sim \mathcal{N} \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \text{ for } v \in C_1 \quad \text{and} \quad \mathbf{h}_v^{(0)} \sim \mathcal{N} \left( \begin{bmatrix} -1 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \text{ for } v \in C_2.$$

We define  $N(v)$  as the 1-hop neighborhood of node  $v$ . Do you expect a one-layer GNN with the message passing step  $\mathbf{m}_v^{(1)}(\mathbf{h}_1^{(0)}, \dots, \mathbf{h}_n^{(0)}) = \frac{1}{|N(v)|} \sum_{u \in N(v)} (\mathbf{W} \mathbf{h}_u^{(0)} + \mathbf{b})$  and the update step  $\mathbf{h}_v^{(1)} = \text{ReLU}(\mathbf{Q} \mathbf{h}_v^{(0)} + \mathbf{p} + \mathbf{m}_v^{(1)})$  to work well for this task? If not, propose a modification to the message passing and/or update step that would solve the problem. Justify your answer.

No, the aggregation scheme assumes that two neighboring nodes have similar embeddings. One could e.g.  $\mathbf{m}_v(\mathbf{h}_1^{(0)}, \dots, \mathbf{h}_n^{(0)}) = \frac{1}{|N_2(v)|} \sum_{u \in N_2(v)} \mathbf{W} \mathbf{h}_u^{(0)} + \mathbf{b}$  where  $N_2(v)$  is the 2-hop neighborhood.

## Problem 10: Semi-Supervised Learning (Version D) (6 credits)

In this problem, we consider a Stochastic Block Model with two ground-truth communities  $C_1$  and  $C_2$ . The SBM has community proportions  $\pi$  and edge probability  $\nu$  given as

$$\pi = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} \quad \text{and} \quad \nu = \begin{bmatrix} 0.1 & 0.8 \\ 0.8 & 0.1 \end{bmatrix}.$$

We consider a sampled graph  $G$  with  $n$  nodes from the SBM defined as above where the node labels are defined as the ground-truth communities of the SBM. The task is now to predict the labels of all nodes of the graphs where only a fraction of the node labels is available for training.

a) Do you expect label propagation with the optimization problem

$$\min \sum_{i,j} w_{ij} (\mathbf{y}_i - \mathbf{y}_j)^T (\mathbf{y}_i - \mathbf{y}_j)$$

to work well for this task? If not, propose a modification of the optimization problem which would solve the problem. Justify your answer.

0  
1  
2  
3

No, this optimization problem assumes a homophilic graph, i.e. connected nodes are likely to have same labels. One could, for example, introduce  $\min \sum_{i,j} w_{ij} (\mathbf{y}_i - \mathbf{y}_j)^T \mathbf{H} (\mathbf{y}_i - \mathbf{y}_j)$  with  $\mathbf{H} = \begin{bmatrix} .2 & .9 \\ .9 & .2 \end{bmatrix}$  which accounts for nodes from different classes to be more likely to connect.

b) The nodes are now assigned node features sampled as

$$\mathbf{h}_v^{(0)} \sim \mathcal{N} \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \text{ for } v \in C_1 \quad \text{and} \quad \mathbf{h}_v^{(0)} \sim \mathcal{N} \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \text{ for } v \in C_2.$$

0  
1  
2  
3

We define  $N(v)$  as the 1-hop neighborhood of node  $v$ . Do you expect a one-layer GNN with the message passing step  $\mathbf{m}_v^{(1)}(\mathbf{h}_1^{(0)}, \dots, \mathbf{h}_n^{(0)}) = \frac{1}{|N(v)|} \sum_{u \in N(v)} (\mathbf{W} \mathbf{h}_u^{(0)} + \mathbf{b})$  and the update step  $\mathbf{h}_v^{(1)} = \text{ReLU}(\mathbf{Q} \mathbf{h}_v^{(0)} + \mathbf{p} + \mathbf{m}_v^{(1)})$  to work well for this task? If not, propose a modification to the message passing and/or update step that would solve the problem. Justify your answer.

No, the aggregation scheme assumes that two neighboring nodes have similar embeddings. One could e.g.  $\mathbf{m}_v(\mathbf{h}_1^{(0)}, \dots, \mathbf{h}_n^{(0)}) = \frac{1}{|N_2(v)|} \sum_{u \in N_2(v)} \mathbf{W} \mathbf{h}_u^{(0)} + \mathbf{b}$  where  $N_2(v)$  is the 2-hop neighborhood.

Sample Solution