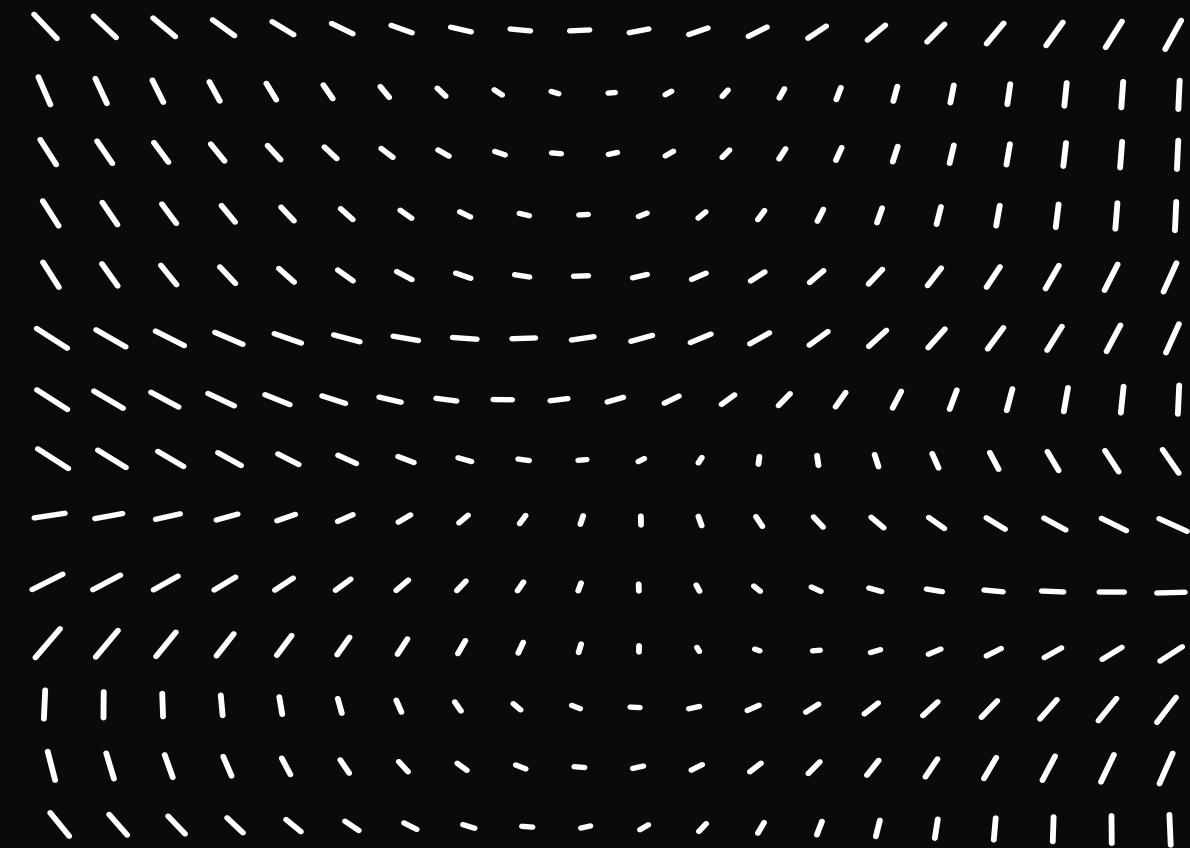


Machine learning for
graph and sequential data

Notes, exercise solutions
and previous year exams



Chapter 1 - Normalizing Flows :

probability distribution $p(x)$

$$\text{i.e. } p(x) \text{ where } x = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix}$$

- ① Data sampling : generate x_i following distr. $p(x)$
 ② Density evaluation : given any x_i compute the prob. density at this point $p(x_i)$

Main Idea

* NF can model flexible distributions for data sampling and density evaluation.

- idea → Normalizing flows are based on change of variable
- useful in simple computation like integration.

The probability in input space should be the same as in the outer space

$$p_2(x) = p_1(z) \left| \frac{\partial z}{\partial x} \right|$$

ex 1) Shift $p_1(z) = \text{Unif}(0, 1)^2$, $f(z) = z + \text{Shift} = x$
 $p_2(x) = p_1(z)$

ex 2) Linear $p_1(z) = \text{Unif}(0, 1)^2$, $f(z) = Mz = x$
 $M = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$

$$p_2(x) = p_1(z) \frac{1}{\det(M)}$$

Change of Variables : if $D \in \mathbb{N}$, $p_1(z)$ a D -dimensional distribution $f(z) = x$ is an invertible and differentiable transformation, then distribution $p_2(x)$ is

$$p_2(x) = p_1(f^{-1}(x)) \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right|$$

- determinant term : accounts for "distortion rate"
- $\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}}$: Jacobian ; $D \times D$ matrix $\frac{\partial \mathbf{f}^{-1}(\mathbf{x})}{\partial \mathbf{x}} = \left(\frac{\partial f_i^{-1}(x)}{\partial x_j} \right)$
- transformation f should be valid (invertible and differentiable).

- # conditions: 1 sufficient condition for valid transform
- invertibility : ↗ if p & o/p mapping have same dimension
 - if $D=1$, sufficient if f is strictly monotonic
 - if the transformation f is linear, its determinant should be non-zero.
 - differentiability ↗ Both f & f^{-1} are continuously differentiable can be violated in practice, but we assume this.

stacking: stacking transformations f_i :

apply change of variable multiple times

- the first variables z_0 is transformed by $z_1 = f(z_0)$
- the variables z_{i-1} is transformed $z_i = f(z_{i-1})$
- the last variable z_{K-1} is transformed by

$$\mathbf{z} = \mathbf{z}_K = f_K(z_{K-1})$$

$$p_K(\mathbf{x}) = p_0(z_0) \prod_{i=1}^K \left| \det \left(\frac{\partial f_i^{-1}(z_i)}{\partial z_i} \right) \right|$$

↙
Log version

$$\log(p_K(\mathbf{x})) = \log(p_0(z_0)) + \sum_{i=1}^K \log \left| \det \left(\frac{\partial f_i^{-1}(z_i)}{\partial z_i} \right) \right|$$

Forward and Reverse Parametrization:

The change of variable does a mapping w/o two distro.

$$p_2(x) = p_1(f^{-1}(x)) \cdot \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right|$$

- **Reverse Parametrization**: evaluate $p_2(x)$ at x
- **Forward Parametrization**: sample points from $p_2(x)$

Reverse Parametrization:

parametrize inverse transform $g = f^{-1}$

- $g_\phi(x) = z$ is computable and parameters ϕ can be learned.

$f = g^{-1}$ might not be easily computable

$$p_2(x) = p_1(g_\phi(x)) \left| \det \left(\frac{\partial g_\phi(x)}{\partial x} \right) \right|$$

Stacking: $g_\phi = g_{\phi_1} \circ \dots \circ g_{\phi_K}$

to compute density:

1. For any $x^{(j)}$, we can set $x^{(j)} = \pi_K$
2. Compute transformations $z_{i=1}^{(j)} = g_{\phi_i}(z_{i-1}^{(j)})$ and $\left| \det \left(\frac{\partial g_{\phi_i}(z_i^{(j)})}{\partial z_i^{(j)}} \right) \right|$
3. $z_0^{(j)}$, we compute $p_0(z_0^{(j)})$ (gaussian or uniform) and thus $p_X(x^{(j)})$

► We aim to learn the transformation g_ϕ

► We can call the distro. $p_\phi(x)$ to make the dependency on ϕ clear.

Learning :

$$\mathcal{D} = \{x^{(j)}\}_{j=1}^N \quad (\text{considering iid})$$

$$\max_{\phi} \log p_{\phi}(D) = \max_{\phi} \frac{1}{N} \sum_{x^{(j)} \in D} \log p_{\phi}(x^{(j)})$$

Forward Parameterization :

- Parameterize the transformation f that we know analytically
 - $f_{\theta}(z) = x$ is computable and θ is learned.

- We know that the inverse f^{-1} exists, but we might not know it analytically
 - $f_{\theta}^{-1}(x) = z$ might not be (easily) computable

- The change of variable

$$p_2(x) = p_1(z) \left| \det \left(\frac{\partial f_{\theta}(z)}{\partial z} \right)^{-1} \right|$$

- The formula only uses the known parameterization function f_{θ}

- given a sample $z^{(j)}$, we can compute samples $x^{(j)} \sim p_2(x) \neq p_1(x^{(j)})$

- We can also stack transformations and forward param.
 $f_{\theta} = f_{\theta_N} \circ \dots \circ f_{\theta_1}$

- forward parameterization enables sampling from distro $p_{\theta}(x)$

1. sampling $z_0^{(j)} \sim p_0(z)$ (eg. gaussian & uniform)

2. Compute the transform.

$$z_i^{(j)} = f_{\theta_i}(z_{i-1}^{(j)}) \text{ and } \left| \det \left(\frac{\partial f_{\theta_i}(z_{i-1}^{(j)})}{\partial z_i^{(j)}} \right) \right|$$

3. For the particular sample $x^{(j)} = z_N^{(j)}$, we can compute $p_{\theta}(x^{(j)})$

Forward pointer ; This is exactly what we need in VI.

- sample x from q , and compute prob. $q(x)$

Jacobian Determinant Computation :

$$Jg = \begin{bmatrix} \frac{\partial g_1(x)}{\partial x_1} & \cdots & \frac{\partial g_1(x)}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n(x)}{\partial x_1} & \cdots & \frac{\partial g_n(x)}{\partial x_n} \end{bmatrix}$$

$$\textcircled{1} \quad |A^{-1}| = \frac{1}{|A^{-1}|}$$

$$\textcircled{2} \quad |A| = \prod_{i=1}^n a_i$$

$$\textcircled{3} \quad \begin{vmatrix} B & 0 \\ 0 & C \end{vmatrix} = |B| |C|$$

Diagonal Jacobian

$$\textcircled{1} \quad Jg = \begin{bmatrix} \frac{\partial g_1(x_1)}{\partial x_1} & \cdots & 0 \\ 0 & \ddots & \vdots \\ 0 & \cdots & \frac{\partial g_D(x_D)}{\partial x_D} \end{bmatrix}$$

$$\textcircled{2} \quad |Jg| = \prod_{i=1}^D \frac{\partial g_i(x_i)}{\partial x_i}$$

Triangular Jacobian

$$\textcircled{1} \quad Jg = \begin{bmatrix} \frac{\partial g_1}{\partial x_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ \frac{\partial g_D(x_1, \dots, x_D)}{\partial x_D} & \cdots & \frac{\partial g_D(x_1, \dots, x_D)}{\partial x_D} \end{bmatrix}$$

$$g(x) = \begin{bmatrix} g_1(x_1) \\ \vdots \\ g_D(x_1, \dots, x_D) \end{bmatrix}$$

$$\textcircled{2} \quad |Jg| = \prod_{i=1}^D \frac{\partial g_i(x)}{\partial x_i}$$

Full Jacobian

$$Jg = L\nu$$

$$|Jg| = |L| |\nu|$$

Questions

1. Is $f(z) = 1 - z$ a valid transformation?

1) $f(z) = 1 - z$ is bijective mapping and differentiable

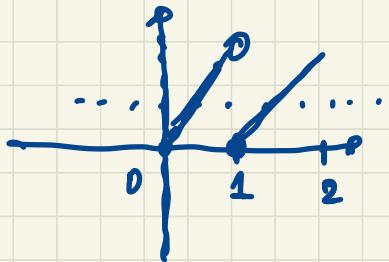
$$z = 1 - f(z) \quad f^{-1}(z) = \underline{f(x) = 1 - x}$$

2. Is $f(z) = 2 - 3z$ a valid transformation?

2) $f(z) = 2 - 3z$ if bijective and is cont. and differentiable $\forall z$.
 \therefore it is a valid mapping

3. Is $f(z) = \begin{cases} -z, z \in [0,1[\\ 1 - z, z \in [1,2] \end{cases}$ a valid transformation?

3)



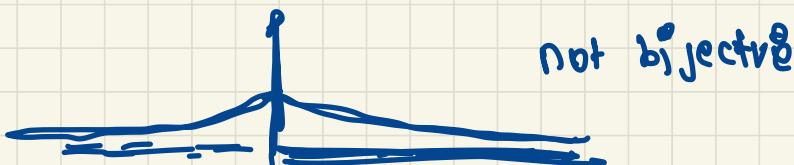
the transform is not a bijection, also it is not cont. & differentiable hence not a valid transform

- For which x is it possible to compute $p(x)$ with the forward parametrization?
- i) the point x can be used to calculate the density for a point \tilde{x} for which corresponding $\tilde{z} = g^{-1}(\tilde{x})$
- Suppose the forward transformation is defined as $f(z) = \exp(-z^n)$. What is the corresponding inverse transformation $g(x)$? Does $g(x)$ exist for any n ?

2)

There will be two cases, $n \rightarrow \begin{cases} \text{even} \\ \text{odd.} \end{cases}$

for $n \rightarrow$ even



- ~~3)~~ Suppose the forward transformation is defined as the sigmoid function. What is the corresponding inverse transformation?

3)

$$f(z) : \frac{1}{1+e^{-z}} \quad \text{continuous and differentiable}$$

$$1+e^{-z} = \frac{1}{f(z)}$$

$$1 - \frac{1}{f(z)} = e^{-z}$$

$$\ln\left(1 - \frac{1}{f(z)}\right) = -z$$

$$z = \ln(f(z)) - \ln(1-f(z))$$

1. Let's assume you get the following Jacobian:

How expensive is it to compute the determinant?
Can you comment on this in the context of NFs?

$$\begin{bmatrix} \frac{\partial g_1(x_1)}{\partial x_1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ \frac{\partial g_D(x_1, \dots, x_D)}{\partial x_1} & \dots & 0 \end{bmatrix}$$

determinant is = 0
 $\therefore 0(1)$

2. What is the complexity to compute the Jacobian determinant of an arbitrary valid transformation?

$O(D^3) \rightarrow$ Complexity of computing an arbitrary determinant.

- 3) Considering high-dimensional data (i.e. D is high), what type of Jacobian would you use?

3) triangular / diagonal Jacobian

Problem 1: Consider the following transformation $f: \mathbb{R}^3 \rightarrow \mathbb{R}^3$

$$f(z) = \begin{bmatrix} 2z_1 \\ e^{z_1} z_2 \\ e^{-z_1-z_2} z_3 \end{bmatrix}.$$

Prove or disprove whether the transformation f is invertible.

$$f\left(\begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix}\right) = \begin{bmatrix} 2z_1 \\ e^{z_1} z_2 \\ e^{-z_1-z_2} z_3 \end{bmatrix}$$

$$f(z_2) = e^{z_1} z_2$$

$$z_2 = e^{-z_1} \cdot f(z_2).$$

$$z_2 = e^{-f(z_1)/2} f(z_2)$$

$$f^{-1}(z_2) = e^{-z_1/2} z_2 \quad \text{--- (2)}$$

$$f(z_1) = 2 z_1$$

$$f^{-1}(z_1) = z_1/2 \quad \text{--- (1)}$$

$$f(z_3) = e^{-z_1-z_2} z_3$$

$$z_3 = e^{z_1+z_2} \cdot f(z_3)$$

$$z_3 = f(z_3) e^{\frac{f(z_1)}{2}} + e^{-\frac{f(z_1)}{2}} \cdot f(z_2)$$

$$f^{-1}(z_3) = z_3 e^{\frac{z_1}{2}} + e^{-z_1/2} z_2 \quad \text{--- (3)}$$

$$f^{-1}\left(\begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix}\right) = \begin{bmatrix} z_1/2 \\ e^{-z_1} z_2 \\ z_3 e^{z_1/2} + e^{-z_1/2} \cdot f(z_2) \end{bmatrix}$$

There inverse is a valid mapping

Problem 2: Consider the following transformation $f: \mathbb{R}^2 \rightarrow \mathbb{R}^2$:

$$f(z) = \begin{bmatrix} z_1^2 z_2 \\ z_2^3 \end{bmatrix}.$$

Prove or disprove whether the transformation f is invertible.

$$f\left(\begin{bmatrix} z_1 \\ z_2 \end{bmatrix}\right) = \begin{bmatrix} z_1^2 z_2 \\ z_2^3 \end{bmatrix}$$

Therefore it is not a valid transformation



violates the one to one mapping condition

$$f(z_2) = z_2^2$$

$$f(z_2) = \pm \sqrt{z_2}$$

not a unique image,

Problem 3: Consider the transformation $f(z) = Az + b$ from \mathbb{R}^2 to \mathbb{R}^2 , where $A \in \mathbb{R}^{2 \times 2}$ and $b \in \mathbb{R}^2$. Under what conditions on A and b is this transformation invertible? Justify your answer.

$$f(z) = Az + b$$

$$\frac{\partial f(z)}{\partial z} = A^T$$

$$\det\left(\frac{\partial f(z)}{\partial z}\right) = \det(A^T) = \det(A)$$

necessary condition of invertibility in case of linear transform is that $\det(J) \neq 0$.

$$\begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} \neq 0 \Rightarrow$$

$$a_{11}a_{22} - a_{12}a_{21} \neq 0.$$

Problem 4: We consider the following forward transformation $f: \mathbb{R}^3 \rightarrow \mathbb{R}^3$

$$\mathbf{x} = f(\mathbf{z}) = \begin{bmatrix} z_1 \\ e^{z_1} z_2 \\ \sqrt[3]{e^{-z_1} z_3 + z_1^2} \end{bmatrix}.$$

We assume a uniform base distribution $p_1(\mathbf{z}) = U([0, 2]^3)$. Evaluate the density $p_2(\mathbf{x})$ at the points

$$\mathbf{x}^{(1)} = \begin{bmatrix} 0 \\ 1 \\ \frac{1}{3} \end{bmatrix} \text{ and } \mathbf{x}^{(2)} = \begin{bmatrix} -1 \\ 2 \\ 1 \end{bmatrix}.$$

$$4) \quad p_2(\mathbf{x}) = p_1(f^{-1}(\mathbf{x})) \left| \det \frac{\partial f^{-1}(\mathbf{x})}{\partial \mathbf{z}} \right|$$

$$\begin{aligned} z_1 &= x_1 \\ z_2 &= x_2 e^{-x_1} \\ z_3 &= (x_3^3 - x_1^2) e^{x_1} \end{aligned}$$

$$f^{-1}(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 e^{-x_1} \\ e^{x_1} (x_3^3 - x_1^2) \end{bmatrix}$$

$$\frac{\partial f^{-1}(\mathbf{x})}{\partial \mathbf{z}} = \begin{bmatrix} 1 & 0 & 0 \\ -x_2 e^{-x_1} & e^{-x_1} & 0 \\ x & x & 3x_3^2 e^{x_1} \end{bmatrix}$$

$$\left(e^{-x_1}, 3x_3^2 e^{x_1} \right) = \underline{\underline{3x_3^2}}$$

$$f^{-1}(x^{(1)}) = \begin{bmatrix} 0 \\ 1 \\ 1/27 \end{bmatrix} = z^{(1)}$$

$$f^{-1}(x^{(2)}) = \begin{bmatrix} -1 \\ 2e \\ 0 \end{bmatrix} = z^{(2)}$$

$$p_2(x^{(1)}) = p_1(f^{-1}(x^{(1)})) \left| \det \frac{\partial f^{-1}(x^{(1)})}{\partial x} \right|$$

$$\therefore \frac{1}{2^3} \cdot 3 \cdot \frac{1}{3^2} = \boxed{\frac{1}{24}}$$

$$p_2(x^{(2)}) = p_2(f^{-1}(x^{(2)})) \left| \det \frac{\partial f^{-1}(x^{(2)})}{\partial x} \right|$$

$$= \frac{1}{2^3} \cdot 2 \cdot 0 \Rightarrow \boxed{0}$$

Problem 5: We consider the following forward transformation $x = f(z) = \sum_{k=1}^K \sigma(kz)$ from \mathbb{R} to $(0, K)$ with $\sigma(z) = \frac{1}{1+e^{-z}}$. We assume a Gaussian base distribution $p_1(z) = \mathcal{N}(0, 1)$. We sampled one point from the base distribution $z^{(1)} = 0$. Compute the corresponding sample $x^{(1)}$ from the transformed distribution and evaluate its density $p_2(x^{(1)})$.

5) $f(z^{(1)}) = \sum_{k=1}^K \sigma(K \cdot 0) = \frac{K_0}{2}$

$$\left| \det \left(\frac{\partial f(z)}{\partial z} \right) \right| = f'(z) = \sum_{k=1}^K K \sigma(Kz) (1 - \sigma(Kz))$$

$$p_2(x^{(1)}) = p_1(z^{(1)}) \left| \det \left(\frac{\partial f(z)}{\partial z} \right) \right|^{-1}$$

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{(0)^2}{2}}$$

$$\frac{1}{\sqrt{2\pi}} \frac{4}{\sum_{k=1}^K k \sigma(k \cdot 0) (1 - \sigma(k \cdot 0))}$$

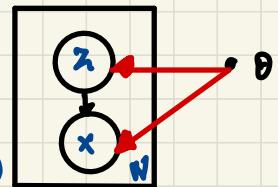
$$= \frac{1}{\sqrt{2\pi}} \frac{4}{\sum_{k=1}^K k}$$

$$\boxed{\frac{8}{\sqrt{2\pi} K(K+1)}}$$

Chapter 2 : Variational Inference

Latent variable Models (LVMs)

- 1) generate latent variable $z \sim p_\theta(z)$
- 2) generate data x conditional $x \sim p_\theta(x|z)$ on z



Joint probability distribution : $p_\theta(x, z) = p_\theta(z)p_\theta(x|z)$

Marginal Distribution : $p_\theta(x) = \int p_\theta(x, z) dz = \underbrace{\left[p_\theta(z) p_\theta(x|z) \right]}_{= \mathbb{E}_{z \sim p_\theta(z)} [p_\theta(x|z)]}$

Tasks in LVMs :

- 1) Inference : $p_\theta(z|x) = \frac{p_\theta(x|z)p_\theta(z)}{p_\theta(x)}$
- 2) Learning : $\max_{\theta} \log p_\theta(x) = \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_\theta(x_i)$
- iid assumption

• Normally, $\int p_\theta(x, z) dz$ doesn't have a closed form solution, numerical integration is infeasible.

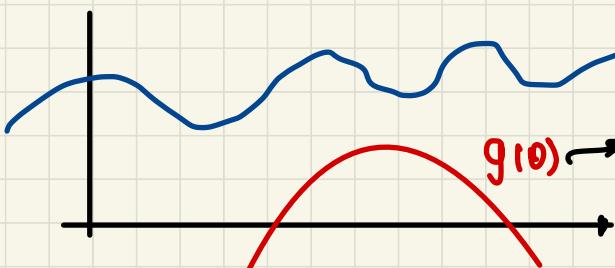
We can use Normalizing flows :

using reverse parametrization with parametric function $g_\theta(x)$ and base distribution p_1 we obtain
↳ to evaluate $p_\theta(x)$ at a point.

reverse parametrization with parametric function $g_\theta(x)$ and base distribution p_1 we obtain

$$\max_{\theta} \log p_\theta(x) = \max_{\theta} \left[\log p_1(g_\theta(x)) + \log \left| \det \left(\frac{\partial g_\theta(x)}{\partial x} \right) \right| \right]$$

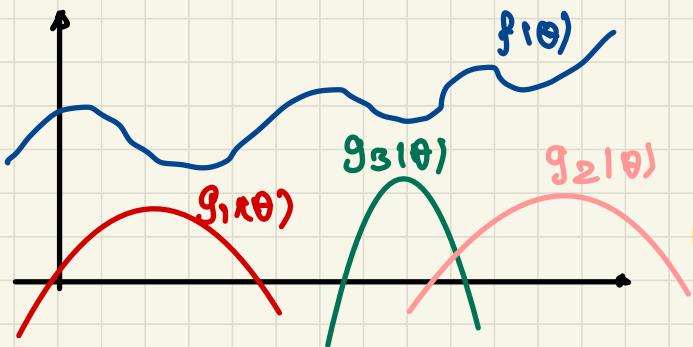
Maximization using Lower Bounds :



$\rightarrow \int_2 \nabla f(\theta)$ are intractable
 $f(\theta)$. Intractable
 find a nice $g(\theta)$ that is lower bound of f .

- Maximizing $g(\theta)$ would give us lower bound on the solution

$$\max_{\theta} f(\theta) \geq \max_{\theta} g(\theta)$$



using multiple lower bounds is even better.

$$\max_{\theta} f(\theta) \geq \max_{\theta} \max_{g \in G} g(\theta)$$

algorithm : approx. solving $\max_{\theta} f(\theta)$ intractable

1. construct lower bound $g(\theta)$ s.t $f(\theta) \geq g(\theta)$
 $+ g \in G$ of θ
2. solve optimization problem $\max_{g \in G, \theta} g(\theta)$

Lower Bound for Marginal log likelihood

$$\log p_\theta(x) = \mathbb{E}_{z \sim q_\theta(z)} \left[\log \frac{p_\theta(x, z)}{q(z)} \right] + \text{KL}(q(z) || p_\theta(z|x))$$

arbitrary dist.

Kullback - Leibler Distribution

$$KL(q(z) || p(z)) = \int q(z) \log \frac{q(z)}{p(z)} dz$$

Properties:

- asymmetric $KL(q(z) || p(z)) \neq KL(p(z) || q(z))$
- non negative
- if $KL(p(z) || q(z)) = 0 \Rightarrow p(z) = q(z)$ almost everywhere

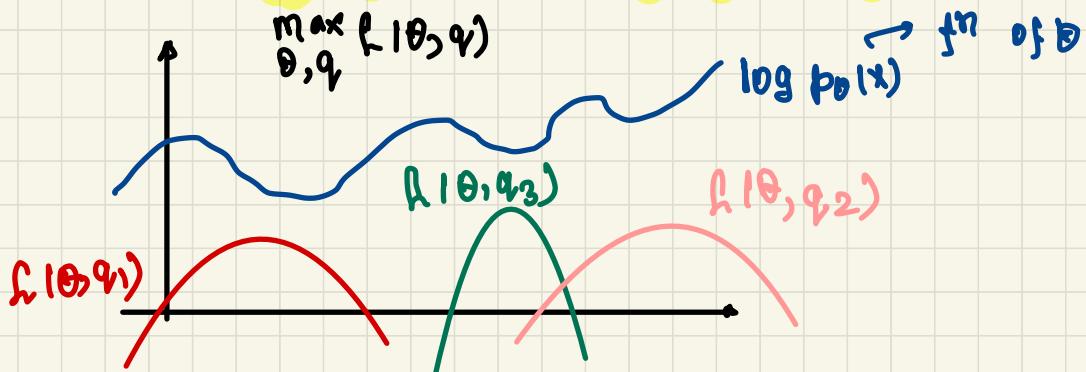
Evidence Lower Bound (ELBO)

$$\log p_{\theta}(x) = \underbrace{\mathbb{E}_{z \sim q(z)} \left[\log \frac{p_{\theta}(x, z)}{q(z)} \right]}_{\text{lower bound}} + \underbrace{KL(q(z) || p(z))}_{\geq 0}$$

* Tightness depends on how close $q(z)$ is to posterior $p(z|x)$ in terms of KL divergence

Note: Different choices of $q(z)$ leads to different lower bounds

objective: param θ and distribution $q(z)$ that maximizes lower bound

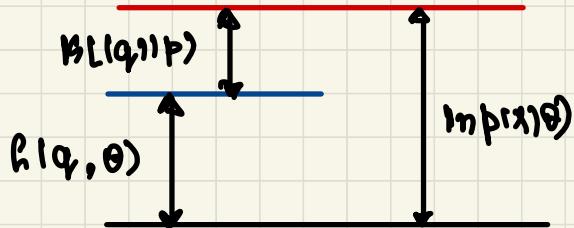
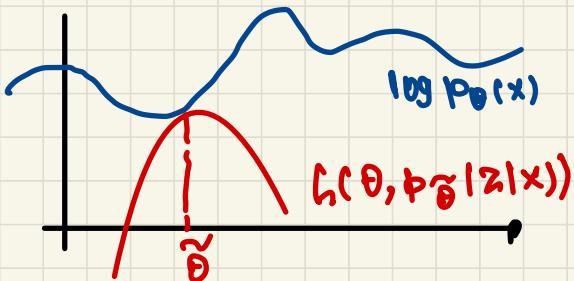


Alternative interpretation of the ELBO

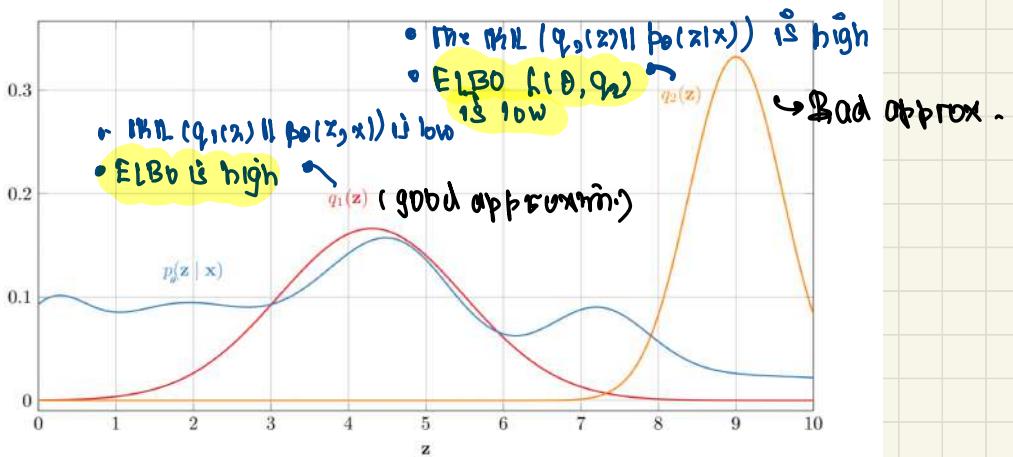
$$\log p_\theta(x) = L(\theta, q) + \text{KL}(q(z) || p_\theta(z|x))$$

$$\Rightarrow L(\theta, q) = -\text{KL}(q(z) || p_\theta(z|x)) + \log p_\theta(x)$$

• for fixed θ , $q(z) = p_\theta(z|x)$
 makes ELBO exactly equal to $\log p_\theta(x)$



for a fixed θ : maximizing the ELBO w.r.t q is equivalent to making $q(z)$ as close as possible to $p_\theta(z|x)$



EM algorithm and Variation Inference

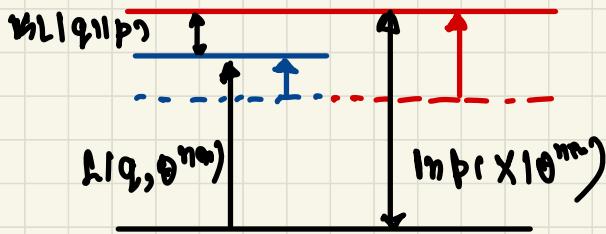
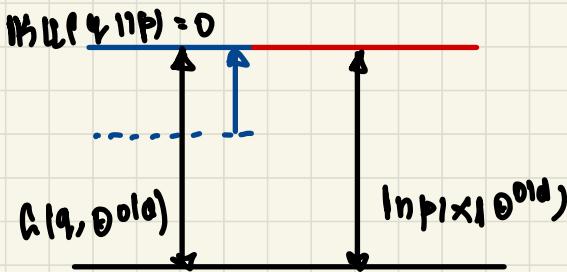
- true posterior $p_\theta(z|x) \rightarrow$ intractable
- EM algorithm
 - E Step $q(z) = p_\theta(z|x)$
 - M Step $\theta^{\text{new}} = \underset{\Theta}{\operatorname{arg\max}} \mathbb{E}_{z \sim q(z)} [\log p_\theta(x, z)]$

E - Step

$$\begin{aligned} q(z) &= p_\theta(z|x) \\ &= \underset{q}{\operatorname{arg\min}} \mathbb{H}[q(z) || p_\theta(z|x)] \\ &= \underset{q}{\operatorname{arg\max}} L(\theta, q) \end{aligned}$$

M - Step

$$\begin{aligned} \theta^{\text{new}} &= \underset{\Theta}{\operatorname{arg\max}} \mathbb{E}_{z \sim q(z)} [\log p_\theta(x, z)] \\ &= \underset{\Theta}{\operatorname{arg\max}} L(\theta, q) \end{aligned}$$



Optimizing the ELBO :

$$\max_{\theta, q} L(\theta, q)$$

$q(z)$ is prob. distr.
 \therefore 1) What is the domain
 2) How to opt. wrt. to a prob. dist.

Parametric family of distributions!

Pick tractable family of distribution \mathcal{Q}

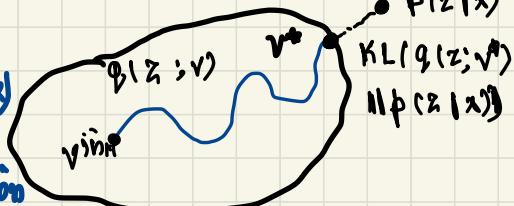
↳ 1) we can draw samples

↳ 2) every distribution in \mathcal{Q} is specified over a parameter vector ϕ .

our problem becomes:

$$\max_{\theta \in \text{IR}^M} \mathbb{E}_{z \sim q_{\theta}(z)} \left[\log p_{\theta}(x, z) - \log q_{\theta}(z) \right]$$

- 1) $\max_{\theta} \mathbb{E}_{z \sim q_{\theta}(z)} \left[\log p_{\theta}(x, z) - \log q_{\theta}(z) \right] = \min_{\theta} \text{KL}(q_{\theta}(z) \| p_{\theta}(z|x))$
2) true posterior intractable
3) Optimizing over $\phi \rightarrow$ distribution $q_{\phi} \in \mathcal{Q}$ that is closest to true posterior $p_{\theta}(z|x)$



$$\max_{\theta, \phi} \mathbb{E}_{z \sim q_{\phi}(z)} \left[\log p_{\theta}(x, z) - \log q_{\phi}(z) \right] = \max_{\theta, \phi} L(\theta, \phi)$$

Use gradient ascent $\nabla_{\theta} L(\theta, \phi)$
need to compute $\nabla_{\phi} L(\theta, \phi)$

In some cases integral can be calculated, but in case it can't be

→ assume ϕ is fixed
④ Approximating $\nabla_{\theta} L(\theta, \phi)$: $\mathbb{E}_{z \sim q_{\phi}(z)} [f_{\theta}(z)] = \frac{1}{S} \sum_{i=1}^S f_{\theta}(z_i)$

$$\nabla_{\theta} \mathbb{E}_{z \sim q_{\phi}(z)} [f_{\theta}(z)] = \nabla_{\theta} \int q_{\phi}(z) f_{\theta}(z) dz$$

$$= \int q_{\phi}(z) \nabla_{\theta} f_{\theta}(z) dz = \mathbb{E}_{z \sim q_{\phi}(z)} [\nabla_{\theta} f_{\theta}(z)] = \frac{1}{S} \sum_{i=1}^S \nabla_{\theta} f_{\theta}(z_i) \quad z_i \sim q_{\phi}(z)$$

② Approximating $\nabla_{\phi} h(\theta, \phi)$: assume θ is fixed

$$\mathbb{E}_{z \sim q_{\phi}(z)} [h_{\phi}(z)] = \int q_{\phi}(z) h_{\phi}(z) dz$$

in our case

$$h_{\phi}(z) = \log p_{\theta}(x, z) - \log q_{\phi}(z)$$

* cannot "push" the gradient inside integral

Reparameterization Trick:

- Sampling from many distributions $q_{\phi}(z) \rightarrow$ can be represented as a deterministic transformation $T(\epsilon, \phi)$ of some base distribution $b(\epsilon)$

- $q_{\phi}(z) = N(z | \mu, R R^T)$ multivariate normal
Sampling $\xrightarrow{\text{Drawing sample from } \epsilon \sim N(0, I)}$ $z = T(\epsilon, \phi) = R\epsilon + \mu$

distribution $b(\epsilon)$ does not depend on ϕ .

$$\mathbb{E}_{z \sim q_{\phi}(z)} [h_{\phi}(z)] = \mathbb{E}_{\epsilon \sim b(\epsilon)} [h_{\phi}(T(\epsilon, \phi))]$$

$$\nabla_{\phi} \mathbb{E}_{\epsilon \sim b(\epsilon)} [h_{\phi}(T(\epsilon, \phi))] = \int b(\epsilon) \nabla_{\phi} h_{\phi}(T(\epsilon, \phi)) d\epsilon$$

$$= \mathbb{E}_{\epsilon \sim b(\epsilon)} [\nabla_{\phi} h_{\phi}(T(\epsilon, \phi))] \approx \frac{1}{S} \sum_{i=1}^S \nabla_{\phi} h_{\phi}(T(\epsilon_i, \phi))$$

$\epsilon_i \sim b(\epsilon)$

Putting it all together

1. We define a latent variable generative model for our data \mathbf{x}

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z}$$

2. We are interested in maximum likelihood estimation of model parameters $\boldsymbol{\theta}$

$$\max_{\boldsymbol{\theta}} \log p_{\theta}(\mathbf{x})$$

3. We can obtain a lower bound on $\log p_{\theta}(\mathbf{x})$ using some distribution $q(\mathbf{z})$

$$\log p_{\theta}(\mathbf{x}) \geq \mathcal{L}(\boldsymbol{\theta}, q) := \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})]$$

4. Our original optimization problem can be approximately solved as

$$\max_{\boldsymbol{\theta}} \log p_{\theta}(\mathbf{x}) \geq \max_{\boldsymbol{\theta}, q} \mathcal{L}(\boldsymbol{\theta}, q)$$

5. We pick a parametric family of variational distributions \mathcal{Q}

$$\mathcal{Q} = \{q_{\boldsymbol{\phi}}(\mathbf{z}) \text{ for } \boldsymbol{\phi} \in \mathbb{R}^K\}$$

6. Our optimization problem now becomes

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

7. We obtain the gradients $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi})$ and $\nabla_{\boldsymbol{\phi}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi})$ using Monte Carlo (with the reparametrization trick)

8. We find $\boldsymbol{\theta}^*, \boldsymbol{\phi}^*$ by maximizing our objective function using gradient ascent

Mean Field assumption

In practice we make a simplifying assumption that $q(\mathbf{z})$ factorizes:

$$q(\mathbf{z}) = \prod_{i=1}^N q_i(z_i)$$

$$\begin{aligned} & \frac{1}{N} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\sum_{i=1}^N \log p_\theta(x_i, z_i) - \log q(\mathbf{z}) \right] \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{z_i \sim q_i(z_i)} \left[\log p_\theta(x_i, z_i) - \log q_i(z_i) \right] \end{aligned}$$

Questions

- Assume that we have an LVM, where $p_\theta(x|z)$ and $p_\theta(z)$ are tractable (i.e. we can compute them). Can it happen, that we can also compute $p_\theta(z|x)$ for this model, but cannot compute $p_\theta(x)$? Why or why not?

1) $p_\theta(z|x) = \frac{p(x|z) p_\theta(z)}{p_\theta(x)}$

No, this cannot happen.

(Baye's
Theorem)

- Why is it always possible to compute $\log p_\theta(x)$ in a latent variable model, where z can take only finitely many values?

2) $\log p_\theta(x) = \sum_{c=1}^C \log p_\theta(x, z=c)$

$$= \sum_{c=1}^C \log [p_\theta(x|z=c) p(z=c)]$$

1) Slide 57: Why is it necessary for all functions $g \in \mathcal{G}$ to be lower bounds on f ? What happens if some functions in \mathcal{G} are not lower bounds?

The $\max_{g \in \mathcal{G}, \theta} g(\theta)$ will be meaningless

2) Slide 57: Can we use a similar approach if we want to approximately minimize some intractable function f ? What changes need to be done in this case?

either we can use $-f$ or we take upper bounds and minimize them.

- 3) Assume that $p_\theta(z|x)$ is a distribution on $[0, \infty)$ (e.g. exponential distribution), and our variational distribution $q(z)$ is a distribution on all of \mathbb{R} (e.g. normal distribution).
- What happens to the ELBO in this case?
 - Why is the optimization problem of maximizing the ELBO ill-defined?
 - How can we fix this problem?

$$\begin{aligned} 3) \quad & \mathbb{E}_{z \sim q(z)} [\underbrace{\log p(z|x) - \log q(z)}_{\text{ELBO}}] + \text{cons:} \\ & = \int_{-\infty}^{\infty} q(z) [\Theta]. d z + \text{cons.} \end{aligned}$$

for $x \leq 0$ $p(z|x)$ is not defined.

: ELBO does not exist

to fix this is to take $q(z)$ s.t. it is defined
in $[0, \infty)$

1. Which of the following conditions have to be satisfied by a distribution $q(\mathbf{z})$, such that it's possible to use it in variational inference (as described in our recipe on slides 80-81)

- a) We can compute the expected value of \mathbf{z} in closed form
- b) We can compute the entropy of $q(\mathbf{z})$ in closed form
- c) We can draw samples from $q(\mathbf{z})$ with reparametrization
- d) We can compute the density $\log q(\mathbf{z})$ for an arbitrary \mathbf{z}
- e) We can compute $\log q(\mathbf{z})$ for a sample \mathbf{z} drawn from $q(\mathbf{z})$
- f) The distribution can be factorized as $q(\mathbf{z}) = \prod_i q_i(z_i)$

i) only c) and e) are necessary as

We need to draw samples for $q(\mathbf{z})$
and we compute $q(\mathbf{z})$

check the ELBO formulation.

2. Think of a probabilistic model with two latent variables $z_1, z_2 \in \mathbb{R}$ and an observed variable $x \in \mathbb{R}$ (i.e. write down $p_{\theta}(x|z_1, z_2)$ and $p_{\theta}(z_1, z_2)$), where the posterior can be factorized as $p_{\theta}(z_1, z_2|x) = p_{\theta}(z_1|x)p_{\theta}(z_2|x)$.

3. Same as question 3, but now the posterior cannot be factorized.

~~Problem 1:~~ Consider the following latent variable model.

$$p_\theta(z) = \text{Exp}(z|\theta) = \begin{cases} \theta \exp(-\theta z) & \text{if } z > 0, \\ 0 & \text{else.} \end{cases}$$

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

where $x \in \mathbb{R}$ is the observed data and $z \in \mathbb{R}_+$ is the latent variable. We have observed a single data point x and now would like to maximize the marginal log-likelihood $\log p_\theta(x) = \log \left(\int p(x|z)p_\theta(z)dz \right)$ w.r.t. the model parameters $\theta \in \mathbb{R}_+$. For this we will use variational inference.

We define the following parametric family of variational distributions

$$\textcolor{red}{q_\phi(z)} = \text{Exp}(z|\phi) = \begin{cases} \phi \exp(-\phi z) & \text{if } z > 0 \\ 0 & \text{else;} \end{cases}$$

that is parametrized by $\phi \in \mathbb{R}_+$. We are interested in solving the optimization problem

$$\max_{\theta>0, \phi>0} \mathcal{L}(\theta, \phi).$$

- Assume that θ is known and fixed. Does there exist a value of ϕ such that the ELBO is tight, i.e. $\log p_\theta(x) = \mathcal{L}(\theta, \phi)$? Justify your answer.
- Write down the ELBO $\mathcal{L}(\theta, \phi)$ for the above probabilistic model $p_\theta(x, z)$ and the variational distribution $q_\phi(z)$ and simplify it as far as you can. Your final answer should be a closed-form expression (no integrals or expectations).
- Compute the gradients of the ELBO $\nabla_\theta \mathcal{L}(\theta, \phi)$ and $\nabla_\phi \mathcal{L}(\theta, \phi)$.

i) $\log p_\theta(x) = \mathcal{L}(\theta, \phi) + \underbrace{\text{KL}\left(q_\phi(z) \parallel p(z|x)\right)}$
as θ is fixed
should be zero
to satisfy the condition
for this to happen $q_\phi(z) = p(z|x)$

$$p(z|x) \propto p(x|z) \cdot p(z)$$

$$\propto \exp\left(-\frac{1}{2}(x-z)^2\right) \exp(-\theta z) \mathbf{1}(z>0)$$

$$\propto \exp\left(-\frac{1}{\theta}z^2 + (x-\theta)z\right) \mathbf{1}(z>0)$$

$$q_{\phi}(z) \propto \exp(-\phi z) \quad (z > 0)$$

$$-\phi z = -\frac{1}{2}z^2 + (x-\theta)z$$

Quadratic term on RHS

$$\begin{aligned} (b) \quad L(\Theta, \Phi) &= \mathbb{E}_{z \sim q_{\phi}(z)} [\log p(x|z) + \log p(z) - \log q(z)] \\ &= \mathbb{E}_{z \sim q_{\phi}(z)} \left[-\frac{1}{2} \log(2\pi) - \frac{1}{2}(x-z)^2 + \log(\Theta) \right. \\ &\quad \left. - \Theta z - \log(\Phi) + \Phi z \right] \\ &= \mathbb{E}_{z \sim q_{\phi}(z)} \left[-\frac{z^2}{2} + xz + \log \Theta - \Theta z - \log \Phi + \Phi z \right] \\ &\quad + \text{const.} \end{aligned}$$

$$\mathbb{E}_{z \sim q_{\phi}(z)} [z] = \frac{1}{\Phi}$$

$$= -\frac{1}{\Phi^2} + \frac{x}{\Phi} + \log(\Theta) - \frac{\Theta}{\Phi} - \log(\Phi)$$

+ const

$$= -\frac{1}{\Phi^2} + \frac{x}{\Phi} + \log(\Theta) - \frac{\Theta}{\Phi} - \log(\Phi) + \text{const.}$$

$$\text{TC) } \frac{\partial}{\partial \theta} L(\theta, \phi) = -\frac{1}{\phi} + \frac{1}{\theta}$$

$$\frac{\partial}{\partial \phi} L(\theta, \phi) = \frac{2}{\phi^3} - \frac{\theta - \theta}{\phi^2} - \frac{1}{\phi}$$

~~* Problem 2:~~ You want to draw samples from an exponential distribution with rate ϕ with reparametrization. Assume that

$$q_\phi(z) = \text{Exp}(z|\phi) = \begin{cases} \phi \exp(-\phi z) & \text{if } z > 0 \\ 0 & \text{else;} \end{cases}$$

where $\phi \in \mathbb{R}_+$.

~~* You have access to an algorithm that produces samples ϵ from an exponential distribution with unit rate, that is~~

$$b(\epsilon) = \text{Exp}(\epsilon|1) = \begin{cases} \exp(-\epsilon) & \text{if } \epsilon > 0 \\ 0 & \text{else.} \end{cases}$$

Find a deterministic transformation $T_\phi: \mathbb{R}_+ \rightarrow \mathbb{R}_+$ that converts a sample $\epsilon \sim b(\epsilon)$ into a sample from $q_\phi(z)$.

~~Now, you have access to an algorithm that produces samples u from a uniform distribution on $[0, 1]$, that is~~

$$b(u) = \begin{cases} 1 & \text{if } u \in [0, 1] \\ 0 & \text{else.} \end{cases}$$

Find a deterministic transformation $S_\phi: [0, 1] \rightarrow \mathbb{R}_+$ that converts a sample $u \sim b(u)$ into a sample $z = S_\phi(u) \sim q_\phi(z)$.

$$\text{a) } b(\epsilon) = \text{Exp}(\epsilon|1) = \begin{cases} \exp(-\epsilon) & \epsilon > 0 \\ 0 & \text{otherwise} \end{cases}$$

To transforms $\epsilon \sim \text{Exp}(1)$ into $z = T_\phi(\epsilon) \sim \text{Exp}(\phi)$

taking into consideration the cumulative density function $F_\epsilon(a) = \Pr(\epsilon \leq a)$ and CDF of z $F_z(a) = \Pr(z \leq a)$

$$\begin{aligned} F_z(a) &= \Pr(z \leq a) = \Pr(T_\phi(\epsilon) \leq a) \\ &= \Pr(\epsilon \leq T_\phi^{-1}(a)) \\ &= F_\epsilon(T_\phi^{-1}(a)) \end{aligned}$$

$$F_Z(a) = 1 - \exp(-\Phi a)$$

$$F_Z(T_\Phi^{-1}(a)) = 1 - \exp(-T_\Phi^{-1}(a))$$

$$T_\Phi^{-1}(a) = \Phi a \Rightarrow$$

$$\boxed{T_\Phi(x) = \frac{x}{\Phi}}$$

\therefore The desired transformation is

$$\boxed{T_\Phi(x) = \frac{x}{\Phi}}$$

(b) given $U \sim b(u) \rightarrow$ uniform distribution in $[0,1]$

find $S_\Phi: [0,1] \rightarrow \mathbb{R}_+$ s.t it converts sample $U \sim b(u)$ into $Z = S_\Phi(U) \sim \underline{\Phi^{-1}(Z)}$

CDF of $U \rightarrow F_U(a) = \Pr(U \leq a)$
CDF of $Z \rightarrow F_Z(a) = \Pr(Z \leq a)$

$$\begin{aligned}\Pr(Z \leq a) &= \Pr(S_\Phi(U) \leq a) \\ &= \Pr(U \leq S_\Phi^{-1}(a)) \\ &= F_U(S_\Phi^{-1}(a))\end{aligned}$$

$$1 - \exp(-\Phi a) = S_\Phi^{-1}(a)$$

$$1 - S_\Phi^{-1}(a) \approx \exp(-\Phi a)$$

$$\ln(1 - S_\Phi^{-1}(a)) \approx -\Phi a \cdot$$

$$\boxed{\sim S_\Phi(a) = -\frac{1}{\Phi} \ln(1 - a)}$$

Problem 3: You are given two distributions $q(\mathbf{z})$ and $p(\mathbf{z})$ over some random vector $\mathbf{z} \in \mathbb{R}^D$. Assume that both distributions can be factorized as

$$q(\mathbf{z}) = \prod_{i=1}^D q_i(z_i)$$

$$p(\mathbf{z}) = \prod_{i=1}^D p_i(z_i).$$

(This is equivalent to saying that each component z_i is independent of z_j for $j \neq i$ under the distributions q and p). Your task is to prove that in this case the following equality holds

$$\text{KL}(q(\mathbf{z}) \| p(\mathbf{z})) = \sum_{i=1}^D \text{KL}(q_i(z_i) \| p_i(z_i)).$$

$$\begin{aligned}
3) \quad & \text{KL}(q(\mathbf{z}) \| p(\mathbf{z})) = \int q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z} \\
& = \int \int \dots \int q_1(z_1) q_2(z_2) \dots q_D(z_D) \log \left(\prod_{i=1}^D \frac{q_i(z_i)}{p_i(z_i)} \right) dz_1 dz_2 \dots dz_D \\
& = \sum_{i=1}^D \left(\int \dots \int q_1(z_1) \dots q_{i-1}(z_{i-1}) \underbrace{\log \frac{q_i(z_i)}{p_i(z_i)}}_{=} dz_1 \dots dz_D \right) \\
& = \sum_{i=1}^D \int q_i(z_i) \ln \frac{q_i(z_i)}{p_i(z_i)} \left(\int \dots \int q_1(z_1) \dots q_{i-1}(z_{i-1}) \underbrace{q_{i+1}(z_{i+1}) \dots q_D(z_D)}_{=} dz_1 \dots dz_{i-1} dz_{i+1} \dots dz_D \right) dz_i \\
& = \boxed{\sum_{i=1}^D \int q_i(z_i) \ln \left(\frac{q_i(z_i)}{p_i(z_i)} \right) dz_i} = 1
\end{aligned}$$

Chapter 3: Variational Autoencoders

In variational inference, our optimization problem is

$$\max_{\theta, \phi} \mathbb{E}_{z \sim q_{\phi}(z)} [\log p_{\theta}(x|z) + \log p_{\theta}(z) - \log q_{\phi}(z)]$$

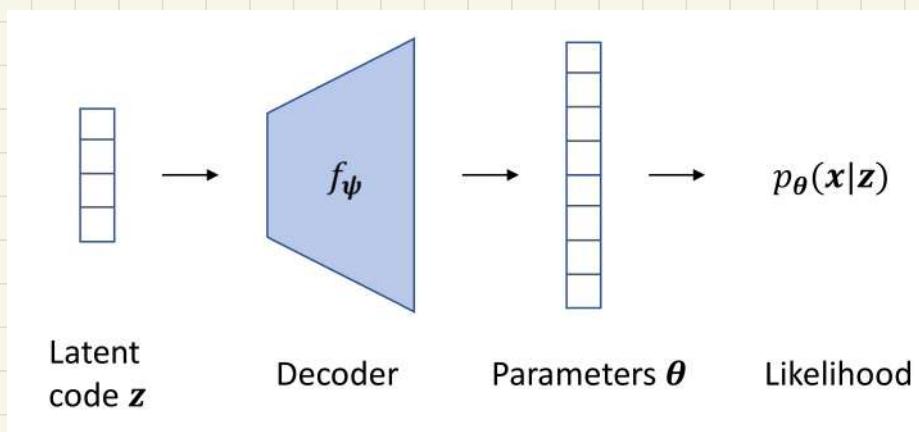
- choosing $p_{\theta}(z)$: pick simple $p(z) = N(z|0, I)$
- choosing $p_{\theta}(x|z)$: idea: pick a parametric distribution $p_{\theta}(x|z)$ whose parameters are produced by some function f_{ψ} that takes z as input.

Choosing $p_{\theta}(x|z)$ involves trade-off between expressiveness and efficiency

$$p_{\theta}(x|z) = N(x|\mu = f_{\psi}(z), I) = \prod_{j=1}^D N(x_j|\mu_j = f_{\psi}(z)_j, 1)$$

The decoder f_{ψ} :

The neural network f_{ψ} is often called "the decoder" since it converts the latent variable z into the parameters θ of $p_{\theta}(x|z)$



* Choosing $q_{\theta}(z)$:

Dataset: N iid samples)

variational distribution is over $z^{(1)} : q_{\theta}(z^{(1)}, \dots, z^{(N)})$

making the mean field assumption

$$q_{\theta}(z^{(1)}, \dots, z^{(N)}) = \prod_{i=1}^N q_{\theta(i)}(z^{(i)})$$

choice of $q_{\theta(i)}(z^{(i)})$

$\theta^{(i)}: q_{\theta(i)}(z^{(i)}) \xrightarrow{\text{diagonal}}$

① multivariate normal: $q_{\theta(i)}(x^{(i)}) = N(z^{(i)} | \mu^{(i)}, \Sigma^{(i)})$

② Normalizing flow → forward parameterization

Why is Normal $q_{\theta}(z)$ a good choice?

$$\begin{aligned} L(\theta, \phi) &= \mathbb{E}_{z \sim q_{\theta}(z)} [\log p_{\theta}(x|z) + \log p(z) - \log q_{\theta}(z)] \\ &= \mathbb{E}_{z \sim q_{\theta}(z)} [\log p_{\theta}(x|z)] - \text{KL}(q_{\theta}(z) || p(z)) \end{aligned}$$

KL divergence b/w two multivariate normal dist:

$$\text{KL}(q_{\theta}(z) || p(z)) = \frac{1}{2} (\text{tr}(\Sigma) + \mu^T \mu - \log |\Sigma| - 1)$$

even simpler for diagonal covariance $\Sigma = \text{diag}(\sigma^2)$

$$\text{KL}(q_{\theta}(z) || p(z)) = \frac{1}{2} \left(\sum_{j=1}^L (\sigma_j^2 + \mu_j^2 - \log \sigma_j^2 - 1) \right)$$

for multiple parameters:

$$L(\theta, \phi) = \frac{1}{N} \sum_i^N L_i(\theta, \phi^{(i)}) = \frac{1}{N} \sum_i^N \left[\mathbb{E}_{z^{(i)} \sim q_{\theta}^{(i)}(z^{(i)})} [\log p_{\theta}(x^{(i)} | z^{(i)})] - \text{KL}(q_{\theta}^{(i)}(z^{(i)}) || p(z^{(i)})) \right]$$

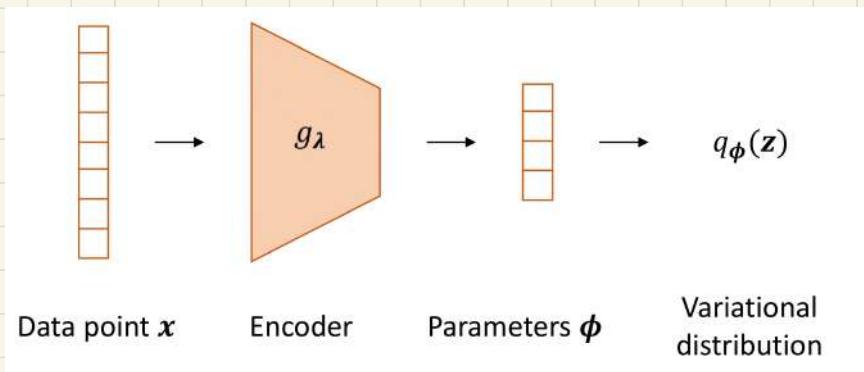
- But we have to learn separate parameters $\phi^{(i)}$ for each instance.

Ammortized Inference :

- we need to find $\phi_{\text{optimal}}^{(i)}$ for every sample $x^{(i)}$
- standard approach : $\phi_{\text{optimal}}^{(i)} = \underset{\phi^{(i)}}{\operatorname{arg\,max}} \ell_i(\theta, \phi^{(i)})$
- Better idea : train a NN g_λ that tries to map every $x^{(i)}$ in the training set to its optimal parameters $\phi_{\text{opt}}^{(i)}$

$$\max_{\lambda} \frac{1}{N} \sum_i^N \ell_i(\theta, g_\lambda(x^{(i)}))$$

Encoder g_λ



"encoder" converts data point x into the parameters ϕ that defines the distribution $q_\phi(z)$ over the latent code z .

Pulling it all together:

ELBO is a single sample

$$L(\psi, \lambda) = \mathbb{E}_{z \sim q_{\phi}(z)} [\log p_{\theta}(x|z)] - \mathbb{H}[L(q_{\phi}(z)||p(z))]$$

where $\Theta = f_{\psi}(z)$ and $\Phi = g_{\lambda}(x)$

Recipe for optimizing the ELBO:

1. Compute $\Phi = g_{\lambda}(x)$
2. Compute an MC estimate of ELBO; often done using a single sample
 - a) Draw $z' \sim q_{\phi}(z)$ with reparametrization
 - b) Compute $\Theta = f_{\psi}(z')$
 - c) ELBO: $L(\psi, \lambda) \approx \log p_{\theta}(x|z') - \mathbb{H}[L(q_{\phi}(z)||p(z))]$
3. Backprop (compute $\nabla_{\psi} L(\psi, \lambda)$ and $\nabla_{\lambda} L(\psi, \lambda)$)
4. Update the NN weights ψ and λ using gradient ascent.

reparametrization:

$$L(\psi, \lambda) = \mathbb{E}_{z \sim q_{\phi}(z)} [\log p_{\theta}(x|z)] - \mathbb{H}[L(q_{\phi}(z)||p(z))]$$
$$\Theta = f_{\psi}(z) \quad \Phi = g_{\lambda}(x)$$

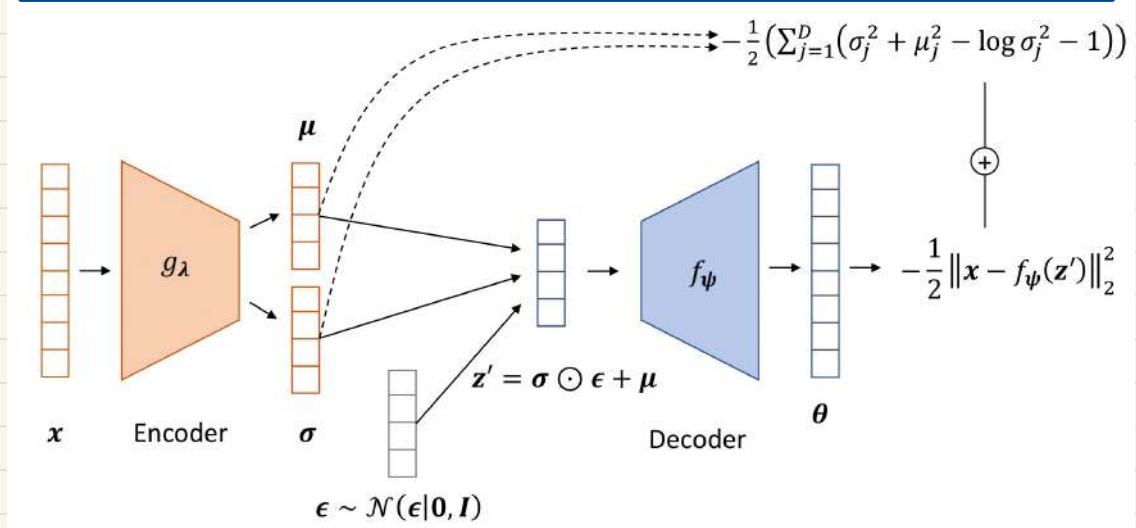
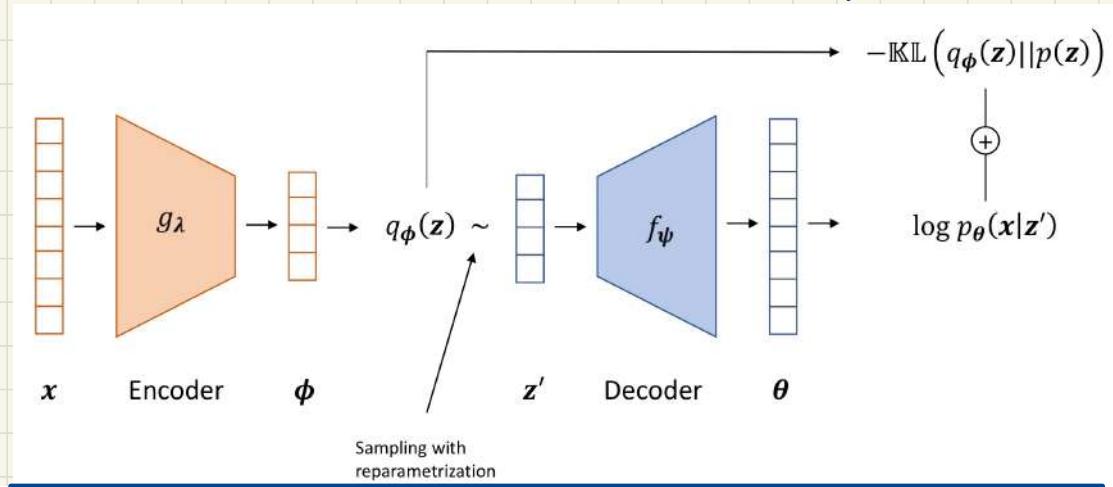
This means we need to sample from $q_{\phi}(z)$ with reparametrization

1. $\epsilon \sim b(\epsilon)$
2. $z' = T(\epsilon, \Phi)$

Variational Autoencoder (VAE)

$$L(\psi, \theta) = \mathbb{E}_{z \sim q_\phi(z)} [\log p_\theta(x|z)] - \text{KL}(q_\phi(z)||p(z))$$

$\phi = g_\lambda(x)$ and $\Theta = f_\psi(z')$



$$\ell(\psi, \lambda) = \log p_\theta(x|z')$$

$$z' = g_\lambda(x) \quad \Theta = f_\psi(z')$$

03 - Generative Adversarial Networks:

idea: Transform given initial noise using flexible function

Resulting Distribution:

- 1) Draw an initial noise vector from the prior $z \sim p(z)$
- 2) Deterministically transform z into the target data/obj space $x = f_\theta(z)$

$$p_\theta(x) = \underbrace{\frac{d}{dx_1} \dots \frac{d}{dx_d}}_{\text{intractable}} \int p(z) dz$$
$$+ f_\theta(z) \leq x^q$$

Sampling is easy; density evaluation is hard.
Without likelihood, many inference and parameter learning
tools become unavailable.

Forward pointer: generative adversarial network:

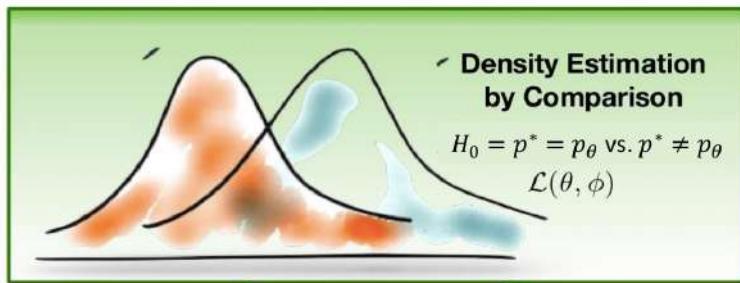
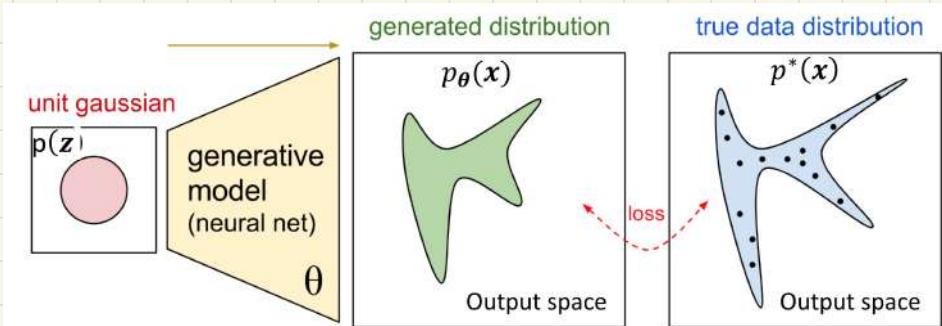
main idea:

- two competing NN models
- generator: takes noise as input and generate "fake" samples
- discriminator: receives samples from both generator and training data and has to distinguish b/w the two "real" or "fake"
- Goal: train the generator in such a way that the discriminator can not distinguish b/w "real" or "fake"

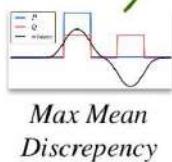
Likelihood free inference:

idea: any method that compares two set of samples
this is called density estimation by comparison.

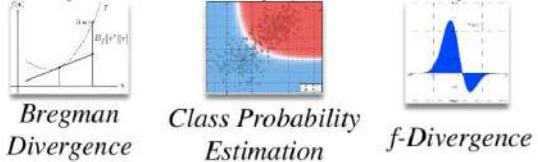
Hypothesis \rightarrow true data distribution $p^*(x)$ and the model distribution $p_\theta(x)$ are equal.



Density Difference

$$r_\phi = p^* - p_\theta$$


Density Ratio

$$r_\phi = \frac{p^*}{p_\theta}$$


Density difference : approach 1

i) Moment matching $\rightarrow p^*$ and p_θ are identical iff the expectations of any test statistic $s(x)$ are identical.

- thus, minimize the gap $\min_{\theta} (\mathbb{E}_{p^*}[s(x)] - \mathbb{E}_{p(z)}[s(f(x))])$

approx. expectation using Monte Carlo

2) Max Mean Discrepancy:

- When $\pi(x)$ are defined within a reproducing Kernel Hilbert Space, we obtain kernel based forms of these objectives
- These kernels are highly flexible and allow easy handling of data such as images, strings & graphs.

2) Ratio Based Approach

- Density ratio $\pi^*(x) = p^*(x) / p_0(x)$

cannot compute closed form approximating means solving
approximate $\pi^*(x)$ by $\pi_\phi(x)$ a learning problem

general principle 1) optimize ratio loss:
approx. true density ratio $\pi^*(x)$ (learning ϕ)

2) optimize generative loss:
drive the density ratio towards 1 (learning θ)

can be solved alternatively,

Following three methods are instantiation of this principle

1) Loss Probability estimation: density ratio classifier that distinguishes b/w observed and generated data. GAN

2) Divergence Maximization: min f-divergence b/w true density and generated density

3) Ratio Matching: Directly minimize the error b/w the true density ratio and estimate of it

$$L = \mathbb{E}_{p_\theta} [(r_\phi(x)) - r^*(x)]^2$$

Learning via Class Probability:

$$p(x|Y=1) = p^*(x) \quad \text{and} \quad p(x|Y=0) = p_\theta(x)$$

true generator

$$P(Y=1) = \pi \quad P(Y=0) = 1 - \pi$$

$$r^*(x) = \frac{p^*(x)}{p_\theta(x)} = \frac{p(Y=1|x)}{p(Y=0|x)} \cdot \frac{1-\pi}{\pi}$$

density ratio estimation = class probability estimation

**Scoring function or
discriminator** $D_\theta(x) = p(Y=1|x)$

$$L_{\theta, \phi} = \mathbb{E}_{(x,y) \sim p(x,y)} \left[-y \log [D_\phi(x)] - (1-y) \log (1 - D_\phi(x)) \right]$$

$$= \pi \mathbb{E}_{p^*(x)} [-\log D_\phi(x)] + (1-\pi) \mathbb{E}_{p(x)} [-\log [1 - D_\phi f_\theta(x)]]$$

i) solving : $\phi^*(\theta) = \underset{\phi}{\operatorname{argmin}} L_{\theta, \phi}$

leads to the "best" discrimination for given generation model.

$$r^*(x) = \frac{p^*(x)}{p_\theta(x)} = \frac{p(Y=1|x)}{p(Y=0|x)} \approx \frac{D_{\phi^*(\theta)}(x)}{1 - D_{\phi^*(\theta)}(x)}$$

(\pi \approx \bar{\pi})

aim to drive $\underline{r^*(x) \rightarrow 1}$

solving $\theta^* \approx \operatorname{argmax}_\theta L_{\theta, \phi^*}(\theta)$

• generator and discriminator play a min-max game

$$\min_{\theta} \max_{\phi} \mathbb{E}_{p_{\theta}(x)} [\log D_{\phi}(x)] + (1-\pi) \mathbb{E}_{p(x)} [\log [1 - D_{\phi}(f_{\theta}(z))]]$$

• Discriminator \rightarrow aims to distinguish b/w $p_{\theta}(x)$
min-cross entropy and $p(x)$ (max)

• generator \rightarrow aims to generate samples that
are indistinguishable. (min)

max (min cross entropy)?

Learning via Class Probability Estimation:

Ratio Loss (discriminator loss) optimization:

$$\min_{\phi} \pi \mathbb{E}_{p_{\theta}(x)} [-\log D_{\phi}(x)] + (1-\pi) \mathbb{E}_{p(x)} [-\log [1 - D_{\phi}(f_{\theta}(z))]]$$

generative loss optimization:

$$\min_{\theta} \mathbb{E}_{p(z)} [\log [1 - D_{\phi}(f_{\theta}(z))]]$$

Questions

1. Assume that each data point is represented by a vector $x \in \{1, 2, \dots, C\}^D$. What distribution would you pick for $p_{\theta}(x|z)$? How can we parametrize this distribution with a neural network?

1) Categorical distribution

$D \times C$ matrix of logits A

$$p(x_i=c | A) = \frac{\exp(A_{i,c})}{\sum_{j=1}^C \exp(A_{i,j})}$$

2. Assume that each datapoint x is represented by a variable-length sequence of real numbers $x_i \in \mathbb{R}^{D_i}$ (D_i might be different for different i 's). What NN architecture could we use for the encoder g_λ in this case?

For variable lengths using RNN would be a good choice → we can use the hidden states as $q(x)$

3. Slide 94: Assume that we choose to model $p_\theta(x|z)$ with a normalizing flow. Should we use forward or reverse parametrization? Why?

We should choose reverse parametrization as we need $p_\theta(x|z)$ for single instances

4. Slide 97: Assume that we choose to model $q_\phi(z)$ with a normalizing flow. Should we use forward or reverse parametrization? Why?

We should use forward parametrization because we want to sample from $q(z)$

5. Slide 106: How can we ensure that the vector σ produced by the encoder g_λ is always positive?

We can use element-wise strictly positive function like softplus or exp.

1. What can we say about the ratio $r^*(x) = p^*(x)/p_\theta(x)$ when:

- a) The generator and the discriminator are optimal
- b) The generator only is optimal
- c) The discriminator only is optimal

1 depends on 1 depends on generator generator

→ cannot make any remark

$$1) r^*(x) = \frac{p^*(x)}{p_\theta(x)}$$

density ratio only depends on the generator

2) For a given data x , what is the probability for the discriminator to be correct

Ambi when:

-90%

- a) The generator and the discriminator are optimal $\leq 50\%$.
- b) The generator only is optimal $\leq 50\%$.
- c) The discriminator only is optimal $\geq 80\%$.



Problem 1: Below we show pseudocode for implementing 3 autoencoder-like neural net architectures. The observed data is denoted as $x \in \mathbb{R}^D$. Here, $g_\lambda : \mathbb{R}^D \rightarrow \mathbb{R}^L$ and $f_\psi : \mathbb{R}^L \rightarrow \mathbb{R}^D$ are fully connected feedforward neural networks with learnable parameters λ and ψ . The output layers of g_λ and f_ψ have no (i.e. have linear) activation functions. \mathcal{N} denotes the normal distribution, I_N is the $N \times N$ identity matrix, and $\mathbf{0}_N$ is the vector of all zeros of length N .

For each of the architectures below, explain whether it's **necessary** to use the reparametrization trick to compute the gradient of the loss \mathcal{L} w.r.t. **both** λ and ψ . Answer "Yes" or "No" and provide a justification. If the answer is "Yes", modify the code to implement the reparametrization trick.

a) Model 1

$$\begin{aligned} z_i &\sim \mathcal{N}(x_i, I_D) \\ h_i &= g_\lambda(z_i) \\ \tilde{x}_i &= f_\psi(h_i) \\ \mathcal{L} &= \|x_i - \tilde{x}_i\|_2^2 \end{aligned}$$

a) No, the parameters do not depend on learnable parameters

b) Yes

$$\begin{aligned} h_i &= g_\lambda(x_i) \\ z_i &\sim \mathcal{N}(h_i, I_L) \\ \tilde{x}_i &= f_\psi(z_i) \\ \mathcal{L} &= \|x_i - \tilde{x}_i\|_2^2 \end{aligned}$$

$$\begin{aligned} e_i &\sim \mathcal{N}(\mathbf{0}_L, I_L) \\ \tilde{x}_i &= h_i + e_i \end{aligned}$$

b) Model 2

c) Model 3

$$\begin{aligned} h_i &= g_\lambda(x_i) \\ z_i &\sim \mathcal{N}(\mathbf{0}_L, I_L) \\ \tilde{x}_i &= f_\psi(h_i + z_i) \\ \mathcal{L} &= \|x_i - \tilde{x}_i\|_2^2 \end{aligned}$$

→ No sampling required
→ No sampling req.
no.

Problem 2: Consider the same setup as in the previous problem. The model specified below is not well defined. Your task is to find the problem with the model and modify the pseudo code to fix it.

In addition, if you think it's **necessary** to use the reparametrization trick, include it in your implementation.

$$\begin{aligned} h_i &= g_{\lambda}(x_i) \\ z_i &\sim \mathcal{N}(0_L, \text{diag}(h_i)) \\ \tilde{x}_i &= f_{\psi}(z_i) \\ \mathcal{L} &= \|x_i - \tilde{x}_i\|_2^2 \end{aligned}$$

covariance matrix should be **pos def**
i.e. Entries should be pos.

$$h_i = \exp(g_{\lambda}(x_i))$$

as x_i samples depends on λ of encoder
 we need to use the **reparametrization technique**
 when sampling z_i .

$$\begin{aligned} \varepsilon_i &\sim N(0_L, I_L) \\ \underline{z_i} &\sim \underline{\varepsilon_i} \odot \underline{h_i} \end{aligned}$$

Problem 3: The loss used in generative adversarial networks (GANs) can be written in the following form:

$$\min_{\theta} \max_{\phi} \mathcal{L}(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{p^*(\mathbf{x})} [\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{p(\mathbf{z})} [\log(1 - D_{\phi}(f_{\theta}(\mathbf{z})))]$$

where $p^*(\mathbf{x})$ is the true data distribution, $p(\mathbf{z})$ is the distribution of the noise, f_{θ} is the generator, and D_{ϕ} is the discriminator.

- a) For a given generator (fixed parameters θ) assume there exists a discriminator $D_{\phi^*}(\mathbf{x})$ with parameters ϕ^* such that for all \mathbf{x} :

$$D_{\phi^*}(\mathbf{x}) = \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + p_{\theta}(\mathbf{x})}$$

where $p_{\theta}(\mathbf{x})$ is the distribution learned by the generator. Show that D_{ϕ^*} is **optimal**, i.e. $\phi^* = \arg \max_{\phi} \mathcal{L}(\theta, \phi)$.

Hint: $\arg \max_y [a \log(y) + b \log(1 - y)] = \frac{a}{a+b}$ for any $a, b \in \mathbb{R}_0^+, a + b > 0$.

$$\max_{\phi} \mathcal{L}(\theta, \phi) = \max_{\phi} \mathbb{E}_{p^*(\mathbf{x})} [\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{p(\mathbf{z})} [\log(1 - D_{\phi}(f_{\theta}(\mathbf{z})))]$$

$$\max_{\phi} \mathbb{E}_{p^*(\mathbf{x})} [\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{p(\mathbf{z})} [\log(1 - D_{\phi}(f_{\theta}(\mathbf{z})))]$$

$$\max_{\phi} \int [p^*(\mathbf{x}) \log D_{\phi}(\mathbf{x}) + p_{\theta}(\mathbf{x}) \log(1 - D_{\phi}(\mathbf{x}))] d\mathbf{x}$$

It's not obvious to find D_{ϕ^*} from this integral but what if we found optimal discriminant for every single x ? This will achieve higher value for own optimizer

$$\leq \int \max_{D_{\phi}} [p^*(x) \log D_{\phi}(x) + p_{\theta}(x) \log (1 - D_{\phi}(x))] dx$$

$$y = D_{\phi}(x) \quad 1 = p^*(x) \quad b = p_{\theta}(x).$$

$$D_{\phi^*} = \arg \max_{D_{\phi}(x)} [p^* \log D_{\phi}(x) + p_{\theta}(x) \log (1 - D_{\phi}(x))]$$

$$= \frac{p^*(x)}{p^*(x) + p_{\theta}(x)}$$

∴ the discriminator maximizes the value for expression, as this expression is greater than our previous expression so D_{ϕ^*} should also maximize that

b) What is value of the optimal $D_{\phi^*}(x)$ when:

- The generator is optimal i.e. $p_{\theta}(x) = p^*(x)$
- The generator assigns a zero probability $p_{\theta}(x) = 0$ to a sample x whereas $p^*(x) \neq 0$
- The generator assigns a non-zero probability $p_{\theta}(x) \neq 0$ to a sample x whereas $p^*(x) = 0$

$$b) D_{\phi^*}(x) = \frac{1}{2}$$

$D_{\phi^*}(x) = 1$ such all samples as real.

$D_{\phi^*}(x) = 0$ all such samples fake.

Chapter 4: Robustness Part 1

A Adversarial Examples - Definition

- Dataset : $(x_i, y_i) \sim \mathbb{P}_{\text{data}}$, $(x_i, y_i) \in \mathbb{R}^d \times \mathcal{Y}$
- Classifier : $f: \mathbb{R}^d \rightarrow \mathcal{Y}$

specify a perturbation set $\mathcal{P}(x)$, i.e. set of perturbations which when applied to x do not change its semantic

we say $\tilde{x} \in \mathcal{P}(x)$ is an adversarial example for f at (x, y)

if
 $\begin{cases} f(x) = y & \text{i.e. } f \text{ correctly classifies } x \\ f(\tilde{x}) \neq y & \text{i.e. } f \text{ fails to correctly classify } \tilde{x} \end{cases}$

• Perturbation should not change semantic content

• L_p norm : $\mathcal{P}_{\epsilon, p}(x) = \{\tilde{x} : \|\tilde{x} - x\|_p \leq \epsilon\}$

Adversarial Attacks: Objective Function

$$\tilde{x}_*^* = \arg \max_{\tilde{x} \in \mathcal{P}(x)} \mathbb{I}_{0/1}(f(\tilde{x}), y)$$

Recall : $\mathbb{I}_{0/1}$ is the 0-1 loss
 However : $\mathbb{I}_{0/1}$ has either zero or undefined gradient

cross entropy loss \mathcal{L} is often :

$$\tilde{x}_*^* = \arg \max_{\tilde{x} \in \mathcal{P}(x)} \mathcal{L}(f(\tilde{x}), y)$$

Projected Gradient Descent:

$$\tilde{x}'_x = \arg \max_{\tilde{x} \in P(x)} L(f(\tilde{x}), y)$$

$$x_{t+1} = \Pi(x_t + \eta_t \nabla_x L(f(x_t), y))$$

- train model : But update data instead of weights

Fast Gradient-Sign Method (FGSM):

$$\tilde{x} = \Pi(x + \eta \cdot \text{sign}(\nabla_x L(f(x), y)))$$

$P(x)$ is a ball with radius ϵ measured by the L_∞ norm, setting $\eta = \epsilon$ yields

Alternative Optimization Problem:

- $\min_{\tilde{x}} D(x, \tilde{x}) \text{ s.t. } \log L(f(\tilde{x}), y) > 0$

- D is large when \tilde{x} is far from x

$$\min_{\tilde{x}} D(x, \tilde{x}) + \lambda L(\tilde{x}, y) \quad (\text{converting to unconstrained optm})$$

an effective loss f^n :

$$L(\tilde{x}, y) = \left[\sum_i \tilde{x}_i y_i - \max_i (\tilde{x}_i)_+ \right]_+$$

- $y \sim \text{original class}$
- $[\tilde{x}]_+ = \max(\tilde{x}, 0)$
- $\tilde{x}_i \approx \log f(\tilde{x}) / \log p_{\text{rob.}} \text{ of class } i$
- loss L is positive if \tilde{x} is classified as y and 0 otherwise

- Fixing a "bad" model does not seem to be a feasible approach when trying to increase the robustness

Robust Training: aim to produce models that are robust to adversarial perturbations

Common theme: optimize "worst case" loss

- let $\ell(y, \hat{y})$ be some loss (for ex. cross entropy)
- the (non-robust) training tries to find an f that minimizes the expected loss

$$R = \mathbb{E}_{(x,y) \in P_{\text{data}}} [\ell(f(x), y)]$$
- The robust version of the problem.

$$R = \mathbb{E}_{(x,y) \in P_{\text{data}}} \left[\sup_{\tilde{x} \in \mathcal{P}(x)} \ell(f(\tilde{x}), y) \right]$$

* || loss achieved by worst-case perturbation in $\mathcal{P}(x)$

Adversarial Training:

Idea: perform stochastic gradient descent (SGD) on the robust loss R_{rob} .

$$R_{\text{rob}} = \mathbb{E}_{(x,y) \in P_{\text{data}}} \left[\sup_{\tilde{x} \in \mathcal{P}(x)} \ell(f(\tilde{x}), y) \right]$$

$f = f_{\theta}$ (NN) parameterized by weights θ

$$\nabla_{\theta} R_{\text{rob}} = \mathbb{E}_{(x,y) \in P_{\text{data}}} \left[\nabla_{\theta} \left| \sup_{\tilde{x} \in \mathcal{P}(x)} \ell(f(\tilde{x}), y) \right| \right]$$

use Daskin's Theorem for calculating the gradient of worst case loss.

Daskin's Theorem : $\Delta(\theta) \rightarrow$ set of \tilde{x} for which supremum is obtained.

if $\Delta(\theta)$ contains a single element, i.e.

$\Delta(\theta) = \{\tilde{x}^*\}$, then sup of differential exists

$$\nabla_{\theta} \left(\sup_{\tilde{x} \in \mathcal{P}(x)} l(f_{\theta}(\tilde{x}), y) \right) = \nabla_{\theta} l(f_{\theta}(\tilde{x}^*), y)$$

Problem: 1) finding worst case perturbed example
 \tilde{x} is intractable.

Idea: create adversarial examples as proxy of the worst case examples.

• Adversarial training algorithm :

- 1) $(x_i, y_i) \sim \mathbb{P}_{\text{data}}$
- 2) using adversarial attack procedure, find an \tilde{x}_i with high loss $l(f_{\theta}(\tilde{x}_i), y_i)$
- 3) update the weights via gradient descent :
$$\theta \leftarrow \theta - \eta \nabla_{\theta} l(f_{\theta}(\tilde{x}_i), y_i)$$

04 - Exact certification:

goal: Develop an algorithm that answers the question
 "Is the classifier f around the sample x adversarial-free (within an ℓ -ball measured by some norm".

- NN \rightarrow sequence of functions
- each layer: $f_i(x) = \sigma(W_i x + b_i)$
- ReLU activation: $\sigma(x) = \max(0, x)$
- overall function: $F: \mathbb{R}^d \rightarrow \mathbb{R}^{l+1}$

$$F(x) = W_L f_{L-1} \circ f_{L-2} \circ \dots \circ f_1(x) + b_L$$

Theorem: Exact certification of NN with ReLU activation function and ℓ_∞ -bounded perturbation is NP Hard.

Mixed integer Linear Programming (MILP):

$$\begin{array}{ll} \min : & c^T x \\ \text{subject to :} & Ax \leq b ; x \geq 0 \end{array} \quad \left. \begin{array}{l} \text{can be solved} \\ \text{efficiently} \end{array} \right\}$$

MILP: some variables constrained to integers
 ↴ and some not.
 NP-complete

Exact certification as MILP:

- suppose our classifier predicts class $c^* = \arg \max_i F(x)_i$
- $m_t = F(x)_{c^*} - F(x)_t$: classification margin c^* b/w c^* and t
- worst case margin:

$$m_t^* = \min_{\tilde{x}} F(\tilde{x})_{c^*} - F(\tilde{x})_t$$

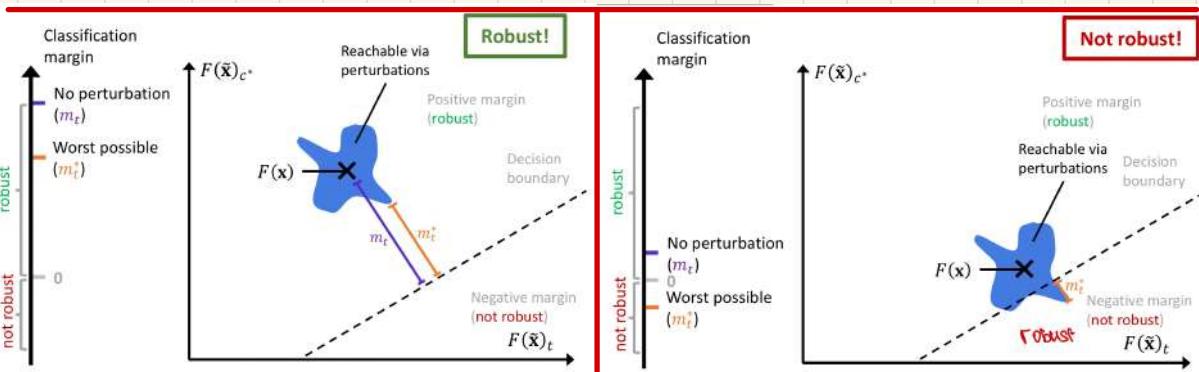
$$\text{subject to } \| \tilde{x} - x \|_p \leq \epsilon$$

$m_t^* > 0 \rightarrow$ prediction cannot be changed from c^* to t

- if $t \neq c^*$, $m_t^* > 0 \rightarrow$ certify robustness
- if for any class $t \neq c^*$, $m_t^* < 0 \rightarrow$ it's adversarial example

(1)

(2)



Optimization Problem:

$$m_t^* = \min_{\tilde{x}} [\tilde{x}^{(2)}]_{c^*} - [\tilde{x}^{(2)}]_t$$

subject to $\|\tilde{x} - x\|_p \leq \epsilon$

$$\tilde{x}^{(0)} = x$$

$$\tilde{x}^{(l+1)} = w_l \tilde{x}^{(l-1)} + b_l, \quad l=1..L$$

$$\tilde{x}^{(0)} = \text{ReLU}(\tilde{x}^{(0)}) \quad l=1..L-1$$

we need to convert to MILP encode:

- 1) L_p constraints \rightarrow adversarial perturbations
- 2) the non linear ReLU constraints

L_p constraints:

$$x_i - \tilde{x}_i \leq \epsilon \quad \forall i$$

$$\tilde{x}_i - x_i \leq \epsilon \quad \forall i$$

L_1 : straightforward

L_2 : mixed integer quadratic program.

ReLU: naive way makes problem nonlinear

$$y_i = \text{ReLU}(x_i) \Leftrightarrow (y_i \leq x_i - b_i(1-a_i)) \wedge (y_i \leq a_i x_i)$$

$$\wedge (y_i \geq 0) \wedge (a_i \in \mathbb{R}, b_i \in \mathbb{R})$$

overall MILP

$$m^*_t = \min_{\tilde{x}, y^{(l)}, \tilde{x}^{(l)}, a_i^{(l)}} [\tilde{x}^{(l)}]_{ct} - [\tilde{x}^{(l)}]_t$$

subject to:

$$\begin{aligned}\tilde{x}_i^{(l)} - \tilde{x}_i &\leq \varepsilon \quad \forall i \\ \tilde{x}_i - x_i &\leq \varepsilon \quad \forall i\end{aligned}$$

$$y_i^{(l)} = \tilde{x}_i$$

$$\tilde{x}^{(l)} = w_2 y^{(l-1)} + b_2 \quad \forall l=1\dots L$$

$$y_i^{(l)} \leq x_i^{(l)} - \ell_i^{(l)} (1 - a_i^{(l)})$$

$$y_i^{(l)} \geq \tilde{x}_i^{(l)}$$

$$y_i^{(l)} \leq u_i^{(l)} \cdot a_i^{(l)}$$

$$y_i^{(l)} \geq 0$$

$$a_i^{(l)} \in \{0, 1\}$$

on the lower and upper bounds: interval arithmetic

$$\begin{aligned}u_i^{(l)} &= [w_2]_+ u^{(l-1)} - [w_2]_- l^{(l-1)} + b_2 \\ l_i^{(l)} &= [w_2]_+ l^{(l-1)} - [w_2]_- u^{(l-1)} + b_2.\end{aligned}$$

$$[w]_+ = \max\{w, 0\}$$

$$[w]_- = \max\{-w, 0\}$$

Lp constraints
 $u_i^{(l)} = x_i + \varepsilon$
 $l_i^{(l)} = x_i - \varepsilon$

stable and unstable units:

$$\begin{aligned}u_i^{(l)} &\geq \ell_i^{(l)} \geq 0 \\ 0 &\geq u_i^{(l)} \geq \ell_i^{(l)} \\ u_i^{(l)} &\geq 0 \geq \ell_i^{(l)}\end{aligned}$$

stably active \Rightarrow can be removed from optim.

stably inactive \Rightarrow unstable

• tighter bounds has no influence on correctness
• they just lead to more stable units \rightarrow fast optimization.

1. Given an arbitrary binary classifier f for an input domain \mathbb{R}^d and the perturbation set $\mathcal{P}_{\epsilon,p}(x)$ as defined before. Is it possible that every $x \in \mathbb{R}^d$ is "robust", i.e. no adversarial example exists?

i) is it possible that every $x \in \mathbb{R}^d$ is "robust"

$$\mathcal{P}_{\epsilon,p}(x) = \{ \tilde{x} : \| \tilde{x} - x \|_p \leq \epsilon \}$$

yes a constant classifi would not be affected by perturbations

2. Will the fast gradient-sign method (FGSM) always find an adversarial example (assuming there exist some in the set of perturbations $\mathcal{P}_{\epsilon,\infty}(x)$)?

no, it is not guaranteed as FGSM only takes the sign of gradient can overshoot.

Problem 1: Suppose we have a trained binary logistic regression classifier with weight vector $\mathbf{w} \in \mathbb{R}^d$ and bias $b \in \mathbb{R}$. Given a sample $\mathbf{x} \in \mathbb{R}^d$ we want to construct an adversarial example via gradient descent on the binary cross entropy loss:

$$\mathcal{L}(\mathbf{x}, y) = -y \log(\sigma(z)) - (1-y) \log(1 - \sigma(z)),$$

where $\sigma(z) = \frac{1}{1+e^{-z}}$ is the logistic sigmoid function, $z = \mathbf{w}^T \mathbf{x} + b$, and $y \in \{0, 1\}$ is the class label of the sample at hand.

- 1) Derive the gradient $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, y)$. How do you interpret the result?

Hint: You may use the relation $1 - \sigma(z) = \sigma(-z)$.

- 2) Provide a closed-form expression for the worst-case perturbed instance $\tilde{\mathbf{x}}^*$ (measured by the loss \mathcal{L}) for the perturbation set $\mathcal{P}(\mathbf{x}) = \{\tilde{\mathbf{x}} : \|\tilde{\mathbf{x}} - \mathbf{x}\|_2 \leq \epsilon\}$, i.e.

$$\tilde{\mathbf{x}}^* = \arg \max_{\|\tilde{\mathbf{x}} - \mathbf{x}\|_2 \leq \epsilon} \mathcal{L}(\tilde{\mathbf{x}}, y)$$

- c) What is the smallest value of ϵ for which the sample \mathbf{x} is misclassified (assuming it was correctly classified before)?
- d) We would now like to perform adversarial training. Provide a closed-form expression of the worst-case loss

$$\hat{\mathcal{L}}(\mathbf{x}, y) = \max_{\|\tilde{\mathbf{x}} - \mathbf{x}\|_2 \leq \epsilon} \mathcal{L}(\tilde{\mathbf{x}}, y)$$

as a function of \mathbf{x} and \mathbf{w} . How do you interpret the results?

a) $z = \mathbf{w}^T \mathbf{x} + b$

$$\begin{aligned} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, y) &:= \frac{-y}{\sigma(z)} \frac{\partial (\sigma(z))}{\partial z} \cdot \nabla_{\mathbf{x}} z - \frac{(1-y)}{\sigma(-z)} \frac{\partial (\sigma(-z))}{\partial z} \nabla_{\mathbf{x}} z \\ &= \frac{-y}{\sigma(z)} \sigma'(z) \sigma(-z) \mathbf{w} + \frac{(1-y)}{\sigma(-z)} \sigma(-z) \sigma'(z) \mathbf{w} \\ &= -y \sigma(-z) \mathbf{w} + (1-y) \sigma(z) \mathbf{w} \end{aligned}$$

gradient is orthogonal and points in wrong direction

b) $\tilde{\mathbf{x}}^* = \mathbf{x} - \frac{\epsilon \mathbf{w}}{\|\mathbf{w}\|_2}$ if $y=1$

$\tilde{\mathbf{x}}^* = \mathbf{x} + \frac{\epsilon \mathbf{w}}{\|\mathbf{w}\|_2}$ if $y=0$

c) for classification to change classification

$$\sigma(\tilde{x}) > 0.5 \Leftrightarrow w^T \tilde{x} + b = 0$$

$$w^T x - w^T \frac{w}{\|w\|_2} + b = 0$$

$$w^T x - \frac{\|w\|_2}{\|w\|_2} w + b = 0$$

$$\frac{1}{\|w\|_2} (w^T x + b) = \varepsilon$$

for misclassification $\varepsilon > \frac{1}{\|w\|_2} (w^T x + b)$

analogous for $y=0$ $\varepsilon > \frac{1}{\|w\|_2} (-w^T x - b)$

$$d) \hat{f}_c(x, y) = \max_{\|x - \tilde{x}\|_2 \leq \varepsilon} f_c(\tilde{x}, y)$$
$$= f_c(\bar{x}^*, y)$$

$$= -y \log \sigma(w^T x - \frac{\|w\|_2}{\|w\|_2} w + b) + (1-y) \log \sigma(-w^T x - \frac{\|w\|_2}{\|w\|_2} w - b)$$

$$= -y \log \sigma(w^T x - \|w\|_2 \varepsilon + b) + (1-y) \log \sigma(-w^T x - \|w\|_2 \varepsilon - b)$$

Problem 2: In the lecture on exact certification of neural network robustness we have considered $K - 1$ optimization problems (one for each incorrect class) of the form (c.f. slide 42):

$$m_t^* = \min_{\hat{\mathbf{x}}, \mathbf{y}^{(t)}, \hat{\mathbf{x}}^{(t)}, \mathbf{a}^{(t)}} [\hat{\mathbf{x}}^{(L)}]_{c^*} - [\hat{\mathbf{x}}^{(L)}]_t \quad \text{subject to MILP constraints.}$$

That is, for each class $t \neq c^*$, we optimize for the **worst-case margin** m_t^* , and conclude that the classifier is robust if and only if

$$\min_{t \neq c^*} m_t^* \geq 0.$$

However, we can equivalently solve the following single optimization problem:

$$m^* = \min_{\hat{\mathbf{x}}, \mathbf{y}^{(t)}, \hat{\mathbf{x}}^{(t)}, \mathbf{a}^{(t)}} \left([\hat{\mathbf{x}}^{(L)}]_{c^*} - y \right) \quad \text{subject to } y = \max_{t \neq c^*} [\hat{\mathbf{x}}^{(L)}]_t \wedge \text{MILP constraints,}$$

where we have introduced a new variable y into the objective function.

Express the equality constraint

$$y = \max(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{K-1})$$

using only linear and integer constraints. To simplify notation, here $\mathbf{x}_k \in \mathbb{R}$ denotes the logit corresponding to the k -th incorrect class, and l_k and u_k its corresponding lower and upper bound.

Hint: You might want to introduce binary variables to indicate which logit is the maximum.

$u_{\max} = \max_k u_k \quad \text{i.e. largest upper bound}$

$$y \leq x_k + (1 - b_k)(u_{\max} - x_k) \quad \forall 1 \leq k \leq K-1$$

$$y \geq x_i \quad \forall 1 \leq i \leq K-1$$

$$b_k \in \{0, 1\} \quad \forall 1 \leq k \leq K-1$$

$$\sum_{k=1}^K b_k = 1 \rightarrow \text{one element = 1 all others are zero.}$$

value assignment of b is to have $b_k = 1$ for (unique) maximum value $x_k = \max(x_1, \dots, x_{K-1})$

to see this, consider the case $b_K = 1$, but x_K is not max., then $y \leq x_K$, however from ② we would get $x_{\max} \geq x_i \forall 1 \leq i \leq K-1$. contradiction

consider the case $b_K = 1$ and corresponding value x_K is indeed the (unique) maximum. from ① & ② $y = x_K$.

O2 - Convex Relaxation:

Relaxed Classification Margin
Instead of solving:

$$m^* := \min_{\tilde{x}} F(\tilde{x})_{C^+} - F(\tilde{x})_{C^-}$$

$$\text{subject to } \|\tilde{x} - x\|_p \leq \epsilon$$

$$y^{(0)} = \tilde{x}$$

$$\tilde{x}^{(l)} = w_l y^{(l-1)} + b_l \quad \forall l=1..L$$

$$y^{(l)} = \text{ReLU}(\tilde{x}^{(l)}) \quad \forall l=1..L-1$$

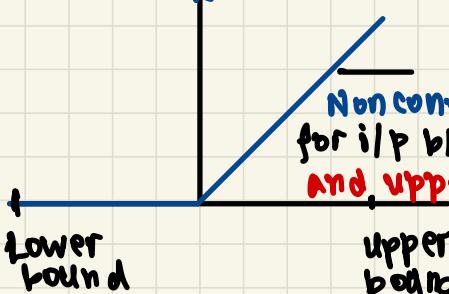
We solve a relaxed problem

- Results in lower bound $m^*(b)$ on the true minimum m^*
 - if $m^* > 0 \rightarrow$ classifier prediction cannot be changed
 - can make optimization possible in polynomial time
 - price → cannot make YES or NO statement in some cases.

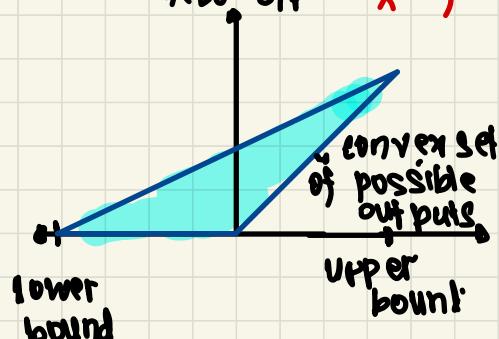
original constraint: $y^{(l)} \in \text{ReLU}(\tilde{x}^{(l)})$

ReLU: Nonconvex opt set

ReLU opt



relaxed constr. $y^{(l)} \in \text{Triangle}$
ReLU opt $\tilde{x}^{(l)}$



Note: Opt of the ReLU activation is no longer deterministic but a variable to optimize over.

Convex Relaxation: formal definition

$$y^{(k)} = \text{Relax}(x_i^{(k)})$$

1) For unstable units ($l_i^{(k)} < 0 \wedge u_i^{(k)} > 0$)

- $y_i^{(k)} \geq 0$
- $y_i^{(k)} \geq \hat{x}_i^{(k)}$ Note: The relaxation is equivalent to relaxing the integer constraint $a_i \in \{0, 1\}$
- $(u_i^{(k)} - l_i^{(k)}) y_i^{(k)} - u_i^{(k)} \hat{x}_i^{(k)} \leq -u_i^{(k)} l_i^{(k)}$ to $a_i \in [0, 1]$

2) For stable active units ($l_i^{(k)} \geq 0$):

- $y_i^{(k)} = \hat{x}_i^{(k)}$

3) For stably inactive units ($u_i^{(k)} \leq 0$):

- $y_i^{(k)} = 0$

Note everything is linear here:

Overall LP:

$$m_t^* = \min_{\tilde{x}, y^{(1)}, \hat{x}^{(k)}} [\hat{x}^{(k)}]_c - [\hat{x}^{(k)}]_t$$

subject to

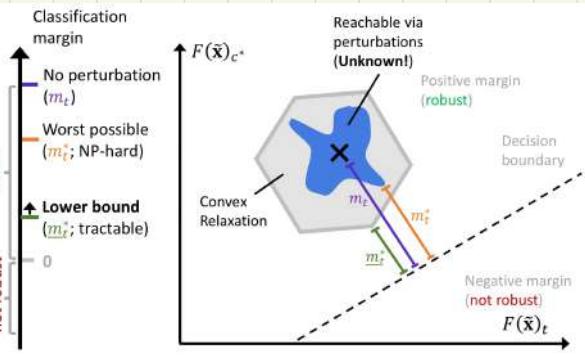
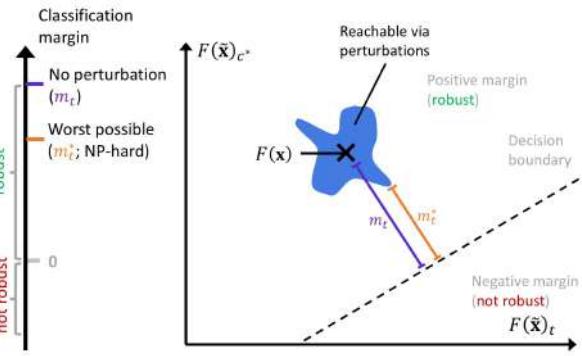
$$\begin{aligned} \tilde{x}_i^{(k)} - \hat{x}_i^{(k)} &\leq \ell \\ \hat{x}_i^{(k)} - x_i &\leq \ell \end{aligned}$$

$$\begin{aligned} \hat{x}_i^{(k)} &= \tilde{x} \\ \hat{x}_i^{(k)} &= w_k y^{(k)} + b_k \quad \forall k=1..L \end{aligned}$$

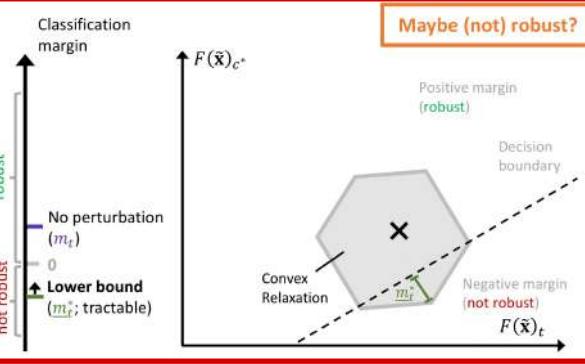
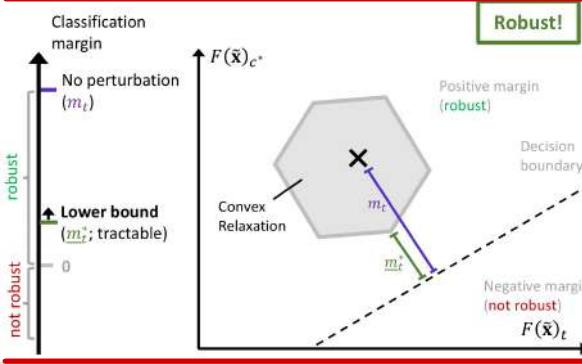
$$\left. \begin{aligned} y_i^{(k)} &\geq \hat{x}_i^{(k)} \\ (u_i^{(k)} - l_i^{(k)}) y_i^{(k)} - u_i^{(k)} \hat{x}_i^{(k)} &\geq 0 \\ u_i^{(k)} - l_i^{(k)} &\leq -u_i^{(k)} l_i^{(k)} \end{aligned} \right\} \begin{aligned} \forall k &= 1..L-1, \\ \forall i &: l_i^{(k)} \leq 0 \wedge u_i^{(k)} > 0 \end{aligned}$$

$$\begin{aligned} y_i^{(k)} &= \hat{x}_i^{(k)} + \ell \quad \forall k=1..L, \forall i : l_i^{(k)} \geq 0 \\ y_i^{(k)} &\leq 0 \quad \forall k=1..L, \forall i : u_i^{(k)} \leq 0 \end{aligned}$$

not robust

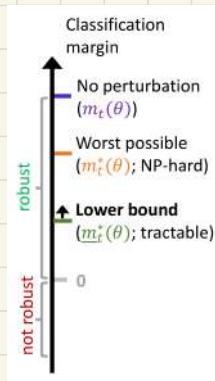


robust



maybe case

Note: tightness of the lower and upper bound influence the quality of relaxation, i.e. how often we return MAYBE.



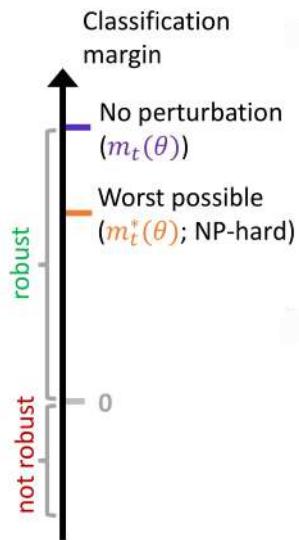
$\Theta \rightarrow$ Weights of NN.

the optim problem's solution depends on this

During robust training we aim to find a good Θ (via gradient descent)

$m_t^*(\Theta)$: dependent on Θ .

challenge is to calculate $\nabla_\Theta m_t^*(\Theta)$ (analytically)



ReLU relaxation via an LP helps in efficiently optimize a robust loss based on the certification

keeping discussion \rightarrow correctly classified simple
 \hookrightarrow misclassified : loss = margin to decision boundary

$$l(f_\theta(\tilde{x}), y) = \max_{t \neq y} \max(F_\theta(\tilde{x})_t - F_\theta(\tilde{x})_y, 0)$$

thus, the supremum evaluates to :

$$\sup_{\tilde{x} \in P(x)} l(f_\theta(\tilde{x}), y) = \max_{t \neq y} (-m_t^*(\theta), 0)$$

To make it tractable, we instead optimize via lower bound
 \hookrightarrow more pessimistic loss

$$\sup_{\tilde{x} \in P(x)} l(f_\theta(\tilde{x}), y) = \max_{t \neq y} (-m_t^*(\theta), 0)$$

$$\leq \max_{t \neq y} (-m_t^*(\theta), 0)$$

challenge: difficult to calculate $T_\theta m_t^*(\theta)$ lower bound

Idea: can we solve another lower bound which does not require us to solve an optimization problem.

Recap: Strong Duality

$$\min_{\mathbf{x}} h_0(\mathbf{x})$$

primal

dual

$$\max_{\alpha} g(\alpha)$$

$$\text{s.t. } h_i(\mathbf{x}) \leq 0 \quad i=1 \dots m$$

In our case, the primal is the LP form slide, overall LP

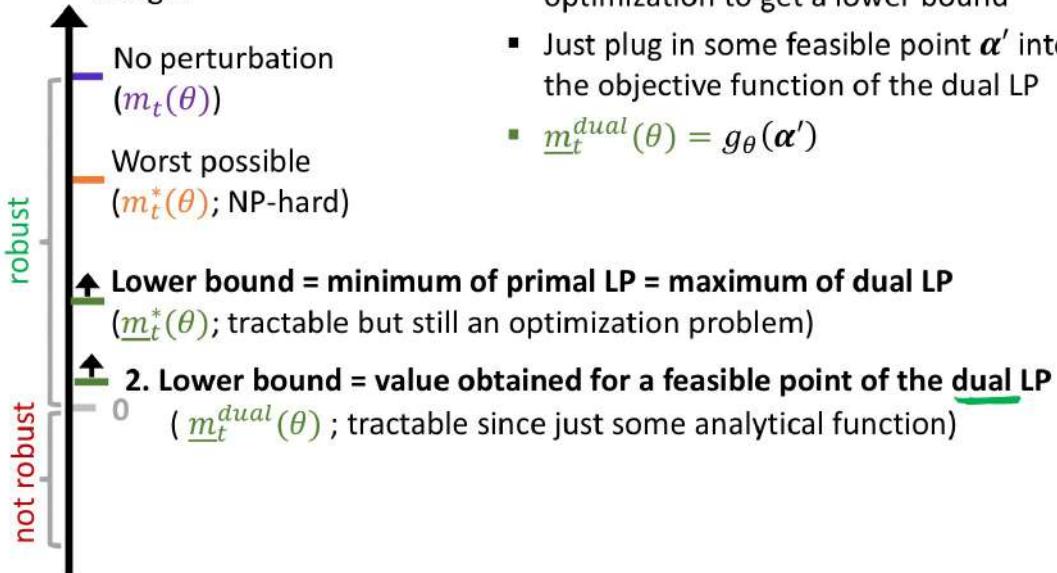
$$\text{s.t. } \alpha_i \geq 0 \quad i=1 \dots M$$

In our case, both primal and dual depend additionally on θ . Thus it would be more accurate to write $g(\alpha)$

$$h_0(\mathbf{x}') \geq h_0(\mathbf{x}^*) = g(\alpha^*) \geq g(\alpha')$$



Classification margin



Classification margin	$\sup_{\tilde{\mathbf{x}} \in \mathcal{P}(\mathbf{x})} \ell(f_\theta(\tilde{\mathbf{x}}), y) = \max_{t \neq y} \max(-m_t^*(\theta), 0)$
No perturbation $(m_t(\theta))$	$\leq \max_{t \neq y} \max(-\underline{m}_t^*(\theta), 0)$
Worst possible $(\underline{m}_t^*(\theta); \text{NP-hard})$	$\leq \max_{t \neq y} \max(-\underline{m}_t^{\text{dual}}(\theta), 0)$
Lower bound = minimum of primal LP = maximum of dual LP $(\underline{m}_t^*(\theta); \text{tractable but still an optimization problem})$	$= \max_{t \neq y} \max(-g_\theta(\alpha'), 0)$
2. Lower bound = value obtained for a feasible point of the dual LP $(\underline{m}_t^{\text{dual}}(\theta); \text{tractable since just some analytical function})$	

We reached our goal: $\nabla_\theta g_\theta(\alpha')$ can be easily computed!

D_B - Lipschitz Continuity for Robustness Certification

a function $f: \mathcal{X} \rightarrow \mathcal{Y}$ is Lipschitz continuous if

$$D_{\mathcal{Y}}(f(x_1), f(x_2)) \leq K \cdot D_{\mathcal{X}}(x_1, x_2) \quad \forall x_1, x_2 \in \mathcal{X}$$

- Idea is to bound how much the output of the network can change for a given perturbation of the input.

- Given a function $f(x) = (\phi_L \circ \phi_{L-1} \cdots \circ \phi_1)(x)$, the Lipschitz constant of f is bounded from above by:

$$L(f) \leq \prod_{i=1}^L L(\phi_i)$$

$$\text{Recall: } F(x) = W_2 f_{L-1} \circ f_{L-2} \cdots \circ f_1(x) + b_L$$

- We can compute the Lipschitz constant of individual layers to upper bound the Lipschitz constant of neural network F .

$$f(x) = Wx + b$$

$$\|Wx_1 + b - (Wx_2 + b)\|_p \leq K \|x_1 - x_2\|_p$$

let $a = x_1 - x_2$; and rearranging

$$L(f) = \sup_{a \neq 0} \frac{\|Wa\|_p}{\|a\|_p} \quad \{ \text{operator norm} \}$$

- $p=1$, $L(f) \rightarrow$ max L1 norm of columns of W
- $p=\infty$, $L(f) \rightarrow$ max L1 norm of rows of W
- $p=2$, $L(f) \rightarrow$ max singular value of W .

[?]
convolution layers:

- computing Lipschitz constant of conv layers is bit trickier in case of $p=2$, but approx efficient.

Activation function

- Typical activation functions (ReLU, softmax,..) have a Lipschitz constant of at worst 1.

ReLU \rightarrow has Lipschitz constant of exactly 1

for feed forward NN F with ReLU activation upper bound is Lipschitz constant

$$L(F) = \prod_{l=1}^L \|W_l\|_F$$

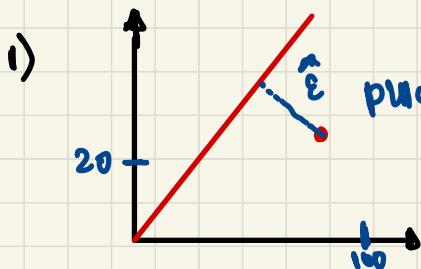
By regularizing norm of weight matrices : we can enforce any Lipschitz constant.

certifying robustness via Lipschitz continuity :

- : Neural net $F \rightarrow$ with Lipschitz constant K .
- given an ℓ_2 sample x and corresponding logits $\hat{y} = F(x)$ we can compute the maximum norm ϵ of a perturbation δ .
 $\delta : \|\delta\|_p \leq \epsilon$ for which we can guarantee robustness of classifier

Approach:

- compute the min. perturbation norm $\hat{\epsilon} = \|\delta\|_p$ in logit space to change classification
- Compute the corresponding perturbation norm $\epsilon = \|\delta\|_p$ based on the Lipschitz constant



$$\text{1) } \|F(x) - F(\tilde{x})\|_p \leq K \|x - \tilde{x}\|_p$$

plugging in the max change we can certify $\hat{\epsilon}$

$$\begin{aligned} \hat{\epsilon} &= \|x - \tilde{x}\|_p \leq \hat{\epsilon}/K \\ \Rightarrow K \|x - \tilde{x}\|_p &\leq \hat{\epsilon} \\ \Rightarrow \|F(x) - F(\tilde{x})\|_p &\leq \hat{\epsilon} \end{aligned}$$

- we can certify that for any $\tilde{x} \in P(x) = \{x : \|x - \tilde{x}\|_p \leq \frac{\epsilon}{\kappa}\}$ the predicted class for sample x does not change
- if κ is small we get more robustness certificates

Partial Considerations

- we want our classifier to have small Lipschitz constant
- However small κ limits the expressiveness of F
- There is a tradeoff b/w expressiveness and robustness

08 - Randomized Smoothing

idea: transform base classifier into a smoothed classifier by randomly adding noise to the input and predicting majority class given many samples.

general idea: for any classifier $f(x) = \arg \max_{c \in \mathcal{C}} f(x)_c$ into a smooth classifier g

$g(x)$: most probable prediction of f under Gaussian noise

Formal definition:

$(g(x))_c$: Probability that f classifies a sample from $N(x, \sigma^2 I)$ (or equivalently $x + \varepsilon$ where $\varepsilon \sim N(0, \sigma^2 I)$) as class c

$$\begin{aligned} P_\varepsilon(f(x + \varepsilon) = c) &= \mathbb{E}_\varepsilon(\mathbb{I}[f(x + \varepsilon) = c]) = \int P(x = z) \mathbb{I}[f(z) = c] dz \\ &= g(x)_c \end{aligned}$$

- output of smooth classifier $g(x)$ is a vector with entries $g(x)_c = P_\varepsilon(f(x + \varepsilon) = c)$
- Now denote with $c^* = \arg \max_c g(x)_c$ the most likely class
 $p_x^* = g(x)_{c^*}$ probability of observing C^* .

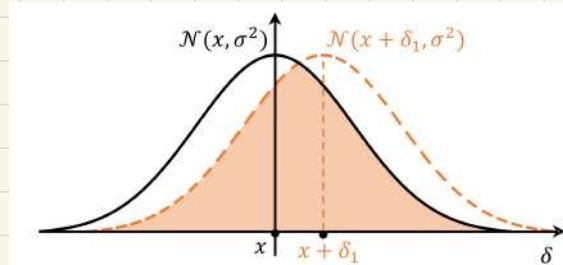
goal: We want to certify that for any admissible perturbation δ it holds

$$\arg \max_c g(x + \delta)_c = c^* \quad \forall \|\delta\|_2 \leq \gamma$$

Perturbation of Input Space - 1D Case

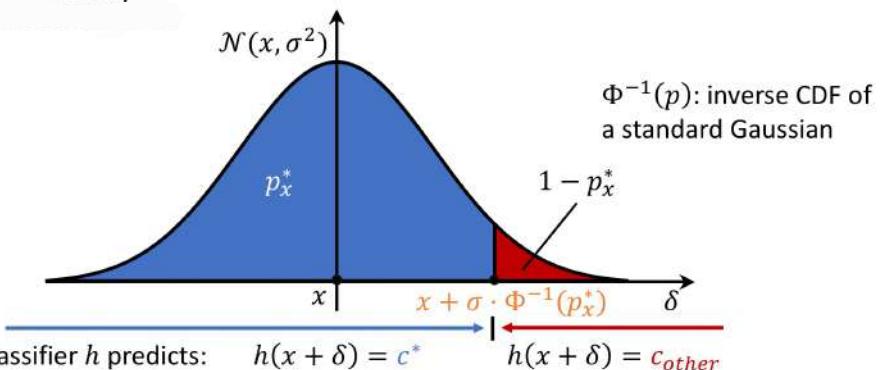
when we slightly perturb x the samples from $\mathcal{N}(x, \sigma^2)$ and $\mathcal{N}(x+\delta), \sigma^2)$ have an overlap $\Rightarrow g(x+\delta)$ and $g(x)$ produce similar output

How large can we make δ and still guarantee that $\text{argmax } g(x+\delta) = \text{argmax } g(x)$



Worst case view on randomized Sampling

- Suppose that with a probability p_x^* the classifier predicts class c^* for a sample from $\mathcal{N}(x, \sigma^2)$.
- In the worst case all the samples from c^* are concentrated on one side, and all samples from other $\neq c^*$ are concentrated on the other side.
- In the worst case: perturbation is orthogonal to the boundary \rightarrow we move straight towards it

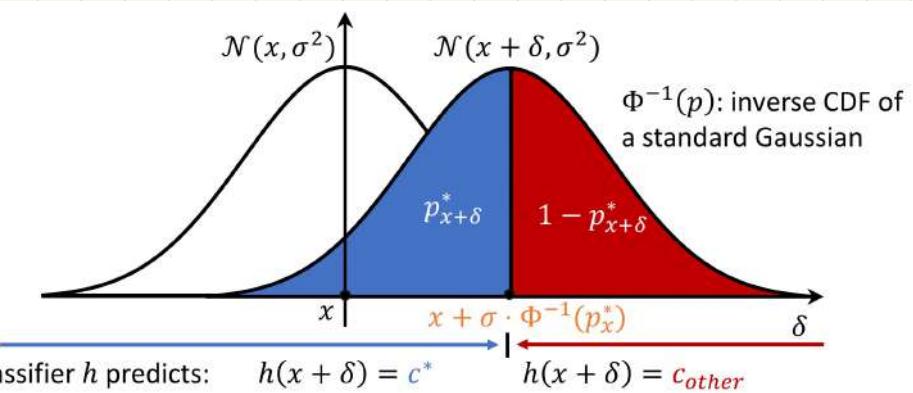


Largest certifiable radius

- If $\|g\|_2 = \sigma^{-1}(p^* x)$: in worst case
- For any perturbation $\|g\|_2 < \sigma^{-1}(p^* x)$,
the smoothed classifier will not change prediction. i.e. $\arg \max g(x+\delta) = \arg \max g(x)$
- Large variance σ^2 could lead to large radii $\|g\|_2$, it could also lead to reproduce $p^* x$ by introducing too much noise

$$p_{x+\delta}^* = 1 - p_{x-\delta}$$

$$= 0.5$$



How to determine $p^* x$?

In neural nets cannot estimate true $p^* x$ to find proportion
 \leftrightarrow solution \rightarrow Use Monte Carlo estimation: sample large number of samples $\overset{\text{from gaussian}}{\sim}$ to estimate p_x as number of samples \rightarrow We get true proportion p_x

- ensure not to overestimate $p^* x$
- $A = I[f(x+\delta) = c^*] \in \{0, 1\}$ for $\delta \sim N(0, \sigma^2)$ denote random variable corresponding to event of observing c^*
- $A \rightarrow$ Bernoulli(Γ) with prob. of observing 1 = $p^* x$
- can compute one-sided $(1-\alpha)$ lower confidence interval for $p^* x$

- If we sample n times and $A=1$ in $\frac{m}{n}$ the PCA $\approx m$
= $\text{Binomial}(n, p^*)$

- We get a lower bound \underline{p}_x^* from confidence interval of the Binomial.

Practical considerations :

For certification, determine \rightarrow most likely class c^*
 \rightarrow Lower bound for \underline{p}_x^*

- We need to take care from a statistical point of view.
- Using a small set of samples we first take a guess at c^*
- Then, using a large set of samples we compute \underline{p}_x^* based on the confidence interval

Training for Randomized Smoothing

data augmentation

instead of training on original data

train on perturbed data

1. Suppose we define a **local** variant of the Lipschitz constant around a given point \mathbf{x}_0 as follows

$$\mathcal{D}_y(f(\mathbf{x}_0), f(\mathbf{x})) \leq k_{\mathbf{x}_0} \cdot \mathcal{D}_{\mathcal{X}}(\mathbf{x}_0, \mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X}$$

How does the local constant $k_{\mathbf{x}_0}$ relate to the global constant k of f ? Which one would provide better guarantees and why?

1) $\mathbf{x}_0 \in \mathcal{X} \rightarrow \text{global const}$
local L-constant provides better guarantees
assuming $P_{\mathcal{X}}(\mathbf{x}_0) \subseteq \mathcal{Y}$ i.e L-const holds & implies -

2. For which class of classifiers does the certificate for the smoothed classifier g equal the certificate for the underlying base classifier f and why?

2)

3. Given two classifiers f_1 and f_2 with Lipschitz constants k_1 and k_2 with $k_1 < k_2$ which one would provide better guarantees? What if we form smoothed classifiers g_1 and g_2 ?

Problem 1: On slide 15 of the robustness chapter, we have defined an optimization problem for untargeted attacks, i.e. we aim to have the sample \hat{x} classified as **any** class other than the correct one:

$$\min_{\hat{x}} \mathcal{D}(\mathbf{x}, \hat{\mathbf{x}}) + \lambda \cdot L(\hat{\mathbf{x}}, y)$$

The loss function is defined as:

$$L(\hat{\mathbf{x}}, y) = \left[Z(\hat{\mathbf{x}})_y - \max_{i \neq y} Z(\hat{\mathbf{x}})_i \right]_+,$$

where $[\mathbf{x}]_+$ is shorthand for $\max(\mathbf{x}, 0)$ and $Z(\mathbf{x})_i = \log f(\mathbf{x})_i$ (i.e. log probability of class i). Here, $L(\hat{\mathbf{x}}, y)$ is positive if $\hat{\mathbf{x}}$ is classified correctly and 0 otherwise.

Provide an alternative loss function to turn this attack into a targeted attack, i.e. we aim to have the sample \mathbf{x} classified as a *specific* target class t .

P1) $\min_{\hat{\mathbf{x}}} D(\mathbf{x}, \hat{\mathbf{x}}) + \lambda L(\hat{\mathbf{x}}, y)$

$$L(\hat{\mathbf{x}}, y) = \left[\max_{i \neq t} z_i(\hat{\mathbf{x}})_i - z_t(\hat{\mathbf{x}})_t \right]$$

Problem 2: Recall from slide 41 the MILP constraints expressing the ReLU activation function:

$$\begin{aligned} y_i &\leq x_i - l_i(1 - a_i), \\ y_i &\leq a_i \cdot u_i, \\ y_i &\geq x_i, \\ y_i &\geq 0, \\ a_i &\in \{0, 1\}, \end{aligned}$$

where $u_i, l_i \in \mathbb{R}$ are upper and lower bounds on the value of the ReLU input x_i .

Show that – for an unstable unit (i.e. $u_i > 0 \wedge l_i < 0$) – a continuous relaxation on a leads to the convex relaxation constraints on slide 54. That is, replacing the constraint $a_i \in \{0, 1\}$ with $a_i \in [0, 1]$ yields

$$(u_i - l_i)y_i - u_i x_i \leq -u_i l_i.$$

$$y_i \leq x_i - j_i(1 - a_i)$$

$$y_i \leq u_i a_i$$

choose a_i s.t. $y_i \leq \min(x_i - j_i(1 - a_i), u_i a_i)$
 it leaves max. $y_i \leq \max_{a_i} \min(x_i - j_i(1 - a_i), u_i a_i)$ negative slope
 free way for optmization. $-ve$ $+ve$ slope

$$y_i \leq a_i u_i \rightarrow \text{plugging}$$

$$\approx y_i(y_i - l_i) - u_i x_i \leq -u_i l_i$$

max at intersection

$$a_i = \frac{x_i - l_i}{u_i - l_i}$$

plugging in orig

~~Problem 3:~~ Convex relaxations of non-linearities are not limited to ReLU. For this exercise, we consider the ReLU6 non-linearity

$$\text{ReLU6}(x) = \min(\max(0, x), 6),$$

which is used in MobileNet models performing low-precision computations on mobile devices.

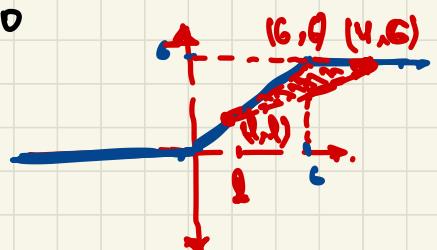
Given input bounds l and u with $l \leq x \leq u$, provide a set of linear constraints corresponding to the convex hull of $\{(x \text{ } \text{ReLU6}(x))^T \mid l \leq x \leq u\}$.

Hint: You have to make a case distinction over different ranges of l and u .

$$\text{ReLU6}(x) = \begin{cases} 0, & x < 0 \\ x, & 0 \leq x < 6 \\ 6, & x \geq 6 \end{cases}$$

If $u_i < 0$ ($l_i < 0$) $\rightarrow y = 0$

If $l_i > 0$, $u_i < 0$ $y = x$



If $l \geq 6$ $y = 6$

$l < 0$ $0 \leq u \leq 6$. (Behaves like unstable)

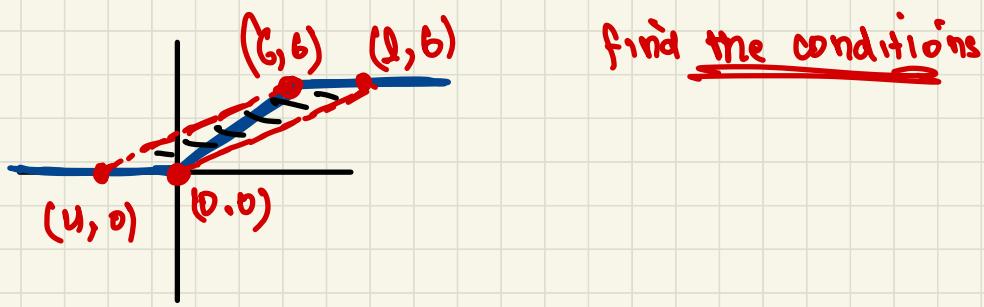
using convex hull to determine $(l, 0)^T, (0, 0)^T$
and $(u, u)^T$

$$y \geq 0 \quad y \geq x \quad y \leq \frac{u}{u-l} (x-l)$$

$0 \leq l < 6$ and $u_i > 6$ convex hull is Δ^2 spanne by

$$(l, l)^T, (6, 6)^T, (u, 6)^T$$

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



Problem 4: Last week, we investigated the adversarial robustness of linear classifiers

$$f(\mathbf{x}) = \mathbb{I}[\mathbf{w}^T \mathbf{x} + b > 0]$$

with weight vector $\mathbf{w} \in \mathbb{R}^d$ and bias $b \in \mathbb{R}$, mapping samples from \mathbb{R}^d to binary labels $\{0, 1\}$.

Given such a linear classifier f , we can define the randomly smoothed classifier $g : \mathbb{R}^d \mapsto \{0, 1\}$ with

$$g(\mathbf{x}) = \operatorname{argmax}_{c \in \{0, 1\}} g_c(\mathbf{x})$$

and

$$g_c(\mathbf{x}) = \Pr_{\boldsymbol{\epsilon}}(f(\mathbf{x} + \boldsymbol{\epsilon}) = c) = \begin{cases} \Pr_{\boldsymbol{\epsilon}}(\mathbf{w}^T(\mathbf{x} + \boldsymbol{\epsilon}) + b \leq 0) & \text{if } c = 0 \\ \Pr_{\boldsymbol{\epsilon}}(\mathbf{w}^T(\mathbf{x} + \boldsymbol{\epsilon}) + b > 0) & \text{else} \end{cases},$$

where $\boldsymbol{\epsilon} \in \mathbb{R}^d$ is a random variable.

For this exercise, we assume that $\boldsymbol{\epsilon}$ follows an isotropic normal distribution, i.e. $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ with elementwise standard deviation $\sigma \in \mathbb{R}_+$.

As discussed in the lecture, evaluating randomly smoothed classifier is typically not tractable and requires sampling. This is however now the case for our simple linear classifier.

Given input $\mathbf{x} \in \mathbb{R}^d$, weights $\mathbf{w} \in \mathbb{R}^d$ and bias $b \in \mathbb{R}$, show that $g_0(\mathbf{x}) = \Phi_{0,1}\left(-\frac{\mathbf{w}^T \mathbf{x}}{\sigma \|\mathbf{w}\|_2} - \frac{b}{\sigma \|\mathbf{w}\|_2}\right)$, where $\Phi_{0,1} : \mathbb{R} \mapsto [0, 1]$ is the cumulative distribution of the standard normal distribution $\mathcal{N}(0, 1)$.

Hint: $\Pr_{\boldsymbol{\epsilon}}(\mathbf{w}^T(\mathbf{x} + \boldsymbol{\epsilon}) + b \leq 0)$ can alternatively be written as:

$$\int_{\mathbb{R}^d} \mathbb{I}[\mathbf{w}^T(\mathbf{x} + \boldsymbol{\epsilon}) + b \leq 0] \mathcal{N}(\boldsymbol{\epsilon} \mid \mathbf{0}, \sigma^2 \mathbf{I}) d\boldsymbol{\epsilon}.$$

P4)

$$g_0(\mathbf{x}) = \Pr_{\boldsymbol{\epsilon}}(\mathbf{w}^T(\mathbf{x} + \boldsymbol{\epsilon}) + b \leq 0)$$

$$= \int_{\mathbb{R}^d} \mathbb{I}[\mathbf{w}^T(\mathbf{x} + \boldsymbol{\epsilon}) + b \leq 0] \mathcal{N}(\boldsymbol{\epsilon} \mid \mathbf{0}, \sigma^2 \mathbf{I}) d\boldsymbol{\epsilon}$$

Integrating over
all possible $\boldsymbol{\epsilon}$

$$\mathbf{z} = \mathbf{w}^T(\mathbf{x} + \boldsymbol{\epsilon}) + b$$

$$\mathbf{y} \sim \mathcal{N}(\mathbf{u}, \Sigma) \quad \text{Affine transform} \quad \mathbf{A}\mathbf{x} + \mathbf{c}$$

$$\mathcal{N}(\mathbf{A}\mathbf{u} + \mathbf{c}, \mathbf{A}\Sigma\mathbf{A}^T)$$

$$\pi \sim N(4, \sigma^2)$$

$$U = w^T \theta + w^T x + b = w^T x + \omega$$

$$\alpha = \sqrt{w^T \sigma^2 I w^T} = \sqrt{w^T w} = \|\omega\|_2$$

$$\begin{matrix} \\ \vdots \\ \vdots \\ \vdots \end{matrix}$$

NTF

Chapter 6 : Autoregressive Models

$$x_t = c + \sum_{i=1}^p \phi_i x_{t-i} + \varepsilon_t$$

autoregressive model of order 'p'

ϕ_1, ϕ_2, \dots → parameters

c → constant

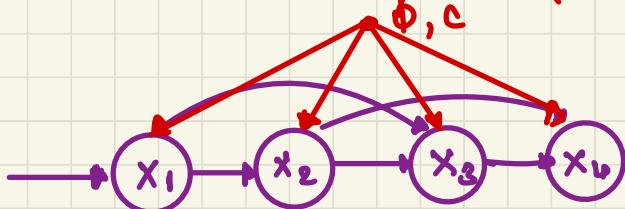
$\varepsilon_t \sim N(0, \sigma^2)$ is white noise

x_{t-i} → lagged value at time $t-i$

- perform regression on $([x_{t-1}, \dots, x_{t-p}])$ are the inputs and are x_t the outputs
- Modifications on x_t have repercussions, $(x_t)_t$ are not independent.

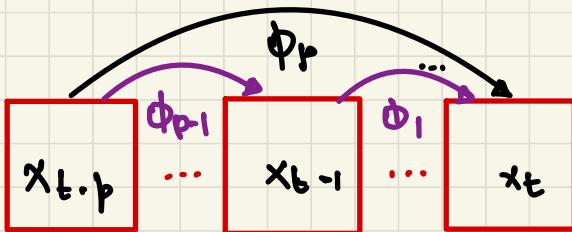
Rewrite AR model :

$$P(x_t | x_{t-1}, \dots, x_{t-p}) \sim N(c + \sum_{i=1}^p \phi_i x_{t-i}, \sigma^2)$$



- mean function : $\mu(t) = \mathbb{E}[x_t]$, depends on t
- autocovariance $\gamma(t, i) = \text{cov}(x_t, x_{t-i})$, by default depends on t and i .
- Pearson's autocorrelation function

$$r(t, i) = \frac{\text{cov}(x_t, x_{t-i})}{\sqrt{\text{var}(x_t)} \sqrt{\text{var}(x_{t-i})}}$$



Stationary Process:

- (i) $\mathbb{E}[x_t] = \mathbb{E}[x_{t-i}] \xrightarrow{\text{mean is constant}} = \mu, \forall t, \forall i$
 - (ii) $\text{cov}(x_t, x_{t-i}) = \gamma_i, \forall t, \forall i \xrightarrow{\text{autocovariance only depends on time lag.}}$
 - (iii) $\mathbb{E}[|x_t|^2] < \infty, \forall t$
- $\xrightarrow{\text{possible to estimate mean and autocovariance by averaging measure over time}}$

• Moments of stationary process AR(p) :

$$- E[X_t] = \mu = \frac{c}{1 - \sum_{i=1}^p \phi_i}, \quad \forall t$$

$$- \text{Var}(X_t) = \gamma_0 = \sum_{j=1}^p \phi_j \gamma_{t-j} + \sigma^2, \quad \forall t$$

$$- \text{Cov}(X_t, X_{t-i}) = \gamma_i = \sum_{j=1}^p \phi_j \gamma_{t-j}, \quad \forall t, X_t, X_i \in \mathbb{R}$$

$$- \rho_i = \frac{\gamma_i}{\gamma_0}$$

conditions for process to be stationary

AR(1)p is stationary iff roots of characteristic polynomial :

4 $\Phi(L) = 1 - \sum_{i=1}^p \phi_i L^i$

lie outside the circle

Example:

- AR(1) : $X_t = c + \phi_1 X_{t-1} + \varepsilon_t$ is stationary
if $|\phi_1| < 1$

- $|\phi_1| < 1$ from quadratic equations

- AR(2) : $X_t = c + \phi_1 X_{t-1} + \phi_2 X_{t-2} + \varepsilon_t$
 $\phi_1 + \phi_2 < 1, \quad \phi_2 - \phi_1 < 1, \quad |\phi_2| < 1$

Parameter Learning

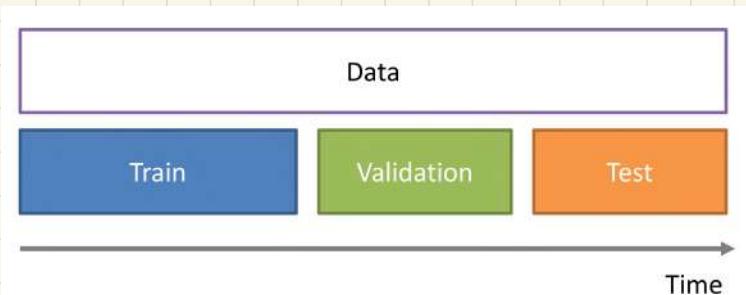
$$\begin{bmatrix} \Phi_1 \\ \Phi_2 \\ \vdots \\ \Phi_p \end{bmatrix} := (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}$$

and $\mathbf{y} = \begin{bmatrix} x_{p+1} & \dots & x_0 \\ x_p & \ddots & x_1 \\ \vdots & \ddots & \vdots \\ x_p & & \vdots \\ x_{p+1} & & \vdots \end{bmatrix}$

Parameter learning from Yule-Walker equations

$$\begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_{p-1} \\ \gamma_p \end{bmatrix} := \begin{bmatrix} \gamma_0 & \gamma_{-1} & \dots & \gamma_{2-p} & \gamma_{1-p} \\ \gamma_1 & \ddots & & & \gamma_{2-p} \\ \vdots & & \ddots & & \vdots \\ \gamma_{p-2} & & & \gamma_{-1} & \gamma_0 \\ \gamma_{p-1} & \gamma_{p-1} & \gamma_{p-2} & \dots & \gamma_1 \end{bmatrix} \begin{bmatrix} \Phi_1 \\ \Phi_2 \\ \vdots \\ \Phi_{p-1} \\ \Phi_p \end{bmatrix}$$

Note when splitting the data always maintain the temporal order.



Questions

1. What is the mean $E[X_t]$ of the following processes:

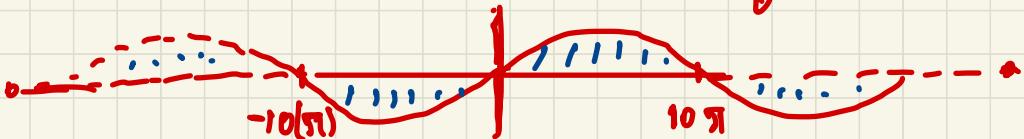
a) $X_t = \sin(t/10) + \varepsilon_t$ where $\varepsilon_t \sim N(0, \sigma)$

b) $X_t = 4 - 0.8 * X_{t-1} - 0.1 * X_{t-2} + \varepsilon_t$ where $\varepsilon_t \sim N(0, \sigma)$

c) $X_t = 4 + X_{t-1} + \varepsilon_t$ where $\varepsilon_t \sim N(0, \sigma)$ and $X_0 \sim N(0, \sigma)$

i) a) $X_t = \underbrace{\sin(t/10)}_{\text{not an autoregressive model.}} + \varepsilon_t \quad \varepsilon_t \sim N(0, \sigma)$

$$E[X_t] = E[\sin(t/10)] + E(\varepsilon_t)$$



We can expect the $\sin(t/10)$ to evenly cover

$$\therefore E\left(\sin\left(\frac{t}{10}\right)\right) = 0$$

b) $X_t = 4 - 0.8 * X_{t-1} - 0.1 * X_{t-2} + \varepsilon_t$

$$\phi_1 = -0.8 \quad \phi_2 = -0.1$$

$$\begin{array}{l} \textcircled{1} \quad \phi_1 + \phi_2 < 1 \\ \textcircled{2} \quad \phi_2 - \phi_1 < 1 \\ \textcircled{3} \quad |\phi_2| = 0.1 < 1 \end{array} \quad \left. \begin{array}{l} -0.9 < 1 \\ 0.7 < 1 \\ 0 \end{array} \right\} \text{satisfied}$$

it autoregressive and stationary process

$$E[X_t] = \frac{c}{1 - \sum \phi_i} = \frac{4}{1 + 0.8 + 0.1} = \frac{4}{1.9} = \boxed{\frac{40}{19}}$$

$$c) X_t = c + X_{t-1} + \epsilon_t$$

This is an autoregressive model.

$\phi_1 = 1$ $|\phi_1| < 1$ for process to be stationary

\therefore The value keeps on increasing unboundedly, so $E[X_t]$ does not exist.

2. Does Yule Walker parameter learning assume a stationary process? Why?

2) Yes, as we need to ensure covariance only depends upon the time lag and not on X_t .

Problem 1: Consider the stationary AR(p) process $X_t = c + \sum_{i=1}^p \phi_i X_{t-i} + \epsilon$ with $\epsilon \sim \mathcal{N}(0, \sigma^2)$. We denote by μ the mean $E[X_t]$ and by γ_i the autocovariance $Cov(X_t, X_{t-i})$. Show:

$$1. \mu = \frac{c}{1 - \sum_{i=1}^p \phi_i}, \text{ for all } t$$

$$2. \gamma_0 = \sum_{j=1}^p \phi_j \gamma_{-j} + \sigma^2$$

$$3. \gamma_i = \sum_{j=1}^p \phi_j \gamma_{i-j}, \text{ for all } t, i \in [1, p]$$

$$0 \quad X_t : c + \sum_{i=1}^p \phi_i X_{t-i} + \epsilon$$

taking expectation both sides

$$\mathbb{E}[X_t] = \underbrace{\mathbb{E}[\epsilon]}_{\text{as this is stationary process}} + \mathbb{E}\left[\sum_{i=1}^p \phi_i X_{t-i}\right] + \underbrace{\mathbb{E}[c]}_{\mu}$$

$$\mu = c + \sum_{i=1}^p \phi_i \underbrace{\mathbb{E}[X_{t-i}]}_{\mu} \rightarrow \mu = c + \sum_{i=1}^p \phi_i \mu$$

$$\mu (1 - \sum_{i=1}^p \phi_i) = c \rightarrow \mu = \frac{c}{1 - \sum_{i=1}^p \phi_i}$$

$$\begin{aligned}
 \textcircled{2} \quad \text{cov}(x_t, x_b) &= \underbrace{\text{cov}(c, x_t)}_0 + \text{cov}\left(\sum_{i=1}^p \phi_i x_{t-i}, x_b\right) \\
 &\quad + \text{cov}\left(\underbrace{\varepsilon_t}_\sigma, x_b\right) \\
 &= \sum_{i=1}^p \phi_i \text{cov}(x_{t-i}, x_b) + \sigma^2 \\
 &= \sum_{i=1}^p \phi_i \gamma_i + \sigma^2
 \end{aligned}$$

take the cov inside

$$\begin{aligned}
 \textcircled{3} \quad \text{cov}(x_t, x_{t-j}) &= \underbrace{\text{cov}(c, x_{t-j})}_0 + \sum_{i=1}^p \phi_i \text{cov}(x_{t-i}, x_{t-j}) \\
 &\quad + \text{cov}(\varepsilon_t, x_{t-j}) \\
 &= \sum_{i=1}^p \phi_i \gamma_{j-i} + \sigma^2
 \end{aligned}$$

Problem 2:

- a) Consider the following AR(1) process $X_t = c + \phi_1 X_{t-1} + \epsilon$ with $\epsilon \sim \mathcal{N}(0, \sigma^2)$. Show that this process is stationary iff the following condition is fulfilled

$$|\phi_1| < 1.$$

- b) Consider the following AR processes:

- $X_t = c + .8 \times X_{t-1} + .1 \times X_{t-2} + \epsilon$ with $\epsilon \sim \mathcal{N}(0, \sigma^2)$
- $X_t = -\sum_{k=1}^p {}_k^p X_{t-k} + \epsilon$ with $\epsilon \sim \mathcal{N}(0, \sigma^2)$

Are these processes stationary?

$$\begin{aligned}
 2) \quad \text{AR(1)} : \quad x_t &= c + \phi_1 x_{t-1} + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, \sigma^2) \\
 \text{for a process to be stationary (AR(1))} \\
 \Phi(L) = 1 - \underline{\phi_1 L} \rightarrow \text{roots should lie outside}
 \end{aligned}$$

unit circle

$$1 - \Phi L = 0 \Rightarrow$$

$$\boxed{L = \frac{1}{\Phi}}$$

$$\frac{1}{\Phi} > 1$$



$$\boxed{|\Phi| < 1}$$

final condition

$$\& \frac{1}{\Phi} < -1$$

b) (i) $x_t = c + 0.8x_{t-1} + 0.1x_{t-2} + \epsilon$

$$\Phi_1 = 0.8 \quad \Phi_2 = 0.1$$

① $\Phi_1 + \Phi_2 < 1 \rightarrow 0.9 < 1$

② $\Phi_2 - \Phi_1 < 1 \rightarrow -0.7 < 1$

③ $|\Phi_2| < 1 \rightarrow 0.1 < 1$

} all conditions
are met, hence
it is a
stationary
process.

(ii) $x_t = - \sum_{k=1}^p \binom{p}{k} x_{t-k} + \epsilon$

$$p(x) = 1 + \sum_{k=1}^p \binom{p}{k} x^k \rightarrow \text{roots should lie outside unit circle}$$

$$= (1+x)^p$$

$x = \underline{-1}$ (not outside
unit circle, \therefore not
stationarity)

Chapter 6: Markov Chains

Markov Property: $P(x_t | x_1, \dots, x_{t-1}) = P(x_t | x_{t-1})$

$t \in \{1, \dots, T\}$ (Discrete values)

$x_t \in \{1, 2, \dots, K\}$

joint distribution

$$P(x_1 = i_1, \dots, x_T = i_T) = P(x_1 = i_1) \prod_{t=1}^{T-1} P(x_{t+1} = i_{t+1} | x_t = i_t)$$

general case of Markov Chain:

$$P(x_1 = i) = \pi_i \quad P(x_{t+1} = j | x_t = i) = A_{ij}^{(t+1)}$$

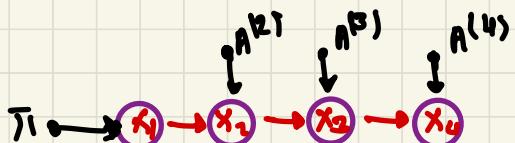
$\pi \in \mathbb{R}^K$ → prior probability

$A^{(t)} \in \mathbb{R}^{K \times K}$ → transition matrix.

consequently joint probability is

$$P(x_1 = i_1, \dots, x_T = i_T) = \pi_{i_1} \times A_{i_1, i_2}^{(1)} \times \dots \times A_{i_{T-1}, i_T}^{(T)}$$

If Parameters = $(K-1) + (T-1)K(K-1)$



Markov Chain - Stationary case

- time homogeneous and stationary Markov Chain

$$P(X_t = i) = \pi_i; \quad P(X_{t+1} = j | X_t = i) = A_{ij}$$

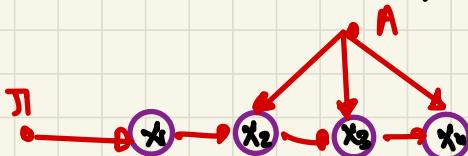
* Transition Matrix $A^{(t)} = A$

↳ INDEPENDENT OF TIME

- joint probability becomes

$$P(X_1 = i_1, \dots, X_T = i_T) = \pi_{i_1} \times A_{i_1, i_2} \times \dots \times A_{i_{T-1}, i_T}$$

$$\# \text{Parameters} = (k-1) + k(k-1)$$



Markov Chains - Learning of Model Parameters

$$P(\text{all}) = \prod_{n=1}^N \Pr(X_n^{(n)}) \prod_{t=1}^{T_{n-1}} \Pr(X_{t+1}^{(n)} | X_t^{(n)}) \cdot \left(\prod_{k=1}^K \pi_k^{L(k)} \right) \left(\sum_{i=1}^k \sum_{j=1}^k N(i,j) A_{ij} \right)$$

$$\log(P(\text{all})) = \sum_{k=1}^K L(k) \log(\pi_k) + \sum_{i=1}^N \sum_{j=1}^N N(i,j) \log(A_{ij})$$

minimizing Φ (all) s.t. $\sum_k \pi_k = 1$ and $\sum_j A_{ij} = 1$

$$A_{ij} = \frac{N(i,j)}{\sum_{j'} N(i,j')}$$

$$\pi_k = \frac{L(k)}{\sum_{k'} L(k')}$$

More insights:

TASK 1 $A_{ij}(n) = \Pr(x_{t+n} = j \mid x_t = i)$

Probability of going from state i to j in exactly n steps

$$P(x_{t+n} = j \mid x_t = i) = \sum_{k=1}^K A_{kj} A_{ik}(n-1)$$

$$A(n) = A(n-1) A \xrightarrow{A \cup \underset{n}{\overbrace{A}} \cdots \cup A} A(n) = A^n$$

Chapman-Kolmogorov Equation

$$A_{ij}(m+n) = \sum_{k=1}^K A_{ik}(m) A_{kj}(n) = A(m+n) = A(m) A(n)$$

TASK 2 $\pi_j(t) = \Pr(x_t = j)$

Probability of reaching state j in t timesteps

$$\Pr(X_t = j) = \sum_{i=1}^K \Pr(X_t = j | X_{t-1} = i) \Pr(X_{t-1} = i)$$

$$= \sum_{i=1}^K A_{ij} \pi_i(t-1)$$

$$\Rightarrow \pi(t) = \pi(t-1) A$$

$$\pi(t) = \pi A^{t-1}$$

$\pi(t)$ and π are row vectors.

Questions

- We assume that $X_t \in \{1, 2, 3\}$. We consider $\pi = \begin{bmatrix} 0.0 \\ 0.5 \\ 0.5 \end{bmatrix}$ and $A = \begin{bmatrix} 0.6 & 0.2 & 0.2 \\ 0.1 & 0.5 & 0.4 \\ 0.4 & 0.1 & 0.5 \end{bmatrix}$.
 - What is the probability to observe the sequence $X^{(1)} = [1, 2, 3]$?
 - What is the probability to observe the sequence $X^{(2)} = [2, 2, 3]$?

$$a) x^{(1)} = [1 \ 2 \ 3].$$

$$= \underbrace{\pi_1}_{0} (1) * A_{12} \ A_{23} = 0.$$

$$b) x^{(2)} : [2, 2, 3]$$

$$= \pi_{12} * A_{22} * A_{23} = 0.5 * 0.5 * 0.4 = \underline{\underline{0.1}}$$

2. We assume that $X_t \in \{1, 2, 3\}$ and we observed three sequences:

- $X^{(1)} = [1, 3, 2]$
- $X^{(2)} = [3]$
- $X^{(3)} = [1, 1, 3, 2]$

What is the MLE of the transition matrix $A \in \mathbb{R}^{3 \times 3}$?

2)

transition occurrences		transition
transition	occurrences	
$1 \rightarrow 1$	# 1	$2 \rightarrow 1$
$1 \rightarrow 2$	# 0	$2 \rightarrow 2$
$1 \rightarrow 3$	# 2	$2 \rightarrow 3$
$2 \rightarrow 1$	0	
$3 \rightarrow 2$	2	
$3 \rightarrow 3$	0	

$$\pi(1) = * 2$$

$$\pi(2) = 0$$

$$\pi(3) = 1.$$

$$\pi = \begin{bmatrix} 2/3 \\ 0 \\ 1/3 \end{bmatrix}$$

$$A = \begin{bmatrix} 1/2 & 0 & 2/3 \\ 1/3 & 1/3 & 1/3 \\ 0 & 1 & 0 \end{bmatrix}$$

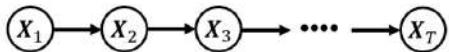
There is no data available for state 2

so, we make a reasonable assumption that transitions are equally likely.

Problem 3: Let \mathbf{X}_t be a 2-D random vector:

$$\mathbf{X}_t = \begin{bmatrix} u_t \\ v_t \end{bmatrix}, \quad \text{where } u_t, v_t \in \{1, 2, \dots, K\}. \quad (1)$$

Consider the following Markov chain.



Model parameters are as follows:

- initial distribution $\pi_x \in \mathbb{R}^{K \times K}$ that parametrizes $Pr(\mathbf{X}_1)$:

$$\bullet \quad Pr(\mathbf{X}_1 = \begin{bmatrix} i \\ j \end{bmatrix}) = \pi_x(i, j). \quad (2)$$

- transition probability matrix $\mathbf{A}_x \in \mathbb{R}^{K \times K \times K \times K}$ that parametrizes $Pr(\mathbf{X}_{t+1} | \mathbf{X}_t)$:

$$Pr(\mathbf{X}_{t+1} = \begin{bmatrix} i_{t+1} \\ j_{t+1} \end{bmatrix} | \mathbf{X}_t = \begin{bmatrix} i_t \\ j_t \end{bmatrix}) = \mathbf{A}_x(i_t, j_t, i_{t+1}, j_{t+1}). \quad (3)$$

The joint probability can be factorized as:

$$Pr(\mathbf{X}_1, \dots, \mathbf{X}_T) = Pr(\mathbf{X}_1) \prod_{t=1}^{T-1} Pr(\mathbf{X}_{t+1} | \mathbf{X}_t).$$

In this task, we refer to this model as "2-D first-order Markov chain".

- a) Does the sequence $[u_1, \dots, u_T]$ (where $u_t \in \{1, 2, \dots, K\}$ is defined in Eq. 1) have the first-order Markov property? Why or why not?

a) u_t depends on both u_{t-1} and v_{t-1}
 and v_{t-1} depends on u_{t-2} & v_{t-2} .
 \therefore the first order Markov property would not hold

- b) Let $[Y_1, \dots, Y_T]$ be a 1-D first-order Markov chain with the following initial and transition probabilities (Y_1, \dots, Y_T are binary-valued).

$$\pi_y = [0.5 \ 0.5], \quad \mathbf{A}_y = \begin{bmatrix} 0.2 & 0.8 \\ 0.5 & 0.5 \end{bmatrix}.$$

- Briefly explain why the sequence $\begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix}, \begin{bmatrix} Y_2 \\ Y_3 \end{bmatrix}, \dots, \begin{bmatrix} Y_{T-1} \\ Y_T \end{bmatrix}$ is a 2-D first-order Markov chain.
- Compute initial and transition probabilities, π_x and \mathbf{A}_x (defined in Eqs. 2 and 3) for the sequence $\begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix}, \begin{bmatrix} Y_2 \\ Y_3 \end{bmatrix}, \dots, \begin{bmatrix} Y_{T-1} \\ Y_T \end{bmatrix}$.

b) (i) $\left\{ \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix}, \begin{bmatrix} Y_2 \\ Y_3 \end{bmatrix}, \dots, \begin{bmatrix} Y_{T-1} \\ Y_T \end{bmatrix} \right\} \rightarrow \text{2D Markov chain}$

$$\Pr \left(\begin{pmatrix} Y_{t-1} \\ Y_t \end{pmatrix} \mid Y_{t-1}, Y_{t-2}, \dots, Y_1 \right) \longrightarrow \textcircled{1}$$

as $\{Y_1, Y_2, \dots, Y_T\}$ is 1-D Markov

$$\Pr(Y_t \mid Y_{t-1}, \dots, Y_1) = \Pr(Y_t \mid Y_{t-1})$$

using the results in $\textcircled{1}$

$$\Pr \left(\begin{pmatrix} Y_{t-1} \\ Y_t \end{pmatrix} \mid Y_{t-1}, \dots, Y_1 \right) = \Pr \left(\begin{pmatrix} Y_{t-1} \\ Y_t \end{pmatrix} \mid \begin{pmatrix} Y_{t-2} \\ Y_{t-1} \end{pmatrix} \right)$$

$$(ii) \pi_x(i, j) = \Pr \left(\begin{bmatrix} Y_2 \\ Y_1 \end{bmatrix} = \begin{bmatrix} j \\ i \end{bmatrix} \right) = \underline{\mathbf{A}_{ij}(j, i)} \underline{\pi_y(i)}$$

$$A_{ij}(i, j, i', j') = \Pr \left(\begin{bmatrix} Y_t \\ Y_{t-1} \end{bmatrix} = \begin{bmatrix} i' \\ j' \end{bmatrix} \mid \begin{bmatrix} Y_{t-2} \\ Y_{t-1} \end{bmatrix} = \begin{bmatrix} i \\ j \end{bmatrix} \right)$$

$$= \begin{cases} 0 & j \neq i' \\ A_{ij}(i, j) & \text{otherwise} \end{cases}$$

Chapter 7 : Hidden Markov Models :

Hidden Markov Model (HMM) is composed of a sequence of hidden / latent variables $[z_1, \dots z_T]$ and a sequence of observed variables $[x_1, \dots x_T]$

- $z_1, z_2, \dots z_T$ satisfy the Markov property

$$P(z_{t+1} | z_t, \dots z_1) = P(z_{t+1} | z_t)$$

transition probability

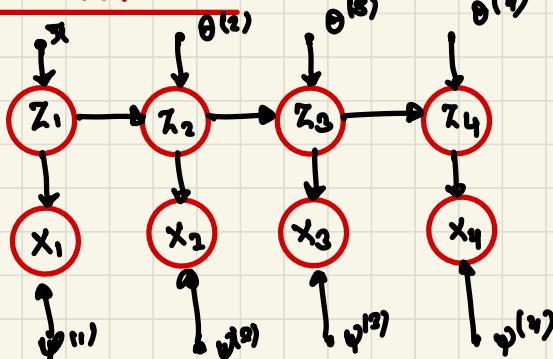
- Distribution of x_t only depends on z_t

$$P(x_{t+1} | z_1, \dots z_T, x_1, \dots x_T) = P(x_{t+1} | z_{t+1})$$

emission probability

- discrete time $t \in \{1, 2, \dots T\}$ and discrete r.v $z \in \{1, 2, \dots k\}$
- observed data can be discrete or continuous

HMM : general case



• joint distribution can be written as:

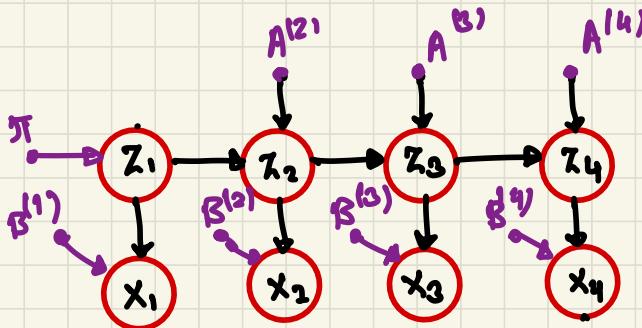
$$P(Z_1 = z_1, \dots, Z_T = z_T, X_1 = x_1, \dots, X_T = x_T)$$

$$= P(Z_1 = z_1 ; \pi) \prod_{t=1}^T P(Z_{t+1} = z_{t+1} | Z_t = z_t; \Theta^{t+1}) \\ * \prod_{t=1}^T P(X_t = x_t | Z_t = z_t; \phi^{t+1})$$

HMM : Discrete Case :

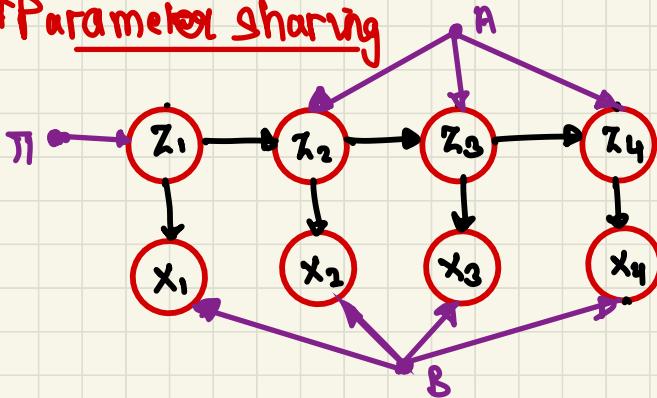
$$X_t \in \{1, 2, \dots, K\}$$

$$P(Z_1 = i) = \pi_i \\ P(Z_{t+1} = j | Z_t = i) = A_{ij}^{(t+1)} \\ P(X_{t+1} = j | Z_{t+1} = i) = B_{ij}^{(t+1)}$$



$$\# \text{ Parameters} \\ = K + (T-1)K^2 + T \times K$$

Parameter sharing



$$\# \text{ Parameters} \\ = K + K^2 + K \times K'$$

From now on assume parameter sharing in the Markov chains

$$P(z_1 = z_1, \dots, z_T = z_T, x_1 = x_1, \dots, x_T = x_T) = P(z_1 = z_1; \pi)^* \prod_{t=1}^{T-1} A_{z_t z_{t+1}} \prod_{t=1}^T B_{z_t x_t}$$

Tasks concerning HMM :

• Inference :

- model parameters are fixed
- info from posterior $P(z_{1:T} | x_{1:T})$

• Parameter Learning :

- we seek to learn the model parameters
- $x_{1:T}$ is observed
- $z_{1:T}$ is (usually) not observed

Inference:

- Filtering : $P(z_t | x_{1:t}) \rightarrow$ incrementally as data streams in
- Smoothing : $P(z_t | x_{1:T})$ offline past and future
 - mode of posterior, viterbi decoding
- MAP Inference : computes $\arg \max_{x_{1:T}} P(z_{1:T} | x_{1:T})$

The Forward Algorithms :

- Goal : incrementally compute $P(z_t | x_1:t)$

Bayes rule gives

$$P(z_t = k | x_1:t) = \frac{P(z_t = k, x_1:t)}{\sum_{j=1}^K P(z_t = j, x_1:t)}$$

- for convenience, we denote

$$\alpha_t(k) \stackrel{\text{def}}{=} P(z_t = k, x_1:t) \text{ and } \alpha_t = \begin{bmatrix} \alpha_t(1) \\ \vdots \\ \alpha_t(K) \end{bmatrix}$$

Hence, we have

$$P(z_t = k | x_1:t) = \frac{\alpha_t(k)}{\sum \alpha_t} = \frac{\alpha_t(k)}{\sum \alpha_t(k)}$$

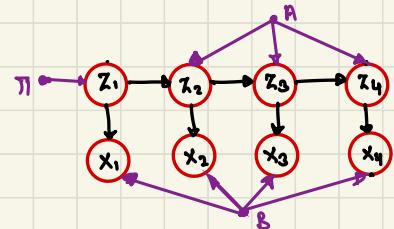
Forward algo computes sequentially

1) computes α_1 (initialization)

2) Given α_t , compute α_{t+1} (recursion)

initialization : computation of parameter α_1

$$\begin{aligned} \alpha_1(k) &= P(z_1 = k, x_1) \\ &= P(z_1 = k) P(x_1 | z_1 = k) \\ &= \pi_k B_{kx_1} \end{aligned}$$



Recursion : given α_t , we can compute α_{t+1}

$$\alpha_{t+1}(k) = P(\pi_{t+1} = k, x_{1:t+1})$$

$$\alpha_{t+1}(k) = \boxed{B_{K|x_{t+1}} \sum_{j=1}^K A_{jk} \alpha_t(j)}$$

$$\alpha_{t+1} = B_{:(x_{t+1})} \odot A^T \alpha_t$$

$O(TK^2)$: linear in T

$$\begin{bmatrix} \alpha_{t+1}(1) \\ \alpha_{t+1}(2) \\ \vdots \\ \alpha_{t+1}(K) \end{bmatrix} = \begin{bmatrix} B(1, x_{t+1}) \\ B(2, x_{t+1}) \\ \vdots \\ B(K, x_{t+1}) \end{bmatrix} \odot \begin{bmatrix} A^T \alpha_t \end{bmatrix} \begin{bmatrix} x_{t+1} \\ \alpha_t(1) \\ \alpha_t(2) \\ \vdots \\ \alpha_t(K) \end{bmatrix}$$

The Forward Backward Algorithm:

incrementally compute $P(\pi_t | x_{1:T})$ full data in time

Bayes Rule :

$$P(\pi_t = k | x_{1:T}) = \frac{P(\pi_t = k, x_{1:T})}{\sum_{j=1}^K P(\pi_t = j, x_{1:T})}$$

for convinience

$$\beta_t(k) = P(x_{t+1:T} | z_t=k) \text{ and } \beta_t = \begin{bmatrix} \beta_{k_1}(1) \\ \vdots \\ \beta_{k_n}(n) \end{bmatrix}$$

using $\alpha_t(k)$ and $\beta_t(k)$

$$P(z_t=k | x_{1:T}) \propto \alpha_t(k) \beta_t(k)$$

Backward algorithm computes recursively

- 1) computes β_T (initialization)
- 2) given $\beta_{t+1} \rightarrow$ compute β_t

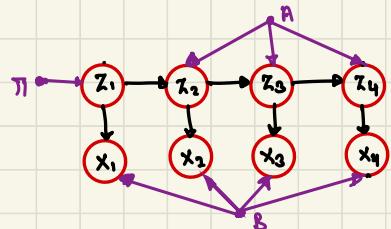
Initialization

$$\beta_T(k) = 1$$

for $t=T$ term is completely captured.

Recursion: given β_{t+1} we compute β_t

$$\begin{aligned}\beta_t(j) &= P(x_{t+1:T} | z_t=j) \\ &= \sum_{k=1}^n A_{jk} B_{kx_{t+1}} \beta_{t+1}(k)\end{aligned}$$



$$\beta_t = \underline{A} (\underline{B}_{:(x_{t+1})} \odot \beta_{t+1})$$

$$\begin{bmatrix} \beta_t(1) \\ \beta_t(2) \\ \vdots \\ \beta_t(K) \end{bmatrix} = \boxed{A} \left(\begin{bmatrix} B(1, x_{t+1}) \\ B(2, x_{t+1}) \\ \vdots \\ B(K, x_{t+1}) \end{bmatrix} \right) \odot \begin{bmatrix} \beta_{t+1}(1) \\ \beta_{t+1}(2) \\ \vdots \\ \beta_{t+1}(K) \end{bmatrix}$$

$\beta_{1:T}$ takes $O(Tk^2)$ operations, linear in T

Forward Backward : Application I

① Prob. of being in state k at time t online:

$$P(z_t=k | x_{1:t}) = \frac{a_t(k)}{\sum_s a_t(s)}$$

via argmax we can get most likely state

② compute the probability of being in state k at time t offline

$$\gamma_t(k) = P(z_t=k | x_{1:T}) = \frac{a_t(k) \beta_t(k)}{\sum_s a_t(s) \beta_t(s)}$$

④ compute probability of two "adjacent" states having specific realizations.

$$q_{t+1}(i,j) = P(z_t=i, z_{t+1}=j | x_1:T) = \frac{\alpha_t(i) A_{ij} \beta_{t+1}(j)}{\sum_{u,v} \alpha_t(u) A_{uv} \beta_{t+1}(v)}$$

MAP inferences in HMM

given $x_1:T \rightarrow$ find most probable hidden state z_1, \dots, z_T

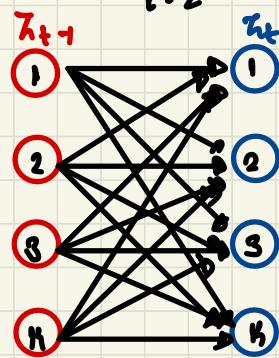
Posterior distribution $P(z_1:T | x_1:T)$

$$\begin{aligned} \arg \max_{z} P(z_1:T | x_1:T) &= \arg \max_z \log \left[P(z_1:T, x_1:T) \right] \\ &= \arg \max_z \log [P(z_1) P(x_1 | z_1)] + \sum_{t=2}^T \log [P(z_t | z_{t-1}) P(x_t | z_t)] \end{aligned}$$

$$\log [P(z_t | z_{t-1}) P(x_t | z_t)]$$

depends on z_{t-1} and z_t

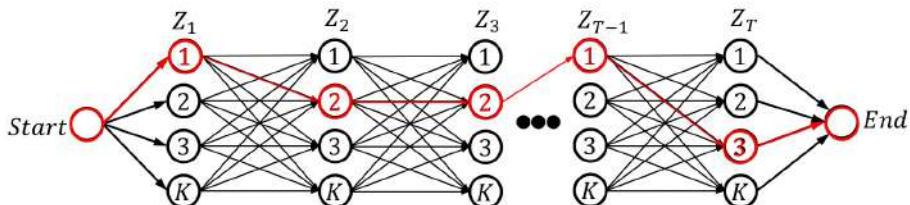
Think of bi-partite graph
weight of the edge $i \rightarrow j$



$$= -\log [P(z_t=j | z_{t-1}=i) P(x_t | z_t=j)]$$

Formulate : viterbi algorithm :

- We can formulate the MAP inference as a shortest-paths problem.
 - weights of edges connected to the *Start* node: $-\log[\Pr(Z_1 = j) \Pr(X_1|Z_1 = j)]$
 - weights of the intermediate layers: $-\log[\Pr(Z_t = j|Z_{t-1} = i) \Pr(X_t|Z_t = j)]$
 - weights of the edges connected to the *End* node: 0



Each directed path corresponds to an assignment to variables $Z_{1:T}$.
 Sum of edge weights = $-\log \Pr(Z_{1:T}, X_{1:T})$

complexity: $O(TK^2)$

Called Viterbi algorithm

Parameter learning :

model parameters $\theta = \{\pi, A, B\}$

goal: Solve $\max_{\theta} \log p_{\theta}(x)$

using ELBO :

$$\begin{aligned} \mathbb{E}_{P(Z_{1:T}|X_{1:T}, \theta^{old})} [\log P(X_{1:T}, Z_{1:T}, \theta)] &= \sum_k P(Z_1 = k | X_{1:T}, \theta^{old}) \log(\pi_k) \\ &\quad + \sum_{i,j} \sum_t P(Z_t = i, Z_{t+1} = j | X_{1:T}, \theta^{old}) \log(A_{ij}) \\ &\quad + \sum_i \sum_t P(Z_t = i | X_{1:T}, \theta^{old}) \mathbb{I}(x_t = j) \log(B_{ij}) \end{aligned}$$

using forward-backward, blue terms can be computed.

Questions

1. Does the sequence $[Z_1, \dots, Z_T]$ fullfill the Markov property ? Why ?

1) yes as $\Pr(z_t | z_{t-1}, \dots, z_1) = \Pr(z_t | z_{t-1})$.

2. Does the sequence $[X_1, \dots, X_T]$ fullfill the Markov property ? Why ?

2) $\Pr(x_t | x_{1:t-1}, z_{1:t}) = \Pr(x_t | z_t)$.
 $z_t \rightarrow z_{t-1} \quad \therefore \text{not Markovian}$

3. In Variational Inference we sometimes need to approximate the ELBO (for example by sampling from the latent Z from the variational distribution). For learning HMM parameters, do we need to sample Z as well?

3) No, as closed form sol'n of forward backward is possible

Problem 1: Consider an HMM where hidden variables are in $\{1, 2\}$ and observed variables are in $\{a, b, c\}$. Let the model parameters be as follows:

$$A = \begin{bmatrix} 1 & 2 \\ 1 & 0.2 & 0.8 \\ 2 & 0.5 & 0.5 \end{bmatrix} \quad B = \begin{bmatrix} a & b & c \\ 1 & 0.2 & 0 & 0.8 \\ 2 & 0.4 & 0.6 & 0 \end{bmatrix} \quad \pi = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

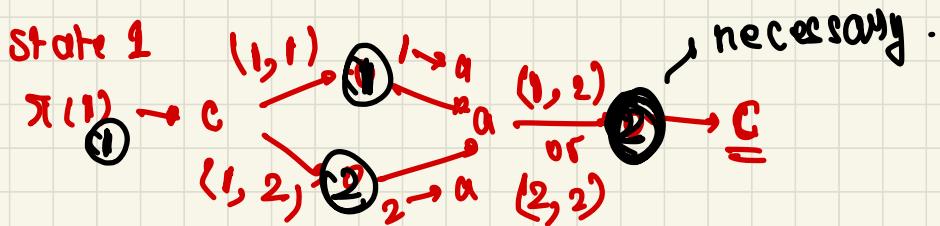
Assume that the sequence $X_{1:5} = [cabac]$ is observed.

1. Filtering: find the distribution $P(Z_3|X_{1:3})$.
2. Smoothing: find the distribution $P(Z_3|X_{1:5})$.
3. Viterbi algorithm: find the most probable sequence $[Z_1, \dots, Z_5]$.

i) Filtering $P(Z_3|X_{1:3})$ use forward algorithm

$$X_{1:5} = [c \underset{\sim}{a} \underset{\sim}{b} a c]$$

$$c \rightarrow a \rightarrow b$$



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

as state 1 $\rightarrow b$ is not possible

naive soln: $\alpha_1 = \pi_1 \Theta B_{1:c} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \Theta \begin{bmatrix} 0.2 \\ 0.8 \end{bmatrix} = \begin{bmatrix} 0.4 \\ 0 \end{bmatrix}$

$\alpha_2 = \Theta B_{2:a} : A^T \alpha_1 = \underline{\quad}$

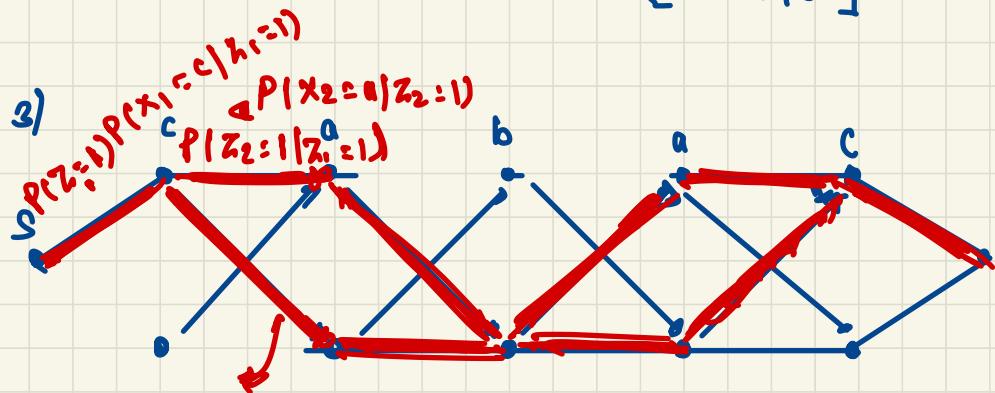
2) smoothing: We will use Forward Backward Algorithm

Clever way is by observation because $x_3 = b$

$$P(z_3 | x_{1:3}) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\beta_3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \beta_4 = A(\beta_{:c} \odot \beta_5) = \begin{bmatrix} 0.16 \\ 0.4 \end{bmatrix}$$

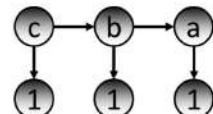
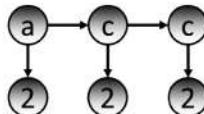
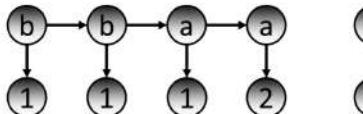
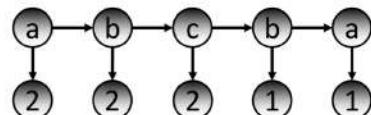
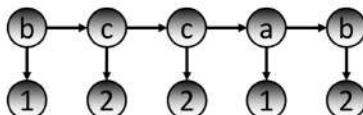
$$\beta_3 = A(\beta_{:a} \odot \beta_4) = \begin{bmatrix} 0.1344 \\ 0.096 \end{bmatrix}$$



$$\Rightarrow P(x_2=a | z_2=1)$$

• calculate the shortest path.

Problem 2: Consider an HMM where states Z_t are in $\{a, b, c\}$ and emissions X_t are in $\{1, 2\}$. Given is the following set of fully-observed instances (two sequences of length 5, one sequence of length 4, and two sequences of length 3):



Learn the parameters of the HMM (i.e. $\pi \in \mathbb{R}^3$, $A \in \mathbb{R}^{3 \times 3}$, and $B \in \mathbb{R}^{3 \times 2}$) using maximum-likelihood estimation.

$$\log P(S_1) = \log P(Z_1) + \sum_{t=2}^{T_1} \log P(Z_t | Z_{t-1}) + \sum_{t=1}^{T_1} \log P(x_t | Z_t)$$

using indicator function and $A, B \in \mathbb{R}^{3 \times 3}$

$$\in \sum_k I(Z_1=k) \log \pi(k) + \sum_{t=2}^{T_1} \sum_{i,j} I(Z_t=j, Z_{t-1}=i) \log A(i,j)$$

$$+ \sum_{t=1}^{T_1} \sum_{i,j} I(X_t=j | Z_t=i) \log B(i,j)$$

$$\log P(S) = \sum_{n=1}^N \sum_k I(Z_1=k) \log \pi(k) + \sum_{n=1}^N \sum_{t=2}^{T_n} \sum_{i,j} I(Z_t=j, Z_{t-1}=i) \log A(i,j)$$

$$+ \sum_{n=1}^N \sum_{t=1}^{T_n} \sum_{i,j} I(X_t=j | Z_t=i) \log B(i,j)$$

$$= \sum_k \underbrace{\sum_{n=1}^N I(Z_1=k) \log \pi(k)}_{L(k)} + \sum_{i,j} \underbrace{\sum_{n=1}^N \sum_{t=2}^{T_n} I(Z_t=j, Z_{t-1}=i)}_{ML(i,j)} \underbrace{\log A(i,j)}$$

$$+ \sum_{i,j} \sum_{n=1}^N \sum_{t=1}^{T_n} I(X_t=j | Z_t=i) \log B(i,j)$$

$$L = \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix} \quad N = \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} \quad M = \begin{bmatrix} 5 \\ 2 \\ 5 \end{bmatrix}$$

Neural sequence Model

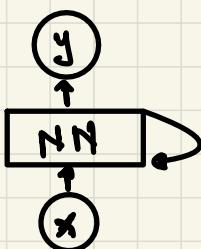
Word2Vec:

- ① CBOW:
 - Predict a word from words surrounding
 - Not good for rare words
- ② skipgram:
 - Predict context given center word
 - good for rare words

Skip grams:

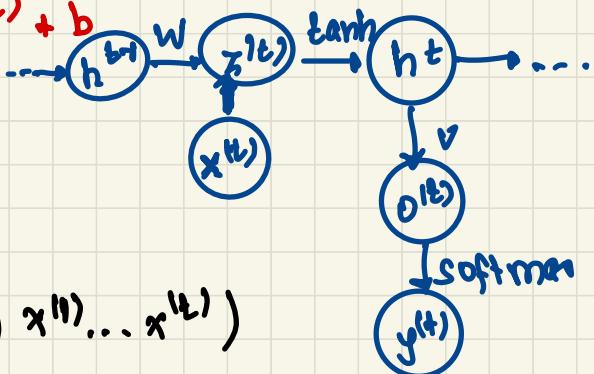
- input: one hot vector with N dimension
- embedding: project the word to D -dimensional space with $U \in \mathbb{R}^{N \times D}$
- Prediction: multiply embedding with $V \in \mathbb{R}^{D \times N}$ and apply softmax
- for $S = \{x_{i-1}, \dots, x_{i-m}, \dots, x_{i+1}, \dots, x_{i+m}\}$ is a window of size l around the word x_i
$$\max \mathbb{E}[P(S|x_i, \theta)] = \min_D (-\mathbb{E}[P(S|x_i, \theta)])$$
where $P(S|x_i, \theta) = \sum_{x_k \in S} P(x_k|x_i, \theta)$ and $P(x_k|x_i, \theta) = \text{softmax}(u_i \cdot v)$

RNN:



- given sequence $\{x^{(1)}, \dots, x^{(n)}\}$ and off $P(y^{(1)}, \dots, y^{(n)}|x^{(1)}, \dots, x^{(n)})$ predict prob

$$\begin{aligned} z^{(t)} &= w^{(t-1)} + v^{(t)} + b \\ h^{(t)} &= \tanh(z^{(t)}) \\ \hat{o}^{(t)} &= v^{(t)} h^{(t)} \end{aligned}$$



weights are shared

$$L = -\log \prod_t p_{\text{model}}(y^{(t)} | x^{(1)}, \dots, x^{(t)})$$

$$= -\sum_t \log p_{\text{model}}(y^{(t)} | x^{(1)}, \dots, x^{(t)}) = -\sum_t L^{(t)}$$

$$\frac{\partial L}{\partial V} = \sum_t (\hat{y}^{(t)} - y^{(t)}) (h^{(t)})^T$$

$$\frac{\partial L}{\partial W} = \sum_t \text{diag}(1 - h^{(t)})^2 \frac{\partial L}{\partial h^{(t)}} (h^{(t-1)})^T$$

$$\frac{\partial L}{\partial U} = \sum_t \text{diag}(1 - (h^{(t)})^2) \frac{\partial L}{\partial h^{(t)}} (x^{(t)})^T$$

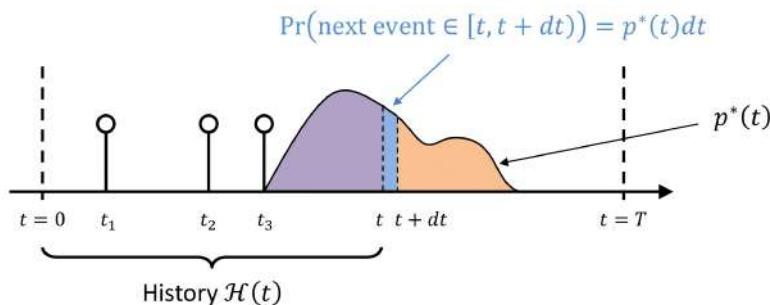
Chapter 9

Temporal point processes

class of probabilistic models that describe the distribution of discrete event sequence in continuous time

defines a generative model for variable-length sequences $t = \{t_1, t_2, \dots, t_n\}$ on the interval $[0, T]$

Modelling the time of next event:



model the distribution $p(\{t_1, \dots, t_n\})$ autoregressively

• predict the time of next event t_i given history $H(t) = \{t_j < t\}$

• important: $H(t)$ depends on the specific sample $\{t_1, \dots, t_n\}$!

conditional density $p^*(t) = p(t | H(t))$

Accumulative density function

$$F^*(t) = \int_{t_{i-1}}^t p^*(u) du = \Pr. \text{ next event happens in } [t_{i-1}, t_i)$$

t_{i-1} is the last event that happened before t_i .

Survival function:

- $S^*(t) = 1 - F^*(t) = \int_t^\infty p^*(u) du \rightarrow \Pr[\text{next event does not take place before } t]$
- or probability that next event happens after t .**

conditional intensity:

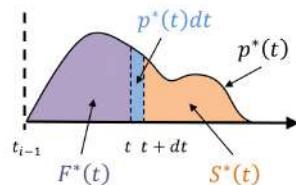
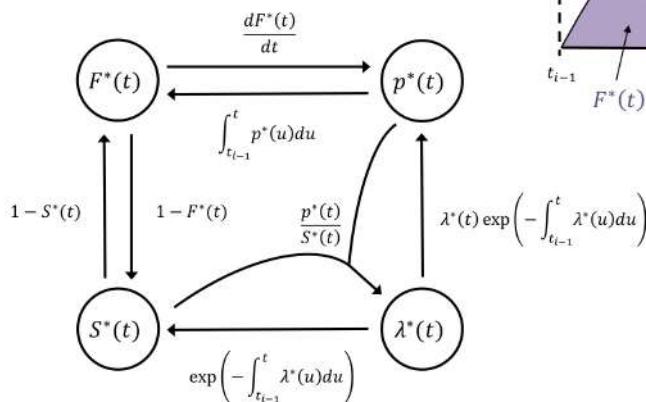
- $\lambda^*(t) dt : \Pr[\text{event in } [t, t+dt)] \text{ given no event in } [t_{i-1}, t)$

$$\lambda^*(t) dt = \frac{\Pr[\text{event in } [t, t+dt) \text{ & no event in } [t_{i-1}, t)]}{\Pr[\text{no event in } [t_{i-1}, t])}$$

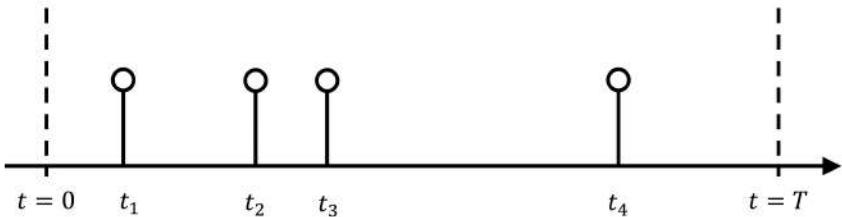
$\cancel{dt+1} \curvearrowleft$

expected # of events / unit time = $\frac{p^*(t) dt}{S^*(t)}$

relationship b/w p^* , F^* , S^* , λ^*



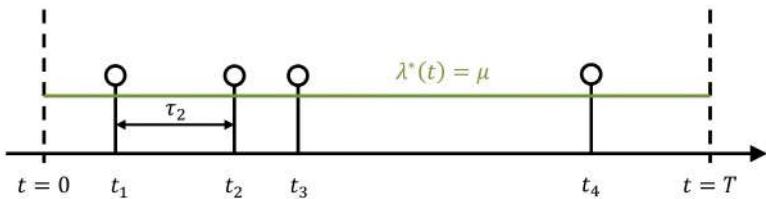
Likelihood of entire sequence



$$\begin{aligned} p(\{t_1, t_2, t_3, t_4\}) &= p^*(t_1) p^*(t_2) p^*(t_3) p^*(t_4) S^*(T) \\ &= \lambda^*(t_1) \lambda^*(t_2) \lambda^*(t_3) \lambda^*(t_4) \exp\left(-\int_0^T \lambda^*(u) du\right) \end{aligned}$$

④ number of events can vary

Homogeneous Poisson Process (HPP) :



- simplest model : constant intensity $\lambda^*(t) = \mu$
- inter event times follow exponential distro.

$$p^*(t) = \mu \exp\left(-\int_{t_{i-1}}^t \mu du\right) = \mu \exp(-\mu(t - t_i))$$

inter-event time t_i

* How to simulate HPP? \rightarrow generating inter-event times

arrival-times = []

$t = 0$

while $t \leq T$:

$\tau_{\text{au}} \sim \text{Exponential}(\lambda)$

$t + \tau_{\text{au}}$

if $t < T$

arrival-times.append(t)

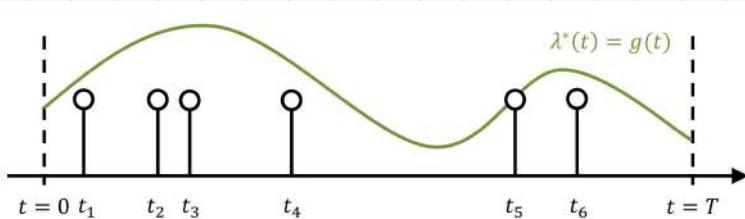
* Sampling from exponential distro: inverse CDF

$$u = F(T) = 1 - \exp(-\lambda T)$$

$u \sim \text{Uniform}(0,1)$

$$T = F^{-1}(u) = -\frac{1}{\lambda} \log(1-u)$$

Inhomogeneous Point Processes (IPP)



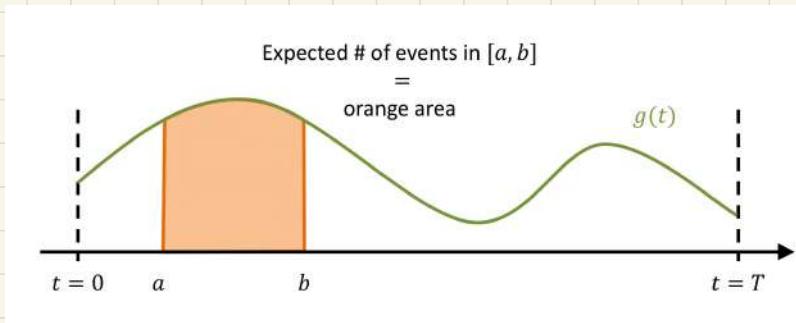
- intensity changes over time: $\lambda^*(t) = g(t) \geq 0$
- intensity is independent of history

Expected number of events

Number of events in $[a, b] \subseteq [0, T]$ follow Poisson

$$N([a, b]) \sim \text{Poisson} \left[\int_a^b g(t) dt \right]$$

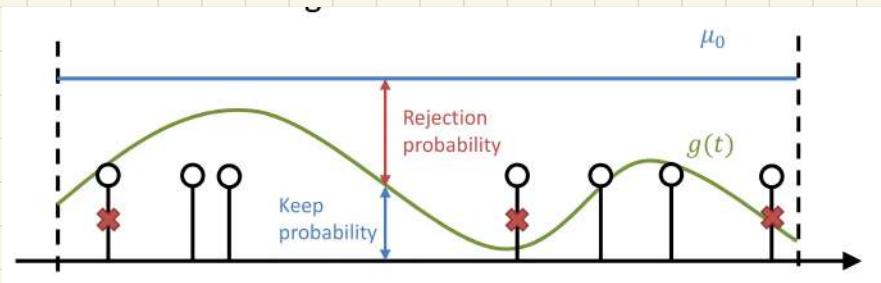
* The expected number of events inside $[a, b]$ is equal to total intensity over the region



combination of two IPPs $\xrightarrow{g(t) \text{ & } h(t)}$ is also an IPP with intensity $g(t) + h(t)$

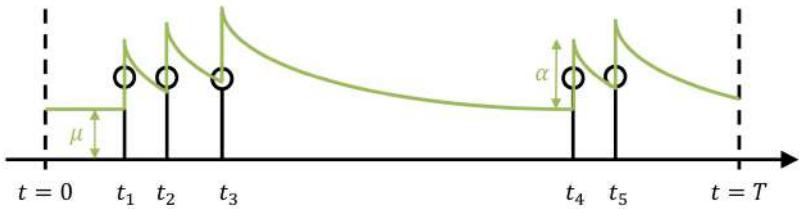
simulating an IPP

integration is hard \rightarrow use the approach: thinning



- 1) Find an upper bound $\mu_0 \geq g(t) \times t$
- 2) simulate candidate events $\{t_i, \dots\}$ from a HPP with rate μ_0 .
- 3) keep each t_i with probability $g(t_i)/\mu_0$.

Hawke's Process



self-exciting process

$$\lambda^*(t) = \mu + \alpha \sum_{t_i \in H(t)} K_w(t - t_i)$$

- Triggering Kernel $K_w(t - t_i) = \exp(-w(t - t_i))$
- Parameters $\mu, \alpha, w \geq 0$
- intensity depends on history

Parameter Estimation in TPPS:

- Parametric conditional intensity $\lambda_\theta^*(t)$
- Maximize the log-likelihood

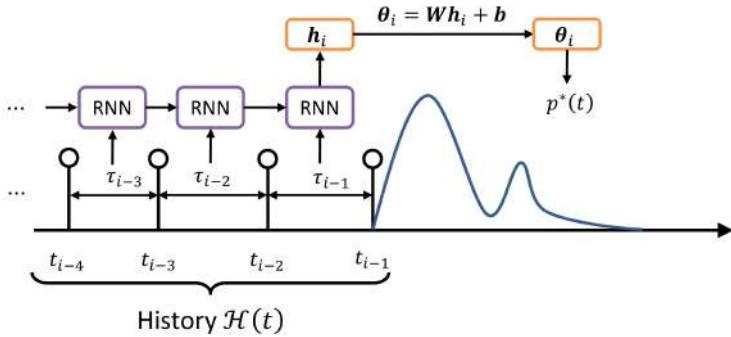
$$\max_{\theta} \sum_{t \in \{t_1, \dots, t_N\}_{\text{train}}} \log p_{\theta}(f|t_1, \dots, t_N)$$

- log likelihood of single sequence

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

different sequences can have different lengths

Modelling TPPs with RNNs :



directly model the conditional distribution $p^*(t)$ using RNN

- 1) Every time an event happens, we feed τ_i into RNN
- 2) hidden state $h_i \in \mathbb{R}^D$ of RNN as history embedding
- 3) use h_i to generate parameters θ_i of distro $p^*(t)$

$$p^*(t) = p(t | H(t)) = p(t|h_i)$$

Modelling $p^*(t)$

Questions :

1. Is it possible to obtain the conditional intensity $\lambda^*(t)$ if you know only the survival function $S^*(t)$ and don't know the conditional PDF $p^*(t)$?

$$1) \quad A^*(t) = \exp \left[- \int_0^t A^*(u) du \right] \asymp$$

69.1

2. Would you use (a) Hawkes process or (b) inhomogeneous Poisson process to model the following event data?

- Customers visiting a supermarket (event = customer enters the supermarket)
- Messages sent by a single user on WhatsApp (event = message sent)
- Taxi rides in a city (event = a trip starts)

2) (i) inhomogenous point process \rightarrow History should not typically excite more occurrences

(ii) Hawke's Process as one message excites and would increase chances of more messages back and forth

(iii) Taxi rides will be inhomogenous point process again as does not depend on history.

3. What can you say about a TPP with the following conditional intensity function? What kind of behavior does it model?

$$\lambda^*(t) = \exp \left(t - \sum_{t_i \in \mathcal{H}(t)} 1 \right)$$

8) $\exp t - \sum_{t_i \in \mathcal{H}(t)} 1$

The intensity decreases as # occurrences increased. This is called inhibiting process.

Problem 1: Consider a temporal point process, where all the inter-event times $\tau_i = t_i - t_{i-1}$ are sampled i.i.d. from the distribution with the survival function

$$S(\tau) = \exp(-(e^{b\tau} - 1))$$

with a parameter $b > 0$.

- Write down the closed-form expression for the conditional intensity function $\lambda^*(t)$ of this TPP. Simplify as far as you can.
- Write down the closed-form expression for the log-likelihood of a sequence $\{t_1, \dots, t_N\}$ generated from this TPP on the interval $[0, T]$. Simplify as far as you can.

$$i) (a) S(\tau) = \exp(1 - e^{b\tau})$$

$$F(\tau) = 1 - \exp(1 - e^{b\tau})$$

$$p(\tau) = b \exp(1 - e^{b\tau}) \exp(b\tau)$$

$$\rho(\tau) = b \exp(1 + b\tau - e^{b\tau})$$

$$\lambda(\tau) = \frac{p(\tau)}{S(\tau)} = \frac{b \exp(1 - e^{b\tau}) \cdot \exp(b\tau)}{\exp(1 - e^{b\tau})} = b \exp(1 + b\tau)$$

$$\lambda(\tau) = b e^{b\tau}$$

$$\lambda^*(t) = b e^{b(t_i - t_{i-1})}$$

$$(b) p(t_1, \dots, t_N) = \left(\prod_{i=1}^N p^*(t_i) \right) S^*(T)$$

$$= \prod_{i=1}^N (b \exp(1 + b\tau - e^{b\tau})).$$

(solve ahead) should be doable

Problem 2: Consider an inhomogeneous Poisson process (IPP) on $[0, 1]$ with the intensity function $\lambda(t) = 2t$. We simulate a sample from this IPP using thinning. For this, we first simulate a *homogeneous* Poisson process (HPP) with intensity $\mu = 4$ and apply the thinning procedure described in the lecture. What is the expected number of events from the HPP that will be rejected when using this procedure?

P2) inhomogeneous Poisson process

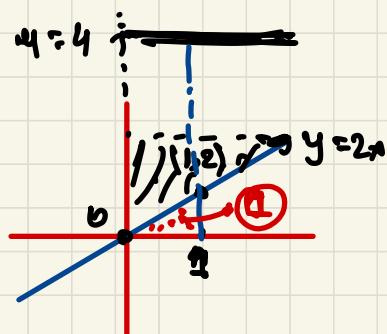
$\lambda(t) = 2t$
approximating using thinning

approximate $\underline{g(t)}$

$$\text{no. of expected events} = \int_{0^+}^1 \lambda(t) dt$$

$$= \int_0^1 2t dt = \boxed{1}$$

$$\text{for HPP} \quad \int_0^1 \lambda(t) dt = \boxed{4}$$



on an average 3 events will be discarded

8

Problem 3: Consider an inhomogeneous Poisson process on $[0, 4]$ with the intensity function $\lambda(t) = \beta t$, where $\beta > 0$ is a parameter that has to be estimated. You have observed a single sequence $\{1, 2.1, 3.3, 3.8\}$ generated from this IPP. What is the maximum likelihood estimate of the parameter β ?

$$P3) \quad \lambda(t) = \beta t \quad S = \{1, 2.1, 3.3, 3.8\}$$

$$S(t) = \exp \left(- \int_{t_{i-1}}^t \beta u du \right)$$

$$= \exp \left(- \int_0^t \beta u^2 du \right) = \frac{e^{-\beta t^2/2}}{\sqrt{\pi}}$$

$$F^*(t) = 1 - e^{-\beta t^2/2} \quad p(t) = \beta e^{-\beta t^2/2}$$

$$P(S(t)) = \prod_{i=1}^4 P^*(t_i) \cdot S(t)$$

$$\begin{aligned} \log(P(S(t))) &= \sum_{i=1}^4 \log(P^*(t_i)) + \log(S(t)) \\ &= \sum_{i=1}^4 \left[\log(\beta e^{-\beta t_i^2/2}) \right] + \log(e^{-8\beta}) \\ &\therefore n \log(\beta) - \frac{\tau^2}{2} \beta \end{aligned}$$

$$\frac{n}{\beta} - \frac{\tau^2}{2} = 0$$

$$\boxed{\beta = \frac{2n}{\tau^2}} \rightarrow \boxed{\frac{1}{2}}$$

#8

Problem 4: Consider a *neural* temporal point process where the conditional intensity function is defined with a neural network. In particular, for a time point t_i , we represent the history $\{t_1, t_2, \dots, t_{i-1}\}$ with a fixed-sized vector $\mathbf{h}_i \in \mathbb{R}^d$. The conditional intensity function $\lambda^*(t)$ is defined as a function of \mathbf{h}_i . We will use the transformer architecture (see previous lecture). We propose the following implementation.

Given the full sequence $\{t_1, t_2, \dots, t_n\}$, we calculate all $\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n\}$ in parallel. We first calculate vectors $\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i \in \mathbb{R}^d$ as a function of t_i . We stack these vectors into matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times d}$. The output of the transformer is: $\mathbf{H} = \text{softmax}(\mathbf{Q}\mathbf{K}^T)\mathbf{V}$, then \mathbf{h}_i is the i th row of \mathbf{H} .

Identify the errors in this implementation compared to the original definition of \mathbf{h}_i . Propose a solution.

- impose right shift and masking (for future)
 $\lambda^*(t) = f(\mathbf{h}_i)$ Masking
 may be get back to this*

Chapter 10 - Graphs & Networks :

Plain / simple graph $G = (V, E)$

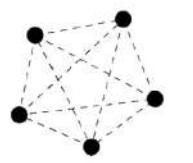
- set of nodes V
- set of edges $E \subseteq V \times V$

May be represented : using binary adjacency matrix

Are real graphs random?

1) Erdős-Renyi Random Graph Model

- start with N (isolated) nodes
- for every pair v_1, v_2 add an edge with probability p
- every edge occurs with equal prob.



• Are real graphs random : No!

Power law distribution:

Definition : - two variables are related by the power law when $y(x) = Ax^{-\gamma}$, A & γ (power law constant) are positive constants.

- A random variable is distributed according to a power law when pdf is given by :

$$p(x) = Ax^{-\gamma}, \gamma > 1$$

Note : Power law distribution looks like a line on log-log.

Characteristic : decay of polynomial is only polynomial

Note: graphs with degree distributions following a power law is called scale free.

Remark: We don't claim that the underlying distro is power law distribution

instead: in many cases a power law distribution is a good description / fit

Gaussian Trap: Be careful while writing algorithm

Patterns and Algorithms: Patterns allow new algos.

real social networks have many triangles



- clustering coefficient: $C = \frac{3 * \text{number of } \Delta s}{\text{number of connected triples}}$
- computing clustering coefficient requires computing Δs . but, triangles are expensive to compute

Efficient triangle counting:

Recap: Eigen decomposition $Ax = \lambda x$

observed: Power law in the eigenvalues & adjacency matrix.

$$\text{- number of triangles} = \frac{1}{6} \text{trace}(A^3) = \frac{1}{6} \sum \lambda_i^3$$

$\lambda_i \rightarrow$ eigenvalue of adjacency matrix

- Eigenvalues follow power law (highly skewed)

Power Iteration: iterative approach to compute a single eigenvalue and eigenvector

$v \leftarrow$ arbitrary vector, $x \leftarrow$ Matrix

iteratively : $v \leftarrow \frac{x \cdot v}{\|x \cdot v\|}$ (until convergence)

& efficient for sparse data

for multiple eigenvalues : Deflate the matrix

eigen decomposition x : $R \Lambda R^T = \sum \lambda_i r_i r_i^T$

$$\hat{x} = x - \lambda_1 r_1 r_1^T$$

* Real graphs are sparse :

characteristic path length :

- for each starting node: shortest path to every other node
- Take the average length of all these paths
- consider avg length of all starting nodes and then take the median.

$$\text{median} \left\{ \frac{1}{|V|} \sum_{v \in V} \text{len}(\text{path}(v, v_j)) \right\}$$

average diameter (mean instead of median in prev)

$$\frac{1}{|V|} \sum_{v \in V} \frac{1}{|V|} \sum_{v_j \in V} \text{len}(\text{path}(v, v_j))$$

Effective diameter / Eccentricity and Hot plot :

- $N_h(u)$: number of nodes reachable from u via h hops

$$N_h(v) = \{ v \in V \mid \text{len}(\text{path}(u, v)) \leq h \}$$
- Total neighbourhood size N_h is the sum over all the starting nodes

$$N_h = \sum_{u \in V} |N_h(u)|$$

Hot plot : plot of N_h vs h .

• effective diameter (or eccentricity) :

- Minimum number of hops in which some fraction of all connected pairs of nodes can reach each other.

$$\min \{ h \in \mathbb{N} \mid N_h \geq 0.9 N^2 \}$$

Chapter 10 : Graphs generated models :

Models assuming (conditional) independent edges

1) Erdős-Renyi Graph Model

$$A_{i,j} \sim \text{Bernoulli}(p)$$

corresponds to Poisson Distro.

• Degree distribution :

• probability of a vertex having degree p is :

$$P_K = \binom{N-1}{K} \cdot p^K \cdot (1-p)^{N-1-K} \approx \frac{e^{-\lambda} \lambda^K}{K!}$$

$\lambda \rightarrow \text{avg degree of graph}$

But in real world we have power-law distributions

• Diameter : $\log(N) / \log(z)$ N : avg node degree

↳ Diameter increase slowly with no. of nodes

But in real world it is constant or it shrinks

- clustering coefficient :
connection probability $p = \frac{2}{N} m$

To generalize to real world data

Idea : $\begin{cases} \text{nodes of same community connect with prob. } p \\ \text{nodes of different community connect with prob. } q, \quad p > q \end{cases}$

Planted Partition Model :

- set of nodes V , partitioned into 2 communities C_1 and C_2
- generate edge b/w every pair of nodes with prob.
 $\Pr(A_{ij} = 1 | z_i, z_j) = \begin{cases} p, & \text{if } z_i = z_j \\ q, & \text{else} \end{cases}$

Stochastic Block Model :

generalizes PPM to graphs with arbitrary numbers and sizes of communities, and varying edge densities

- random variables :
 - $z_i \in \{1, \dots, K\}$: node i belongs to block l community
 - $A \in \{0, 1\}^{N \times N}$: adjacency matrix.

• Model Params :

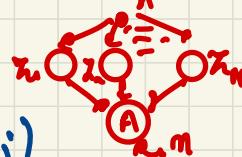
- $\Pi = [\pi_1, \dots, \pi_K]$: community parameters
- π_{uv} : edge probabilities b/w two nodes that are in community $u \& v$

$$\pi = [\pi_1, \dots, \pi_K]$$

• conditional distributions :

- $\Pr(z_i = k) = \pi_k$
- $\Pr(A_{ij} | z_i, z_j) = \text{Bernoulli}(\pi_{z_i z_j})$

$$\pi = \begin{bmatrix} \pi_1, \dots, \pi_K \\ \pi_{K+1}, \dots, \pi_{NK} \end{bmatrix}$$



Preferential Attachment Model :

- ER / PPM / SBM assume edges are generated independently. nodes at beginning; edges Bernoulli distro.

Now generate network based on two processes:

- (i) growth: instead of starting with all nodes, start with a small set of nodes and let the network grow over time by adding new nodes and edges.
- (ii) preferential treatment: "rich get richer" probability of connecting nodes is proportional to the current degree of the nodes.

Note: "rich get richer" principle leads to power law

Initial attractiveness is an extension of power law degree distribution.

algorithm:

- 1) start with m_0 nodes
- 2) add new node w
- 3) simultaneously insert m directed edges (u, v)
probability that the endpoint of an edge (u, v)
corresponds to $A_v = A + \text{indeg}(v)$
- 4) A : initial attractiveness
 $\text{indeg}(v)$ is the number of currently 1.1 incoming edges.
- 4) goto 2 until required number of nodes are obtained

Initial attractiveness : Properties :

Params: m (new edges per step) A (initial attractiveness)

Degree distribution: follows a power law with exponent $\gamma = 2 + A/m$

→ Matches real world data

Diameter : Diameter grows $\Theta(\frac{\log N}{\log \log N})$

- Matches small world effect
- But: still slightly increasing in contrast to real world data.

Average degree: remain constant over time

Configuration Models:

generating networks with arbitrary specified degree distribution.

- 1) assign degree to each nodes, represented as stubs or half-links.
- 2) Randomly select a stub pair and connect them
 - Depending on the order and way in which the stubs were chose, we obtain different

Preserves degree structure



$$k_1 = 3 \quad k_2 = 2 \quad k_3 = 2 \quad k_4 = 1$$

$\nearrow p_1 \quad \nwarrow p_2 \quad \nearrow p_3 \quad \nwarrow p_4$

Deep Generative Networks:

all classical approaches \rightarrow "hand crafted features"
need something to learn from data.

NetGAN : First GAN for networks.

- How to adopt GAN ideas to graph?

Naive $O(N^2)$ solution: learning a distribution over adjacency matrix. Not permutation invariant

solution:

• Learn distribution over random walks

Random Walk :

- Graph $G = (N, E)$, walk length t
- Define $W^t v_i = (w_j : 0 \leq j < t; w_0 = v_i)$ a random walk with starting node v_i
- w_j is the j -th node in the random walk, and the nodes are sampled from

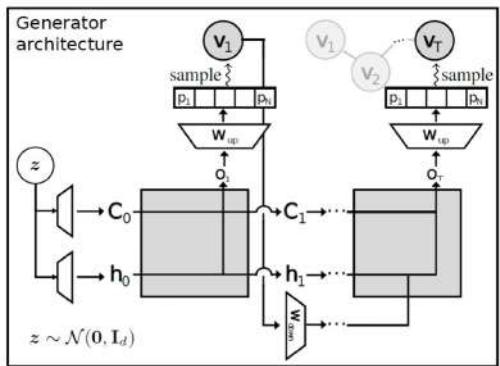
$$\Pr(w_{j+1} = v_u | w_j) = \begin{cases} 1/d_{w_j} & \text{if } v_u \in N(w_j) \\ 0, & \text{else} \end{cases}$$

$\forall j : 0 \leq j \leq t$

- $N(v_u)$ is the set of neighbours of node v_u & d_{v_u} is degree

NetGAN : generator :

$$\begin{aligned} z &\sim \mathcal{N}(0, I_d) \\ m_0 &= g_{\theta'}(z) \\ v_1 &\sim \text{Cat}(\sigma(p_1)), \quad (p_1, m_1) = f_{\theta}(m_0, 0) \\ v_2 &\sim \text{Cat}(\sigma(p_2)), \quad (p_2, m_2) = f_{\theta}(m_1, v_1) \\ &\vdots \\ v_T &\sim \text{Cat}(\sigma(p_T)), \quad (p_T, m_T) = f_{\theta}(m_{T-1}, v_{T-1}) \end{aligned}$$



generative process is as follows

- 1) sample latent noise from normal distribution $z \sim \mathcal{N}(0, I_d)$
- 2) Pass z through NN to obtain starting states c_0 and h_0 for LSTMs
- 3) Run the LSTM from this state to generate a sequence of nodes (random walk)

Problem: Categorical distribution is not differentiable with respect to parameters π .
 - same as variational inference.

We need

- a) Discrete samples to feed into discriminator
- b) gradients to train generators

Solution: Categorical reparametrization trick

- Random variable g is said to have standard Gumbel distribution if $g = -\log(-\log(u))$ with $u \sim \text{Uniform}[0,1]$
- Let v be discrete random variable $P(v=k) \propto \pi_k$ and let $\{g_k\}_{k \in K}$ be i.i.d sequence from Gumbel random variables.

$$v = \arg \max_k [g_k + \log \pi_k]$$

gumbel trick

- Recipe for (reparametrization) sampling from categorical distribution.
- 1) Draw Gumbel noise by transforming uniform samples
- 2) Add it to $\log \pi_k$ which only has to be known up to a normalizing constant.
- 3) Take the value k that produces the maximum

Problem: $\arg \max$ in $v = \arg \max_k [g_k + \log \pi_k]$ is still not differentiable.

Solution: Approximate it with softmax + temperature parameter T .

$$\text{so } v^* = \sigma \left(\frac{\log \pi_k + g}{T} \right)$$

$$\text{as } T \rightarrow 0, v^* \rightarrow v$$

\rightarrow Note: this is only for backward pass
 forward pass is $\arg \max_T (v)$

Questions :

- How much memory do you need to store the edges of a graph with 1000 nodes and 10,000 edges in a dense adjacency matrix? How much for a sparse matrix?

in dense $\rightarrow 10^3 \times 10^3 = 1000 \times 1000 = 10^6$
By Bytes / memory $\approx 4\text{MB}$

sparse
1 entry $\approx 2\text{ pointers} \rightarrow (10^4 \text{ edges} \times 16)$
 $8\text{ Byte per pointer} \approx 0.16\text{ MB.}$

- What is the average degree in an Erdös-Renyi graph with edge probability p ? And in a real world sparse graph with $O(E) = O(N^\alpha)$?

each vertex can have upto $(N-1)$ edges each with a probability p .

Avg degree: $(N-1)p$

In real world : $O(E) = \underline{O(N^\alpha)}$

No. of vertex per edge = 2

$$\frac{2 \cdot N^\alpha}{2} = \underline{2 N^{\alpha-1}}$$

- Can the Erdös-Renyi model generate all graphs that the Stochastic Block Model can generate?

yes as all the edges have an equal probability p .

There is a possibility that ER model generates all possible combinations in SBM model.

- Is the Initial Attractiveness model equal to the Erdös-Renyi model as $A \rightarrow \infty$?

Doubt

- Why does π not need to be normalized for the Gumbel trick to work?

let $\hat{\pi} = c \cdot \pi$ $\log(\hat{\pi}) = \log(c) + \log(\pi)$

argmax is insensitive to constant additions

Problem 1: An unweighted undirected graph represented by an adjacency matrix $A \in \{0, 1\}^{N \times N}$ is given. Prove that the number of triangles in the graph is equal to $\frac{1}{6} \text{trace}(A^3)$ and that this term is in turn equal to $\frac{1}{6} \sum_i \lambda_i^3$ where λ_i are the eigenvalues of the adjacency matrix A . Hint: Show first that A_{ij}^k is the number of walks of length k from node i to node j .

P1) base case $A^0 = I_N$ is true

$$A_{ij}^k = \sum_{v=1}^N A_{iv} A_{vj} = \sum_{v=1}^N A_{vj} \stackrel{k-1}{=} \begin{cases} 1 & \text{if } \exists \text{ edge between } v \\ 0 & \text{otherwise} \end{cases}$$

for A_{ii}^3 ; number of walks of length 2 from node i to i .

In undirected graph this walk is a triangle

so each triangle is counted ~~2x2~~

\swarrow
2 nodes

\searrow
2 directions

$$\frac{1}{6} \sum_i A_{ii}^3 = \frac{1}{6} \text{trace}(A^3)$$

is number
of triangles

$$\text{trace}(A) = \sum_i \lambda_i \quad (\text{sum of eigenvalues})$$

If v is eigenvector of A then for A^3

$$\boxed{\sum_i A_{ii}^3 v} = \lambda^3 v$$

v is eigenvector with
 $\lambda = \lambda^3$

Problem 2: Given is an Erdős-Renyi graph consisting of N nodes, with the edge probability $p \in [0, 1]$. Derive the probability p_k that a node in the graph has degree equal to exactly k .

$$d \sim \text{Binomial}(N-1, p)$$

$$p_k = \binom{N-1}{k} p^k (1-p)^{N-1-k}$$

Problem 3: Given is an Erdős-Renyi graph consisting of N nodes with edge probability $p \in [0, 1]$. What is the expected number of triangles in this graph?

total number of possible triangles $\binom{N}{3}$
 we want this set \mathcal{T}
 let $x_t = 1$ if triplet $t \in \mathcal{T}$
 else 0

$x_t \sim \text{Bernoulli}$ with prob $1-p^3$

$$\mathbb{E}_A \left[\sum_{(i,j,k) \in \mathcal{T}} x_{(i,j,k)} \right] = \sum_{(i,j,k) \in \mathcal{T}} \mathbb{E}[x_{(i,j,k)}] = \sum_{(i,j,k) \in \mathcal{T}} p^3 = \binom{N}{3} p^3$$

$$2! \left(\binom{\frac{N}{2}}{3} a^3 + \binom{\frac{N}{2}}{2} \binom{\frac{N}{2}}{1} ab^2 \right)$$

Problem 5: Compare the two following graph generation processes.

- Graph G_1 is generated by a stochastic block model. It consists of N nodes partitioned into $K=2$ communities. Both communities consist of exactly $N/2$ nodes, and $\eta = \begin{pmatrix} a & b \\ b & a \end{pmatrix}$.
- Graph G_2 is an Erdős-Renyi graph of N nodes and edge probability p . $\rightarrow Nc_3 p^3$ (no. of Δ^3)

Given the probabilities a and b , for which values of p will the expected number of triangles in G_2 be larger than the expected number of triangles in G_1 ?

SBB $\rightarrow N$ nodes $K=2$ partitions $\eta = \begin{pmatrix} a & b \\ b & a \end{pmatrix}$

$$\leftarrow \binom{N}{3} p^3 > 2 \left[\left(\binom{N/2}{3} a^3 + \binom{N/2}{2} \binom{N/2}{1} ab^2 \right) \right]$$

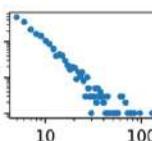
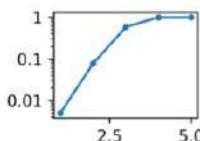
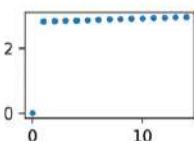
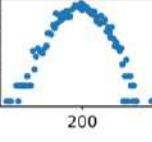
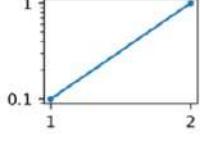
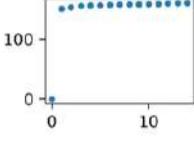
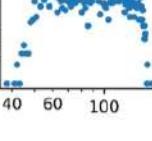
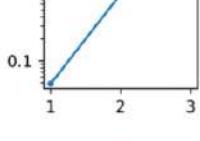
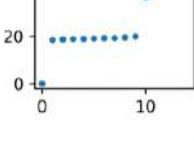
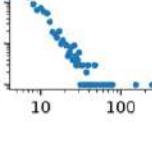
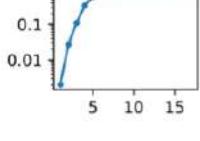
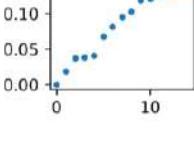
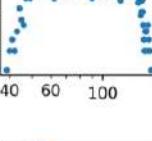
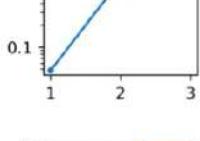
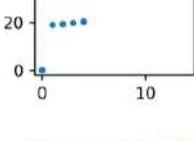
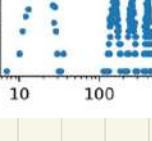
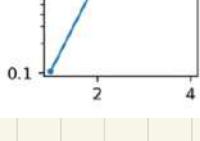
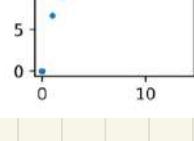
Discuss tomorrow .

Problem 4: Given are 6 graphs $\{G_1, \dots, G_6\}$, which exhibit the properties listed in Table 1. Five of them have been synthetically generated, while one is a real graph. Assign the graphs $\{G_1, \dots, G_6\}$ to the following models (one each) and briefly justify each answer!

- a) Erdős-Renyi model
- b) Stochastic block model with 5 clusters
- c) Stochastic block model with 10 clusters
- d) Stochastic block model with core-periphery structure
- e) Initial attractiveness model
- f) Real graph

These two follow power law G_1, G_4 ?

Table 1: Graphs $\{G_1, \dots, G_6\}$

ID	Degree distribution	Hop plot	Smallest eigenvalues	Clustering coeff.
G_1				0.013
G_2				0.100
G_3				0.275
G_4				0.278
G_5				0.145
G_6				0.191

Chapter 11 : Graph Clustering :

input graph: $G = (V, E)$

aim : find clusters of vertices in the graph .

clustering in network data :

Basic aim

- vertices in the same clusters must be connected by many edges
- only few edges b/w different clusters

Two main clustering algorithms are .

cluster

- Partition approaches: each vertex to exactly one ,
- Non Partition approaches: clusters can overlap ,
outliers that belong to no cluster possible .

A Partitioning Approaches :

idea : constrained optimization problem
given $G = (V, E)$ partition vertex set into clusters

$$C = \{C_1, \dots, C_K\} \text{ s.t.}$$

- given objective function $\Phi(C)$ is optimized
- subject to constraints
 - $C_1 \cup C_2 \cup \dots \cup C_K = V$ (Every vertex from V belongs to a cluster)
 - $\forall 1 \leq i < j \leq K : C_i \cap C_j = \emptyset$ (clusters are disjoint ; no overlap)
- usually NP-hard \rightarrow discrete optimization
- CF - K Means : objects are partitioned such that Total distance is maximized.

Global Minimum Cut :

First Idea: minimize number of edges/weights b/w clusters
 - small value of cut

Task: Partition \mathcal{V} into two sets C_1 and C_2 , such that the sum of inter cluster edge weights $\text{cut}(C_1, C_2) = \sum_{\substack{v_i \in C_1 \\ v_j \in C_2}} w(v_i, v_j)$ is minimized

here : undirected edges, i.e. $w(a, b) = w(b, a)$

Note: $s-t$ cut (where source and sink are given)
 can be done in Ford Fulkerson. (Polynomial time)

However in graph clustering no source | sink is given.
 any partitioning is possible

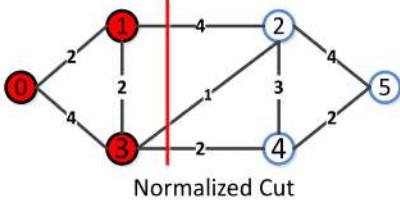
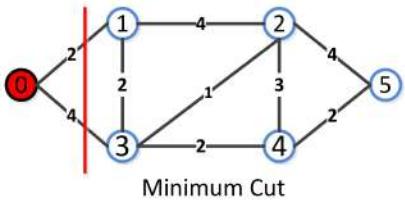
Normalized cut criteria : **Drawbacks of Minimum Cut**
 → cut small vertex set
 → only inter cluster edges are considered.

∴ normalized cut criterion were introduced

- Ratio cut : minimize $\text{cut}(C_1, C_2) + \text{cut}(C_2, \mathcal{V})$
 $\times \text{no. of nodes} \sim \frac{|C_1|}{|C_2|}$

- Normalized cut : minimize $\frac{\text{cut}(C_1, C_2)}{\text{vol}(C_1)} + \frac{\text{cut}(C_2, \mathcal{V})}{\text{vol}(C_2)}$

Note: $\text{vol}(C_i) = \text{assoc}(C_i, \mathcal{V}) = \text{cut}(C_i, \mathcal{V}) =$
 $\sum_{\substack{v_i \in C_i \\ \forall j \in \mathcal{V}}} w(v_i, v_j) \propto \sum_{v_i \in C_i} \deg(v_i)$



only minimize inter
cluster interactions
hence produces a
not so useful cut
as small clusters would
form

$$\frac{\text{cut}(C_1, C_2)}{|C_1|} + \frac{\text{cut}(C_2, C_1)}{|C_2|}$$

A Multiway Graph Partitioning :

- generalization to $K \geq 2$ clusters.
- Partition \mathcal{V} into K C_1, \dots, C_K such that:
 - cut : $\min_{C_1, \dots, C_K} \sum_{i=1}^K \text{cut}(C_i, V \setminus C_i)$
 - Ratio cut : $\min_{C_1, \dots, C_K} \sum_{i=1}^K \frac{\text{cut}(C_i, V \setminus C_i)}{|C_i|}$
 - Normalized : $\min_{\text{cut}} \sum_{i=1}^K \frac{\text{cut}(C_i, V \setminus C_i)}{\text{vol}(C_i)}$

Optimal solution is NP Hard.

A Graph Laplacian :

Definition : Laplacian matrix $L = D - W$

- W = (weighted) adjacency matrix, D = degree matrix

- i.e $L_{UV} = \begin{cases} -w_{UV} & \text{if } (U, V) \in E \\ \deg(V) & \text{if } U = V \\ 0 & \text{otherwise} \end{cases}$

* - Observation 1: For any vector f we have

$$f^T \cdot L \cdot f = \frac{1}{2} \sum_{(u,v) \in E} W_{uv} (f_u - f_v)^2$$

Graph Laplacian :

- discrete analogous of the Laplacian $\sum \frac{\partial^2 f}{\partial x_i^2}$
measures extent a function differs at a point from its value at nearby points
- (discrete) Laplacian: operator on n-dim vector space of functions $f: V \rightarrow \mathbb{R}$
- Laplacian transforms f to another function g
$$g(v) = (Lf)(v) = \sum_{(u,v) \in E} W_{uv} \cdot [f(v) - f(u)]$$
- L is symmetric and positive semi-definite
(as D & W are symmetric) L is PSD (quadratic)
- smallest eigen value of L is 0 - corresponds to $\vec{1}$
- L has n - non negative, real valued eigenvalues
$$0 = \lambda_1 \leq \lambda_2 \dots \leq \lambda_n$$
 - in general symmetric matrices has real-valued eigenvalues
 - $\lambda_1 = 0$ (smallest ent rest are larger)
- Laplacian is additive $L_{G \cup H} = L_G + L_H$
(G and H are 2 graphs on the same vertex set with disjoint edge sets)

- Laplacian of a graph is sum of Laplacians for each edge
 - $L_G = \sum_{(u,v) \in E} L(u,v)$
- allows us to prove properties for single edge and add them up.

Laplacian for an edge i.e. graph with single edge (u,v)

$$L(u,v) = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Observation: No. of eigenvectors of L with eigenvalue 0 corresponds to the number of connected components

* corollary: if the graph is connected: $\lambda_2(L) > 0$

Properties of the Graph Laplacian:

Algebraic connectivity of a graph is $\lambda_2(L)$

- also known as Fiedler eigenvalue.
- magnitude represents how well connected graph is
- $\lambda_2(K_n) = n$, where K_n is complete graph with n nodes

- for every $S \subset V$ we have

$$\Theta(S) = \frac{\text{cut}(S, \bar{S})}{|S|} \geq \lambda_2 \left(1 - \frac{|S|}{|V|} \right)$$

- $\Theta(S)$ is isoperimetric ratio of S

- $\Theta_G \stackrel{\text{def}}{=} \min_{S \subset V} \Theta(S) \geq \lambda_2/2$

- Θ_G is called the Cheeger constant of graph or conductance of a graph

Inequality implies if λ_2 is big, graph is well connected boundary of each small set of vertices is at least λ_2 times a value close to 1.

conclusion: spectrum of L encodes useful info in graph

unfortunately: co-spectral graphs exist

- graphs that are not isomorphic but share same spectrum
- implies that the spectrum doesn't completely characterize the graph.

so spectrality can be useful: find sparse graph with (approx) the same spectrum.

Graph Laplacian and the minimum cut:

- for $k=2$ clusters let $f \in \{-1, 1\}^n$ be indicator vector

$$f_{C_1}[i] = \begin{cases} +1 & \text{if } v_i \in C_1 \\ -1 & \text{if } v_i \in V \setminus C_1 = C_2 \end{cases}$$

- Then $f_{C_1}^T L f_{C_1} = 4 \cdot \text{cut}(C_1, C_2)$

- Thus we can minimize $f_{C_1}^T L f_{C_1}$ to minimize the cut.

But min cut is not ideal so we expand this to ratio cut.

$$\text{ratio-cut } (k=2): \min_{C_1 \subset V} \frac{\text{cut}(C_1, C_2)}{|C_1|} + \frac{\text{cut}(C_2, C_1)}{|C_2|}$$

where $C_2 = V \setminus C_1 = \bar{C}_1$

indicator vectors

$$f_{C_1}[i] = \begin{cases} + \sqrt{\frac{|C_1|}{|C_2|}} & \text{if } v_i \in C_1 \\ - \sqrt{\frac{|C_1|}{|C_2|}} & \text{else} \end{cases}$$

- $\sum_i f_{C_1}[i] = 0$ f_{C_1} is orthogonal to $\vec{1}$.
- $f_{C_1}^T \cdot f_{C_1} = \|f_{C_1}\|^2 = 191$ length is constant
- $f_{C_1}^T \cdot L \cdot f_{C_1} = \dots = 191 \cdot \left[\frac{\text{cut}(C_1, C_2)}{|C_1|} + \frac{\text{cut}(C_2, C_1)}{|C_1|} \right]$
 $\rightarrow 191 \cdot \text{ratio_cut}$.

Minimizing the ratio cut is equivalent to :

$$\min_{C \in V} \{ f_C^T \cdot L \cdot f_C \} \text{ subject to } f_C \perp \vec{1} \text{ and } \|f_C\|_2 = \sqrt{191}$$

and f_{C_1} is defined as before

Discrete optim \rightarrow NP Hard

idea : constraint relaxation
drop the discreteness condition (i.e f_C can take any values).

• Result : $\min_{\substack{f_C \in \mathbb{R}^{|V|} \\ \text{in real space}}} \{ f_C^T \cdot L \cdot f_C \}$ subject to $f_C \perp \vec{1}$
Relaxed constraint $\quad \quad \quad$ and $\|f_C\| = \sqrt{191}$

• For symmetric matrix L we have :

$$-\lambda_1 = \min_{\substack{x_1 \in \mathbb{R}^{|V|} \\ \|x_1\|=1}} x_1^T L x_1 \quad \lambda_2 = \min_{\substack{x_2 \in \mathbb{R}^{|V|} \\ \|x_2\|=1}} x_2^T L x_2 \dots$$

$x_2 \perp x_1$

• f_{C_1} is the second smallest eigenvector of L

- $\vec{1}$ is smallest

$$- L \cdot f_{C_1} = \lambda_2 \cdot f_{C_1} \Leftrightarrow f_{C_1}^T \cdot L \cdot f_{C_1} = 191 \lambda_2$$

in case of 2 clusters: sign of f_{C1} gives clusters

general case (spectral clustering):

↳ ($K > 2$): clusters C_1, C_2, \dots

$$h_{C_i}[i] = \begin{cases} \frac{1}{\sqrt{|C_i|}} & \text{if } v_i \in C_i \\ 0 & \text{else} \end{cases}$$

$$- \text{let } H = [h_{C_1}; h_{C_2}; \dots; h_{C_K}]$$

• observations:

$$- H^T H = Id$$

$$- h_{C_i}^T \cdot L \cdot h_{C_i} = \frac{\text{cut}(C_i, \bar{V} \setminus C_i)}{|C_i|} \quad \text{and} \quad h_{C_i}^T \cdot L \cdot h_{C_i} = (H^T L H)_{ii}$$

$$\text{Ratio cut } (C_1, \dots, C_K) = \sum_{i=1}^K \frac{\text{cut}(C_i, \bar{V} \setminus C_i)}{|C_i|} = \sum_{i=1}^K (H^T L H)_{ii} = \text{trace}(H^T L H)$$

• Minimize ratio-cut

$$\min_{C_1, \dots, C_K} \text{trace}(H^T L H) \quad \text{s.t. } H^T H = Id$$

constraint: allow arbitrary value for H

$$\text{Result: } \min_{H \in \mathbb{R}^{N \times K}} \text{trace}(H^T L H) \quad \text{s.t. } H^T H = Id$$

- To find clusters based on eigen value
- Represent vertex by a vector of its corresponding components (spectral embeddings)
- clustering in embedding space yields final results.

Also known as spectral layout.

Remarks :

- Spectral clustering using unnormalized Laplacian :
 $L = D - W$
 \rightarrow approx. of ratio cut
- Spectral clustering using normalized Laplacian
 $L_{\text{sym}} = D^{-1/2} L D^{-1/2}$ → approx. of normalized cut.

Probabilistic Approaches :

- consider graph as realisation (sample) drawn from a generative model
- generative process is controlled by a set of parameters
- communities are modelled with generative process.

① PPM : $\Pr(A_{ij}=1 | z_i, z_j) = \begin{cases} p & , \text{ if } z_i = z_j \\ q & , \text{ if } z_i \neq z_j \end{cases}$

- likelihood of a community assignment

$$\Pr(A|z) = \prod_{i < j} [p^{A_{ij}} (1-p)^{1-A_{ij}}]^{\mathbb{I}[z_i = z_j]} [q^{A_{ij}} (1-q)^{1-A_{ij}}]^{\mathbb{I}[z_i \neq z_j]}$$

- assume p & q are known.

$$\text{simplest soln (MLE)} \quad z^* = \arg \max_{z \in \{-1, 1\}^N} \Pr(A|z)$$

PPM and Spectral Clustering:

ML estimate for z for planted partition model?
 Let us assume that communities are balanced ($G_1 \approx G_2 = N/2$)

$$E_{\text{cut}}(z) = |\{(i,j) \in E \text{ s.t. } z_i \neq z_j\}| = \sum_{(i,j) \in E} I(z_i \neq z_j)$$

$$\text{likelihood: } \Pr(A|z) \propto \left(\frac{(1-p)q}{(1-q)p} \right)^{E_{\text{cut}}(z)}$$

since $p > q$, MLE is equivalent to finding a minimum balanced cut

How can we solve? Spectral clustering

② SBM generalizes PPM with arbitrary numbers of size of communities

$\underbrace{\pi}_{\text{r.g.}}$ $\rightarrow z_i \in \{1, \dots, k\}$ node i belongs to block z_i
 $A \in \{0, 1\}^{N \times N}$, adjacency matrix

Model Params: $\pi = [\pi_1, \dots, \pi_k]$: community proportions
 π_{uv} : edge prob. b/w two nodes.

conditional distn: $\Pr(z_i = x) = \pi_x$

$$\Pr(A_{ij}|z_i, z_j) = \text{Bernoulli}(\pi_{z_i z_j})$$

Inference

- if η and π are given

$$\Pr(z|A, \eta, \pi) = \frac{\Pr(A|z, \eta, \pi) \Pr(z|\pi)}{\Pr(A|\pi)}$$

intractable $O(4^m)$ ops.

use variational inference:

find $q(z)$ similar to true Posterior $\Pr(z|A, \eta, \pi)$
(mean field assumption)

$$\Pr(z|A, \eta, \pi) \approx q(z|\psi) = \prod_{i=1}^N q(z_i|\psi_i)$$

- if both x and A are observed

$$\pi_k^{\text{MLE}} = \frac{\# \text{ nodes in cluster } k}{N} = \frac{N_k}{N}$$

$$\eta_{uv}^{\text{MLE}} : \frac{\text{observe # edge between } u \text{ & } v}{\text{possible # edges between } u \text{ & } v} = \sum_{(i,j) \in E} \frac{\Pr(z_i=u) \Pr(z_j=v)}{P_{uv}}$$

$$P_{uv} = \begin{cases} \binom{N_u}{2} & u=v \\ N_u \cdot N_v & u \neq v \end{cases}$$

no. possible edges

- if only A is observed:

use EIBD wrt $q(z)$ as well as η and π

Questions

- How can you find the connected components of a graph from its Laplacian?

of eigenvalues which have value = 0
are the number of connected components

- Consider a graph with n arbitrarily connected nodes and k disconnected nodes. What are the first $k+1$ clusters that spectral clustering finds? Why?

Spectral clustering would return $\frac{N}{k}$ solo nodes and the connected graph.

Problem 2: Consider minimizing the ratio cut on a graph with two clusters C_1 and C_2 and N nodes in total. The indicator vector

$$f_{C_1,i} = \begin{cases} +\sqrt{\frac{|C_1|}{|C_1|}} & \text{if } v_i \in C_1 \\ -\sqrt{\frac{|C_1|}{|C_1|}} & \text{otherwise} \end{cases}$$

is defined as in the lecture. Prove the following three properties about f_{C_1} .

a) $1^T f_{C_1} = \sum_i f_{C_1,i} = 0$

b) $f_{C_1}^T f_{C_1} = \|f_{C_1}\|_2^2 = |V|$

c) $f_{C_1}^T L f_{C_1} = |V| \left[\frac{\text{cut}(C_1, C_2)}{|C_1|} + \frac{\text{cut}(C_1, C_2)}{|C_1|} \right]$

2) a) $1^T f_{C_1} = \sum_i f_{C_1,i} = \sqrt{M} \sqrt{\frac{N-M}{M}} - \sqrt{N-M} \sqrt{\frac{M}{N-M}}$
 $= \boxed{0}$

b) $f_{C_1}^T f_{C_1} = \|f_{C_1}\|_2^2 = \sum_i f_{C_1,i}^2$
 $= \sum_i \begin{cases} \sqrt{\frac{|C_1|}{|C_1|}} & \text{if } v_i \in C_1 \\ \sqrt{\frac{|C_1|}{|C_2|}} & \text{otherwise} \end{cases}^2 \cdot \frac{N-M}{M} + \frac{M}{N-M} \cdot \frac{M}{M}$
 $= \boxed{N} = \boxed{1}$

$$c) \quad f_{c_1}^T L f_{c_1} = \frac{1}{2} \sum_{(u,v) \in E} w_{uv} (f_{c_1,u} - f_{c_1,v})^2$$

Ch-12 - Graph Node Embeddings:

Learning Node Representations:

- transform the graph s.t each vertex is represented by a vector.
- $\Phi: \mathcal{V} \rightarrow \mathbb{R}^d$ maps each node to \mathbb{R}^d
- Nodes close in the embedding space are 'similar' w.r.t. graph structure.

Types of embeddings:

- difference come from how we define similarity b/w nodes in embedding space.

Spectral clustering: (Recap)

- construct the graph Laplacian L
- compute first k eigenvectors v_1, v_2, \dots, v_k of L in $\mathbb{R}^{n \times k}$.
- represent the i -th node as i -th row of H
- cluster the vector representation, for ex. K-means

Spectral clustering: embedding view

- Spectral clustering is based on eigenvectors of the graph Laplacian L .
 - L encodes structural behaviour of graph G ,
 - $|V| \times |V|$ adjacency matrix is transformed and reduced to $|V| \times k$ - matrix H .
- Relation to PCA and dimensionality reduction.
 - transformation is based on eigenvectors of data - matrix and one retains only eigenvectors with smallest / largest eigenvalues.

12 "Deep" Node embedding approaches:

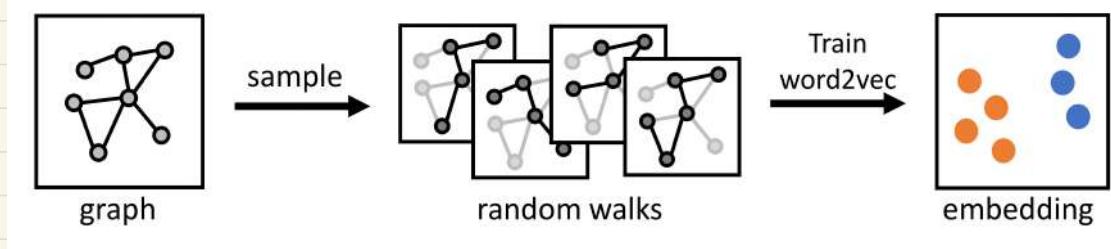
goal : Try to capture more complex structures than spectral clustering.

- thus, hopefully getting better results for specific downstream tasks.

DeepWalk :

idea : transform graph into a set of random walks and learn a word2vec model.

- For every node v_i : sample multiple random walks
- Random walks corresponds to sentences in word2vec both are sequences of discrete tokens
- Train this collection of sentences



Results: nodes that are close to each other in graph and share many neighbours get similar vector representations (embeddings)

Graph2Gauss

Idea : map each node to gaussian distro. in the embedding space such that a node's 1st neighbours are closer than its 2nd neighbours and so forth.

- Learn a mapping : $f_\theta(v) \rightarrow N(\mu_\theta(v), \Sigma_\theta(v))$

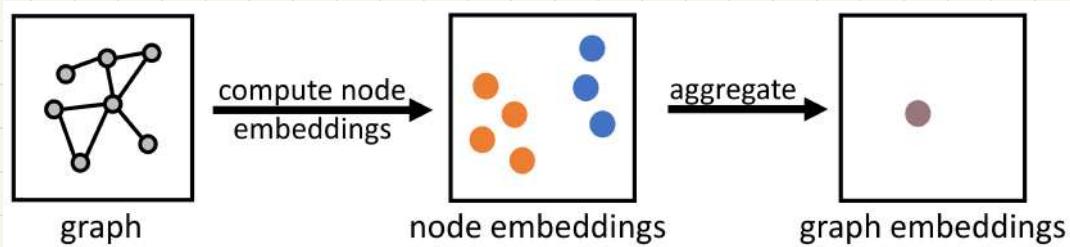
- for each node u define a loss $\sum_{v \sim u} \Sigma_{uv} + e^{-f_u(v)}$
 where $E_{uv} = \text{KL}(f_u(u) || f_u(v))^{(u,v)}$
 u and v are all node pairs such that v is closer to u than t (ranking loss)

Result: distance in embedding space correlates with distances in the graph and the Gaussian variances express how certain the model is about the embedding.

Graph Embedding (vs Node embedding) :

- idea: Leverage node embeddings to embed graphs
 problem: \rightarrow graphs have different number of nodes
 \rightarrow simple concatenation produces incompatible dimensionalities

solution: Aggregate all node embeddings into one graph embedding, for example with mean pooling



#12: Graph Ranking

core idea: A page is important if many important pages point to it.

"voting" principle:

- each page votes for importance of pages it points to
- link's vote & importance of source page
- if page j , with importance r_j has n -outlinks each get r_j/n votes
- page j 's own importance is sum of votes on in-links.

$$\text{rank of page } j : r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i} \quad d_i \rightarrow \text{out degrees}$$

Page Rank: Matrix Formulation:

Stochastic adjacency matrix: M

- if $i \rightarrow j$, then $M_{ji} = 1/d_i$ else $M_{ji} = 0$
- M : column stochastic matrix

Rank vector: \mathbf{r}

- r_i is the importance score of page i
- $\sum_i r_i = 1$

Equation $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$ can be written as

$$\mathbf{r} = M \cdot \mathbf{r}'$$

Notes: 1) r' is eigenvector of M

- eigenvalue 1
- largest eigenvalue is 1 / M is column stochastic

2) Finding π = finding eigenvector of M corresponding to largest eigenvalue.

Computation

Power Iteration : $\bullet \pi \leftarrow \frac{M \cdot \pi}{\|M \cdot \pi\|}$ until convergence.

- required for page rank : $\sum_i \pi_i = 1$
- let $y = Mx$ $\sum x_i = 1$ • M is column stochastic
 $\sum y_i = 1$
- no need to normalize
- start with random (normalized vector) π'

iterate
 $\pi \leftarrow M\pi'$

Iterate till convergence. $\pi_j \leftarrow \sum_{i \rightarrow j} \pi_i / d_i$

Random Walk interpretation :

surfer's path $\{x_1, x_2, \dots\} \rightarrow$ Markov chain
 $\Pr(x_t = i) = \pi_i \quad \Pr(x_{t+1} = j | x_t = i) = M_{ji}$

transition probability matrix of
Markov chain is $B \equiv M^T$

Under some technical conditions :

rank score of page $i = \pi_i = \lim_{t \rightarrow \infty} \Pr(x_t = i)$

or in vector form : $\pi = \lim_{t \rightarrow \infty} \Pi(t)$

$\Pi(t) = \Pi B^{t-1} \Rightarrow$ limit of sequence
 $\Pi B, (\Pi B)B, \dots$

- $\lim_{t \rightarrow \infty} \underline{\pi(t)}$ is called limiting distribution

- $t \rightarrow \infty$, B^t converges to a matrix whose rows are the same
one row of $\lim_{t \rightarrow \infty} B^t \rightarrow$ limiting distr

probability of reaching a node does not depend on starting point.

$$\lim_{t \rightarrow \infty} B^{t-1} = \begin{bmatrix} a & b & c \\ a & b & c \\ a & b & c \end{bmatrix} \Rightarrow \lim_{t \rightarrow \infty} \pi(t) = \lim_{t \rightarrow \infty} \pi B^{t-1} = \begin{pmatrix} \pi_1, \pi_2, \pi_3 \\ \hline a & b & c \\ a & b & c \\ a & b & c \end{pmatrix}$$

$$+ [a \ b \ c]$$

Stationary distribution: $\pi(\infty)$ is called stationary distribution

$$\pi(\infty) = \pi(\infty)B$$

- * $\pi(\infty)$ if exists = transpose of rank vector π .
 $\pi(\infty)$ can be computed by
 - getting eigenvector associated with $\lambda=1$
 - normalizing it to 1.

under "technical conditions": Markov chain has a limiting distribution = unique stationary distribution

technical conditions: irreducible & aperiodic

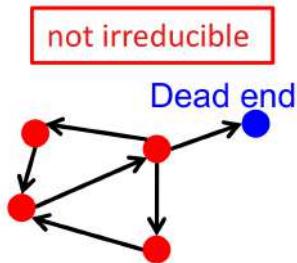
possible to get any state from any state $\exists n \text{ s.t. } \forall n \geq n$

$$\Pr(X_n=i | X_0=i) > 0$$

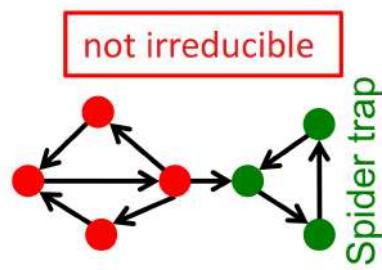
* we might have to prove this

- Markov chain is aperiodic if every state is aperiodic
- irreducible markov chain only need one aperiodic state to imply all states are aperiodic.

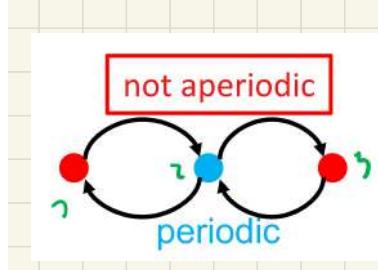
Problems in Page Rank



random walk has nowhere to go
no outlinks



- random walk stuck
- eventually absorbs all importance
- spider trap



if we start at state , we will return to the state in fixed periods

solution : random teleportation

two options

- β follow a link at random
- $1-\beta$ jump to some random page

Page rank equation :

$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1-\beta) \frac{1}{N} \quad \text{if } i \neq j$$

$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + \sum_i (1-\beta) \frac{1}{N}$$

in matrix notation

$$A = \beta M + (1-\beta) \left[\frac{1}{N} \right]_{N \times N}$$

$$\text{final solution : } r = Ar$$

$\underbrace{\quad}_{(N \times N) \text{ matrix with all entries}} = \frac{1}{N}$

Notes on computation

- attention : The matrix A is dense! (N^2 non zero entries)
 - you should never compute r in such way

- consider the teleport by adding constant penalty to each term
- iterate $\pi_j \leftarrow \sum_{i \sim j} \frac{\pi_i}{d_i} + (1-\beta) \frac{1}{N}$ until convergence
only neighbours need to be considered
- To maintain sparsity in matrix form multiply βM
then add a vector

$$r = \beta M r + (1-\beta) \begin{bmatrix} \frac{1}{N} \end{bmatrix}_{N \times N}$$

- vertex oriented computation
 - each vertex performs local computation.

crucial aspects of vertex oriented programming

- GAS principle - gather, apply, scatter
- | | | | |
|--|-----|-----|-----|
| | (a) | (b) | (c) |
|--|-----|-----|-----|
- a) gathers info from the adjacent vertex / edge
 - b) apply transforms
 - c) scatters info to adjacent vertices

PageRank only requires a) and b)

Problems related to page rank :

- Measure generic popularity of a page
 - Biased against topic-specific authorities

solution : Topic-sensitive Page Rank
- Susceptible to link spam
 - artificial link topographies created in order to boost PageRank.

Solution : use trust rank
- use single measure of importance

sol : Hubs and Authority

Topic Sensitive Page Rank:

- goal: Evaluate pages based not just on their popularity, but by how close they are to a particular topic.
- allow search queries to be answered based on the interest of the user
- core idea: Bias the random walk when the walker teleports \rightarrow pick page from S .
 - standard page rank: $S = \text{all pages}$
 $\text{call pages with equal probability}$
 - topic sensitive page rank: $S = \text{set of relevant pages}$
for each teleport get a different vector π_S .
- topic sensitive page rank takes the following:
$$M = \beta M \pi + (1-\beta) \pi \quad \pi_i = \begin{cases} \frac{1}{|S|} & i \in S \\ 0 & \text{otherwise} \end{cases}$$
- generalizing to arbitrary teleports
$$M = \beta M \pi + (1-\beta) \pi \quad \sum_i \pi_i = 1$$

exact solution is $r = (1-\beta)(I - \beta M)^{-1} \pi$

Discovering relevant set:

- create page ranks of different topics
- use context of the query

- **“Normal” PageRank:**
 - Teleports uniformly at random to any node
 - All nodes have the same teleport probability of surfer landing there:
 $\pi = (0.1 \ 0.1 \ 0.1 \ 0.1 \ 0.1 \ 0.1 \ 0.1 \ 0.1 \ 0.1 \ 0.1)^T$
- **Topic-Sensitive PageRank:**
 - Teleports to a topic specific set of pages
 - Nodes can have different probabilities of surfer landing there:
 $\pi = (0.1 \ 0 \ 0 \ 0.2 \ 0 \ 0.5 \ 0 \ 0 \ 0 \ 0.2)^T$
- **Personalized PageRank (Random Walk with Restarts):**
 - Teleport is always to the same node:
 $\pi = (0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0)^T$

Questions :

Ch-13 : Semi supervised Learning

Collective classification :

- Labels are known for some labeled nodes
- goal is to classify the unlabeled nodes

standard assumption : Homophily (aka assortativity)

Label Propagation :

- consider the binary case
- Formal definition, seeds or labelled instances
 - Nodes $\mathcal{V} = S \cup U \rightarrow$ Unlabeled instance
 - symmetric weighted edge matrix $W \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$
 - $y_i \in \{0, 1\}$ for $i \in S$ // given
 - $y_i \in \{0, 1\}$ for $i \in U$ // predict

Idea : Energy minimization $E(y) = \frac{1}{2} \sum_{i,j} w_{ij} (y_i - y_j)^2$

- Smoothness : adjacent nodes should have same class label

Energy minimization as MAP :

- goal : $\min_y E(y) = \min_y \frac{1}{2} \sum_{i,j} w_{ij} |y_i - y_j|^2$
- This is equivalent to optimization problem
 $\arg \min_{y \in \{0, 1\}^N} E(y) = \arg \max_{y \in \{0, 1\}^N} \exp(-E(y))$
- Consider the following approximation to posterior

$$p(y_U | W, y_L) = \frac{1}{Z} \exp(-E(y))$$

$$Z = \sum_y \exp(-E(y)) \rightarrow \underline{\text{normalizing const}}$$

Label Propagation

Two aspects \rightarrow smoothness
Matching the seed labels

$$\min_{y \geq 0, y \in \mathbb{R}^n} \frac{1}{2} \sum_{i,j} w_{ij} (y_i - y_j)^2 \text{ s.t. } y_i = y_i^* \forall i \in S$$

constrained integer optimization problem

We know how to rewrite the problem

$$\min_{y \geq 0, y \in \mathbb{R}^n} y^T L y \text{ s.t. } y_i = \hat{y}_i \forall i \in S.$$

$L = D - W$: graph Laplacian

Solution:

without loss of generality assume that Laplacian matrix is partitioned into blocks for labeled and unlabeled nodes

$$L = \begin{bmatrix} L_{SS} & L_{SU} \\ L_{US} & L_{UU} \end{bmatrix}$$

accordingly let $y = \begin{bmatrix} y_S \\ y_U \end{bmatrix} = \begin{bmatrix} \hat{y}_S \\ y_U \end{bmatrix}$

vector of labels to be learned.

$$y_U = -L_{UU}^{-1} \cdot L_{US} \cdot \hat{y}_S$$

Label Propagation: generalization:

What if we have n labels?

- one hot encoding $y: x: \begin{cases} 1 & \text{if node } i \text{ is of class } \\ 0 & \text{else} \end{cases}$
- energy function: $E(Y) = \sum_{i,j} w_{ij} (y_i - y_j)^T \alpha y_i - y_j$

• Other types of effects encode with compatibility matrix H

• Energy function $E(Y) = \sum_{i,j} w_{ij} (y_i - y_j)^T H (y_i - y_j)$

Label propagation vs SBM :

LP is a discriminative model \rightarrow models the conditional distribution of labels given the similarity graph $p(Y|W)$

SBM is a generative model \rightarrow models $\text{Pr}(A|Z)$ and $\text{Pr}(Z)$

For SBM we get posterior using Baye's Theorem.

• SBM and LP solve a different problem.

Estimate parameters generated
given graph A (unsupervised) $\xrightarrow{\text{predict labels of the}}$
nodes in U given observed
labels W.

Non Graph Data : if we only have vector data

- construct a graph connecting similar points (ex. kNN)
- apply label propagation
- Features are used to construct the graph, but then only graph is used for classification

Transductive Learning :

LP is a special case of Transductive Learning

given : i) set of labelled data $T = \{(x_i, y_i)\}_{i=1..N} \subset \mathcal{X} \times \mathcal{Y}$
ii) set of unlabeled data $U = \{(x_i)\}_{i=1..M} \subset \mathcal{X}$

goal : predict labels only for unlabeled instance u

(AND NOT SOLVE HARDER INTERMEDIATE PROBLEM)

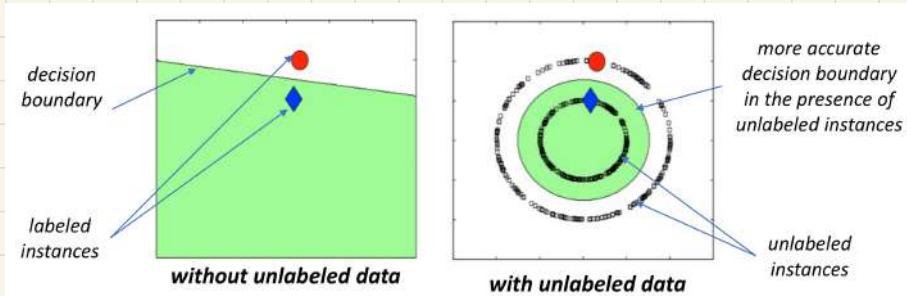
supervised learning is inductive.

• general idea: **Semi supervised learning**

definition: given labelled data (T) and unlabeled data (V)

main idea: use both T and V to learn mapping!

How does semi supervised learning help?



(Semi supervised) Deep Learning on Graphs:

Differential message passing framework

given:

- $G = (V, E)$ nodes V , edges E , $E_{uv} \rightarrow$ edge weights
- each node v has features $x_v \in \mathbb{R}^d$
- subset of labelled nodes $S \subseteq V$, $y_v \rightarrow$ labels of node $v \in S$

• $h_v^{(k-1)}$ hidden representation of node v in layer $(k-1)$

• For each node do:

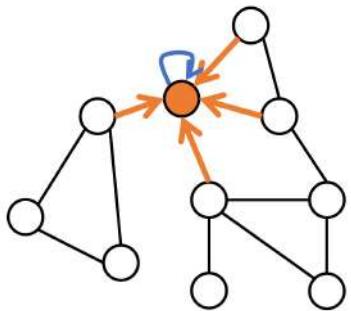
1) gather message from all neighbours

$$m_v^{(k)} = \sum_{u \in N(v)} M(h_v^{(k-1)}, h_u^{(k-1)}, E_{vu})$$

2) update the hidden representation

$$h_v^{(k)} = U(h_v^{(k-1)}, m_v^{(k)})$$

M & U are differentiable fn \rightarrow Neural Nets.



1. Gather message from all neighbors
2. Update hidden representation

instantiation :

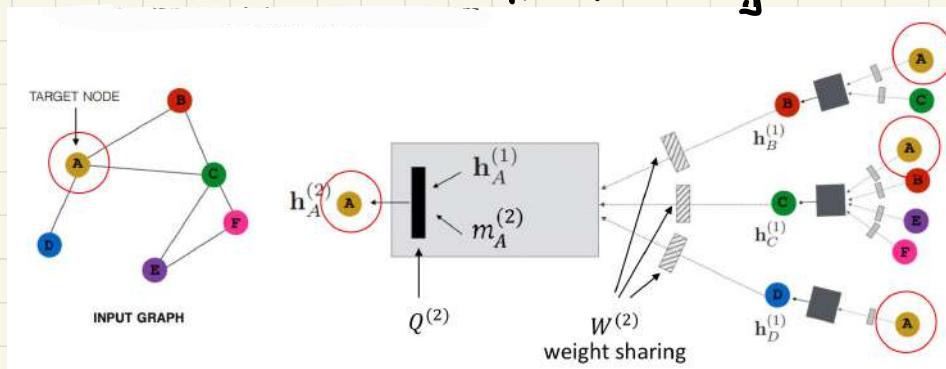
- let $h_v^{(0)} = x_v \rightarrow$ at the first layer the representation are node features
- set of message aggregate functions to be the average over the neighbors

$$m_v^{(k)} = \sum_{u \in \text{neighbors}(v)} \frac{1}{d_v} (W^{(k)} h_u^{(k-1)} + b^{(k)})$$

- calculate hidden representations as simple NNs with a non-linearity

$$h_v^{(k)} = \text{relu}(\theta^{(k)}) h_v^{(k-1)} + p^{(k)} + m_v^{(k)}$$

$W^{(k)}, b^{(k)}, \theta^{(k)}, p^{(k)}$ are trainable parameters of the k -th layer.



How do Neighbours Influence a Given Node

Observation 1: K hidden layers equal to K steps of message passing

- at step 1, node v aggregates info of 1-hop neighbour
 - at step 2, node v (implicitly) aggregates info from neighbour of 1-hop neighbour = 2-hop
- ↓
at step K : implicitly K -hop

Observation 2: K hidden layers \rightarrow representation $h_v^{(K)}$ has info-based on K -hop neighbours.

How to perform Semi-supervised Node Classification

- $h_v^{(K)}$ at final layer $\rightarrow K$ logits
- $p_v = \text{softmax}(h_v^{(K)})$ $\xrightarrow{\text{probability that node } v \text{ belongs to class } c}$ probability of node

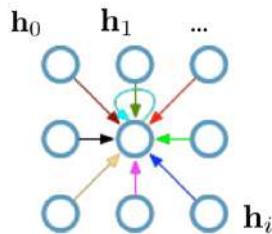
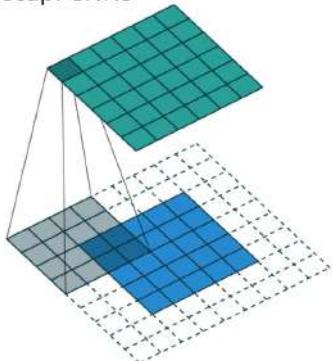
$$\min_{\{w^{(k)}, b^{(k)}, \mu_k, \nu_k\}} - \sum_{v \in S} \sum_{c \in C} y_{vc} \log p_{vc}^{\text{one hot encoded}}$$

\uparrow one hot encoded set of classes

$\min_{\{w^{(k)}, b^{(k)}, \mu_k, \nu_k\}}$
 $\forall k=1 \dots K$
 labelled nodes

y_{vc} \uparrow $v \in S$
 $c \in C$ \uparrow $v \in S$
 $\log p_{vc}$ \uparrow $v \in S$
 \uparrow $c \in C$

- Alternative view of message passing framework is that we are essentially performing a “graph convolution”
- Recap: CNNs

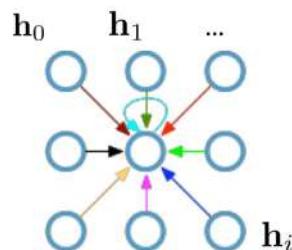
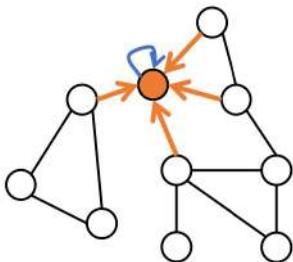


Update for a single pixel:

- Transform messages individually $\mathbf{W}_i \mathbf{h}_i$
- Add everything up $\sum_i \mathbf{W}_i \mathbf{h}_i$

- Images are a special kind of graph: every pixel is a node connected to 8 other nodes (up, down, left, right, etc. pixels)

- Note: Unlike sequences/images, where the convolution operation is clearly defined, there are various versions for graphs. Indeed, one often distinguishes between spatial and spectral approaches.



In spatial domain, graph conv. updates each node's features by considering the local neighbourhood.

- In effect it is some kind of message passing

$$\mathbf{x}_i^{(t+1)} = \mathbf{x}_i^{(t)} + \sum_{(i,j) \in E} \mathbf{x}_j^{(t)}$$

$$H^{(l+1)} = \sigma\left(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} H^{(l)} W^{(l)}\right)$$

● Non-linearity ● Message Passing ● Feature Transformation

Normalizing propagation matrix avoids gradient explosion
 $\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$ where $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}_n$ $D_{ii} = \sum_j A_{ij}$ is degree matrix \mathbf{A}

How deep are GNN?

furthest distance that graph propagates info to

$$H^{L+1} = \sigma(\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{1/2} H^L W)$$

GNNs apply 1 transformation per message passing step

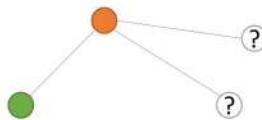
Transformation depth and propagation depth are orthogonal in GNNs.

- We could transform the messages with a multi-layer network in each step or propagate without transforming for multiple steps

Multigraph

- Given:
 - A set of graphs $\mathcal{G} = \{G_i = (V_i, E_i)\}_{i=1..N}$
 - A subset of labeled graphs $\mathcal{H} \subseteq \mathcal{G}$, with y_{G_i} denotes the graph level target of the graph $G_i \in \mathcal{H}$
- Obtain the hidden representation of the last layer for all nodes in all graphs according to the message passing framework
 - Let $H_{G_i} = [h_1^{(K)}, h_2^{(K)}, \dots, h_{|V_i|}^{(K)}]$ be a matrix where we stacked the final representations **for all nodes** for graph G_i
- Define an aggregation function $R(H_{G_i})$ over the node representations of a given graph that produces a representation for the **entire** graph $h_{G_i} = R(H_{G_i})$
- For classification: consider h_{G_i} as the logits and train a model using standard cross-entropy loss, i.e. $\mathcal{L}(y_{G_i}, \text{softmax}(h_{G_i}))$
 - For regression, e.g. predicting functional properties of molecules you can use squared loss
- Example agg. function $R(H_{G_i}) = \frac{1}{|V_i|} \sum_{v \in V_i} h_v^{(K)}$ is the average of the node embedding

- Consider the graph below. What is the influence of the green node on the unlabeled nodes in Label Propagation? Why? How about GNNs with $K = 2$?



- LP does not care about connections b/w labelled nodes.
 - In GNN., 2 hop influence *
- Does semi-supervised learning exist outside of learning on graphs?
yes on vector data as well

- Can you apply GNNs to vector data without a specified graph structure?

GNN is then equivalent to MLP which classifies nodes based on own features

Problem 1: The goal in Label Propagation is to find a labeling $\mathbf{y} \in \{0, 1\}^N$ that minimizes the energy $\min_{\mathbf{y}} \frac{1}{2} \sum_{ij} \mathbf{w}_{ij} (y_i - \hat{y}_j)^2$ subject to $y_i = \hat{y}_i \forall i \in S$ where the set of nodes V has been partitioned into the labeled nodes S and the unlabeled nodes U , $w_{ij} \geq 0$ is the non-negative edge weight and \hat{y}_i are the observed labels.

Following from the first observation regarding the Laplacian, the minimization problem can be rewritten and then relaxed to $\min_{\mathbf{y} \in \mathbb{R}^N} \mathbf{y}^T \mathbf{L} \mathbf{y}$ subject to the same constraints. Show that the closed form solution is

$$\mathbf{y}_U = -\mathbf{L}_{UU}^{-1} \cdot \mathbf{L}_{US} \cdot \hat{\mathbf{y}}_S$$

where w.l.o.g. we assume that the Laplacian matrix is partitioned into blocks for labeled and unlabeled nodes as

$$\mathbf{L} = \begin{pmatrix} \mathbf{L}_{SS} & \mathbf{L}_{SU} \\ \mathbf{L}_{US} & \mathbf{L}_{UU} \end{pmatrix}.$$

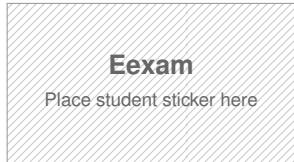
$$\begin{aligned} \text{P1)} \quad \mathbf{y}^T \mathbf{L} \mathbf{y} &= \mathbf{y}^T \begin{pmatrix} \mathbf{L}_{SS} & \mathbf{L}_{SU} \\ \mathbf{L}_{US} & \mathbf{L}_{UU} \end{pmatrix} \mathbf{y} \\ &= \hat{\mathbf{y}}_S^T \mathbf{L}_{SS} \hat{\mathbf{y}}_S + \hat{\mathbf{y}}_S^T \mathbf{L}_{SU} \mathbf{y}_U + \mathbf{y}_U^T \mathbf{L}_{US} \hat{\mathbf{y}}_S \\ &\quad + \mathbf{y}_U^T \mathbf{L}_{UU} \mathbf{y}_U \end{aligned}$$

Laplacian is symmetric

$$\therefore \hat{\mathbf{y}}_S^T \mathbf{L}_{SS} \hat{\mathbf{y}}_S + 2\mathbf{y}_U^T \mathbf{L}_{US} \hat{\mathbf{y}}_S + \mathbf{y}_U^T \mathbf{L}_{UU} \mathbf{y}_U = f(\mathbf{y}_U)$$

$$\begin{aligned} \frac{\partial f}{\partial \mathbf{y}_U} &= 2\mathbf{L}_{US} \hat{\mathbf{y}}_S + (\mathbf{L}_{UU} + \mathbf{L}_{UU}^T) \mathbf{y}_U = 2\mathbf{L}_{US} \hat{\mathbf{y}}_S + \\ &\quad 2\mathbf{L}_{UU} \mathbf{y}_U = 0 \end{aligned}$$

$$\mathbf{y}_U = -\mathbf{L}_{UU}^{-1} \mathbf{L}_{US} \hat{\mathbf{y}}_S$$



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: Wednesday 5th August, 2020

Examiner: Prof. Dr. Stephan Günnemann

Time: 11:30 – 12:45

Working instructions

- This exam consists of **14 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 43 credits.
- Detaching pages from the exam is prohibited.
- Allowed resources:
 - all materials that you will use on your own (lecture slides, calculator etc.)
 - **not allowed are any forms of collaboration between examinees and plagiarism**
- You have to sign the code of conduct.
- Make sure that the **QR codes are visible** on every uploaded page. Otherwise, we cannot grade your exam.
- Only write on the provided sheets, **submitting your own additional sheets is not possible**.
- Last two pages can be used as scratch paper.
- All sheets (including scratch paper) have to be submitted to the upload queue. Missing pages will be considered empty.
- **Only use a black or blue color (no red or green)!**
- Write your answers only in the provided solution boxes or the scratch paper.
- **For problems that say "Justify your answer" you only get points if you provide a valid explanation.**
- **For problems that say "Prove" you only get points if you provide a valid mathematical proof.**
- If a problem does not say "Justify your answer" or "Prove" it's sufficient to only provide the correct answer.
- Exam duration - 75 minutes.

Left room from _____ to _____ / Early submission at _____

Problem 1 Normalizing Flows (4 credits)

We consider two transformations $f_1(\mathbf{z}) = \begin{bmatrix} z_1 \\ z_2^{1/3} \end{bmatrix}$ and $f_2(\mathbf{z}) = \begin{bmatrix} z_1(|z_2| + 1) \\ z_2 \end{bmatrix}$ from \mathbb{R}^2 to \mathbb{R}^2 .

The respective inverse transformation are $f_1^{-1}(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2^3 \end{bmatrix}$ and $f_2^{-1}(\mathbf{x}) = \begin{bmatrix} \frac{x_1}{|x_2|+1} \\ x_2 \end{bmatrix}$.

The respective Jacobians are

$$J_{f_1} = \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{3}z_2^{-\frac{2}{3}} \end{bmatrix} \quad J_{f_2} = \begin{bmatrix} |z_2| + 1 & \text{sign}(z_2)z_1 \\ 0 & 1 \end{bmatrix}$$

$$J_{f_1^{-1}} = \begin{bmatrix} 1 & 0 \\ 0 & 3x_2^2 \end{bmatrix} \quad J_{f_2^{-1}} = \begin{bmatrix} \frac{1}{|x_2|+1} & \frac{-\text{sign}(x_2)x_1}{(|x_2|+1)^2} \\ 0 & 1 \end{bmatrix}$$

- 0 We assume a Gaussian base distribution $p_1(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$. We observed one point $\mathbf{x}_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$. info
 1 We propose to stack the transformations f_1, f_2 to transform the base distribution p_1 into the distribution p_2 with
 2 normalizing flows. Compute the likelihood for \mathbf{x} under the transformed distribution p_2 if the order of transformations
 3 is f_1 followed by f_2 .
 4 Hint: You might use the density of the unit variate Gaussian $p = \mathcal{N}(0, 1)$ at the following points: $p(1/2) = 0.3521$,
 $p(1/3) = 0.3774$, $p(1/9) = 0.3965$, $p(5) = 1.4867e^{-6}$, $p(8) = 5.0523e^{-15}$, $p(10) = 7.6946e^{-23}$

$$f_2(f_1(\mathbf{z})) = \mathbf{x} \quad \mathbf{x} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$\mathbf{z} = f_1^{-1}(f_2^{-1}(\mathbf{x})) \quad \text{--- ii)}$$

$$f_2^{-1}(\mathbf{x}) : \begin{bmatrix} 1/3 \\ 2 \end{bmatrix} \quad f_1^{-1} \circ f_2^{-1}(\mathbf{x}) : \begin{bmatrix} 1/3 \\ 8 \end{bmatrix}$$

$$p_2(\mathbf{x}) = p_1 \left[\begin{bmatrix} 1/3 \\ 8 \end{bmatrix} \right] \times \left(\frac{1}{12+1} \right) \times (3 \times 4)$$

$$= p_1 \left[\begin{bmatrix} 1/3 \\ 8 \end{bmatrix} \right] \times 4 = 7.6266 e^{-15}$$

~~Problem 2~~ Variational Inference (5 credits)

We are performing variational inference in some latent variable model $p_\theta(x, z)$ using the following family of variational distributions $\mathcal{Q}_1 = \{\mathcal{N}(z|\phi, 1) : \phi \in \mathbb{R}\}$.

- a) Assume that the variational distribution $q \in \mathcal{Q}_1$ is fixed, and we are trying to maximize the ELBO w.r.t. θ using gradient ascent. Is it necessary to use the reparametrization trick in this case? If yes, explain how to do it for our family of distributions \mathcal{Q}_1 ; if not, provide a justification.

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

No, we don't need the reparametrization trick for calculating ELBO when performing gradient descent with θ .

We can estimate

$$\begin{aligned} \mathbb{E}_{z \sim q_{\phi}(z)} [f_{\theta}(z)] &\approx \frac{1}{S} \sum_{i=1}^S f_{\theta}(z_i) \\ \nabla \mathbb{E}_{z \sim q_{\phi}(z)} [f_{\theta}(z)] &= \mathbb{E}_{z \sim q_{\phi}(z)} [\nabla f_{\theta}(z)] \\ &\approx \frac{1}{S} \sum_{i=1}^S \nabla_{\theta} f_{\theta}(z_i) \end{aligned}$$

- 4) Consider another family of distributions $\mathcal{Q}_2 = \{\mathcal{N}(z|0, s^2) : s \in (0, \infty)\}$. Which of the following statements is true?
Justify your answer.

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

1. $\max_{\theta, q \in \mathcal{Q}_1} \text{ELBO}(\theta, q) < \max_{\theta, q \in \mathcal{Q}_2} \text{ELBO}(\theta, q)$
2. $\max_{\theta, q \in \mathcal{Q}_1} \text{ELBO}(\theta, q) = \max_{\theta, q \in \mathcal{Q}_2} \text{ELBO}(\theta, q)$
3. $\max_{\theta, q \in \mathcal{Q}_1} \text{ELBO}(\theta, q) > \max_{\theta, q \in \mathcal{Q}_2} \text{ELBO}(\theta, q)$
4. It's impossible to tell without additional information.

- impossible to tell without any knowledge of $p(x, z)$
- Depends on the model.

Problem 3 Robustness of Machine Learning Models (6 credits)

Suppose we have trained a binary classifier $f : \mathbb{R}^d \rightarrow \{0, 1\}$ and want to certify its robustness via randomized smoothing. Therefore, the *smoothed classifier* $g_{\sigma^2}(\mathbf{x}) = \mathbb{E}[\mathbb{I}[f(\mathbf{x} + \varepsilon) = 1]]$, where $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$.

Fact: $\Phi^{-1}(g_{\sigma^2}(\mathbf{x}))$ is $1/\sigma$ -**Lipschitz** w.r.t. \mathbf{x} and the L_2 norm, where $\Phi(z)$ denotes the cumulative distribution function (CDF) of the standard normal distribution.

- 0 a) Using the above fact about the Lipschitz-continuity of $\Phi^{-1}(g_{\sigma^2}(\mathbf{x}))$, show that the largest certifiable L_2 radius r around a sample \mathbf{x} is identical to the result shown in the lecture. More precisely, show that

$$r = \sigma \Phi^{-1}(g_{\sigma^2}(\mathbf{x})).$$

Hint: You may assume we can evaluate $g_{\sigma^2}(\mathbf{x})$ in closed form and you may use the following results: $\lim_{z \rightarrow 0} \Phi^{-1}(z) = -\infty$, $\Phi^{-1}(0.5) = 0$, $\lim_{z \rightarrow 1} \Phi^{-1}(z) = \infty$.

$$\|\Phi^{-1}(g_{\sigma^2}(\mathbf{x})) - \Phi^{-1}(0.5)\|_2 \leq \frac{1}{\sigma} \|\mathbf{x} - \tilde{\mathbf{x}}\|_2$$

$$\Phi^{-1}(g_{\sigma^2}(\mathbf{x})) \leq \frac{1}{\sigma} \|\mathbf{x} - \tilde{\mathbf{x}}\|^2$$

$$\|\mathbf{x} - \tilde{\mathbf{x}}\|^2 \geq -\Phi^{-1}(g_{\sigma^2}(\mathbf{x}))$$

- 0 b) A fellow student has a promising idea: by letting $\sigma \rightarrow \infty$ we can make the Lipschitz constant of the smoothed classifier arbitrarily small, leading to arbitrarily large certifiable radii, i.e. a very robust model. Is this a good idea? Why or why not?

- ① if σ goes to infinity then the variance of the noise is arbitrarily large; and signal / noise ≈ 0 ; we lose all our info
- ② for large σ we sample over really huge space for which large fraction of samples belong to other states

Problem 4 Markov Property (3 credits)

We consider the following sequences of random variables U_0, U_1, \dots, U_t .

- a) $U_t = \begin{bmatrix} X_t \\ Z_t \end{bmatrix}$ where X_t are observed variables and Z_t are latent variables of an Hidden Markov Model. Does the sequence of variables U_t fulfill the Markov property i.e. $P(U_t|U_{t-1}) = P(U_t|U_{t-1}, \dots, U_0)$? Justify your answer.

0
 1

$$P(V_t | U_{t-1}, \dots, U_1, V_0) = P\left(\begin{bmatrix} X_t \\ Z_t \end{bmatrix} | P\left(\begin{bmatrix} X_{t-1} \\ Z_{t-1} \end{bmatrix}, \dots, \begin{bmatrix} X_0 \\ Z_0 \end{bmatrix}\right)\right)$$

$$P(X_t, Z_t | X_{t-1}, Z_{t-1}, \dots, X_0, Z_0) = P(X_t, Z_t | X_{t-1}, Z_{t-1})$$

all info is retained, hence this is true

- b) We consider an AR(p) process X_t . Under what condition on p and k does the sequence of variables $U_t = [X_{t-1}, \dots, X_{t-k}]$ fulfill the Markov property i.e. $P(U_t|U_{t-1}) = P(U_t|U_{t-1}, \dots, U_0)$? Justify your answer.

0
 1

- c) We consider a recurrent neural network which produces X_t . Does the sequence of variables $U_t = X_t$ fulfill the Markov property i.e. $P(U_t|U_{t-1}) = P(U_t|U_{t-1}, \dots, U_0)$? Justify your answer.

0
 1

Recurrent neural networks have all the past info encoded in the hidden state h_t and both x_t and x_{t-2} depends on h_{t-2} hence not markov

Problem 5 Markov Chain (3 credits)

0
1
2
3

We consider a Markov chain X_t in $\{1, C\}$ with parameters π, \mathbf{A} . We assume we observed the sequence $S_k = [\underbrace{v_0, \dots, v_0}_{k \text{ times}}, \underbrace{v_1, \dots, v_1}_{k \text{ times}}, \dots, \underbrace{v_T, \dots, v_T}_{k \text{ times}}]$ where each value is observed k times. The parameter k can be seen as a discretization parameter of the time space.

Compute the likelihood of the sequence under the parameters π, \mathbf{A} i.e. $P_{\pi, \mathbf{A}}(S_k)$. What happens to this quantity if you increase the discretization parameter from k to $k' > k$ but keep the same model parameter π, \mathbf{A} ?

$$P_{\pi, \mathbf{A}}(S_k) = \pi^{\#_0} \times \prod_{t=0}^{T-1} \mathbf{A}^{\#_{t+1}, v_{t+1}} \sum_{t=0}^{k-1} \mathbf{A}^{\#_t, v_t}$$

Since $A_{i,j} < 1 \rightarrow \text{likelihood decreases as } k \text{ increases}$

Problem 6 Temporal Point Process (6 credits)

Consider an inhomogeneous Poisson process (IPP) on the interval $[0, 4]$ with the intensity function

$$\lambda(t) = \begin{cases} a & \text{if } t \in [0, 3] \\ b & \text{if } t \in (3, 4] \end{cases}$$

where $a > 0, b > 0$ are some positive parameters.

- a) Assume that you observed a sequence of events $\{0.2, 1.0, 1.5, 2.9, 3.1, 3.8\}$ generated by the above IPP. What is the maximum likelihood estimate of the parameters a and b ?

0
1
2
3
4

$$\begin{aligned} \log P(\{t_1, \dots, t_n\}) &= \sum_{i=1}^n \log \lambda(t_i) - \int_0^T \lambda(t) dt \\ &= 4 \log(a) + 2 \log(b) - (3a + b). \\ (\text{LL}) &= 4 \log(a) + 2 \log(b) - 3a - b. \\ \frac{\partial \text{LL}}{\partial a} &= \frac{4}{a} - 2 \Rightarrow a^* = \frac{4}{2}. \\ \frac{\partial \text{LL}}{\partial b} &= \frac{2}{b} - 1 \Rightarrow b^* = 2. \end{aligned}$$

- b) Assume that $a = 1$ and $b = 5$. What is the expected number of events generated by the IPP in this case?

0
1
2

$$\begin{aligned} \text{Expected number of events} &: \int_0^T \lambda^*(t) dt \\ &= \int_0^3 a dt + \int_3^4 b dt \\ &= 3 \cdot 1 + 5 = 3 + 5 = 8 \end{aligned}$$

Problem 7 Clustering with the Planted Partition Model (4 credits)

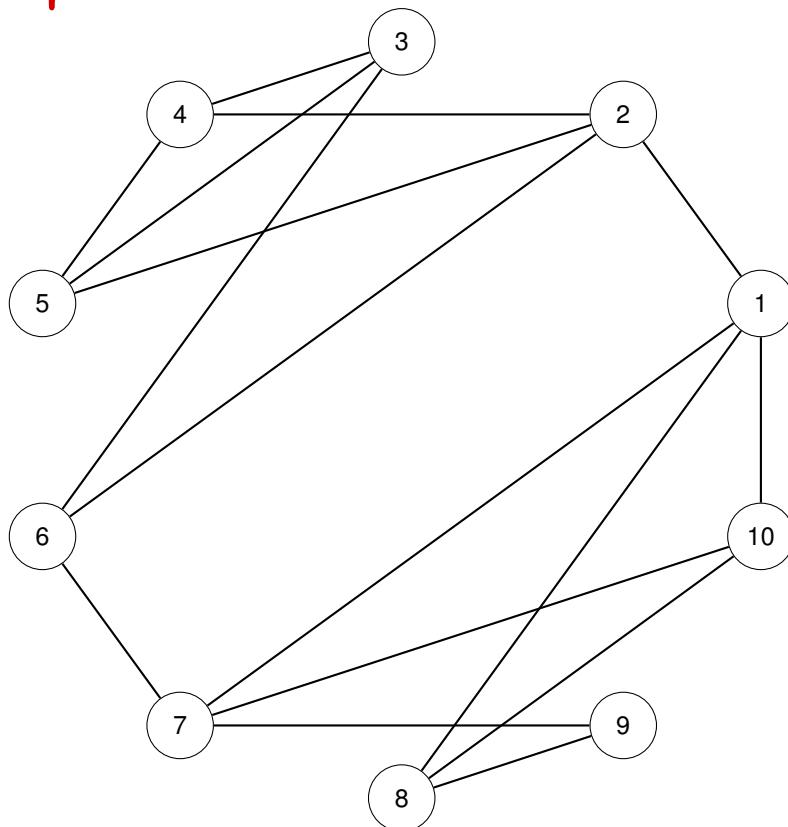
0
1
2
3
4

The following graph has been generated from a planted partition model with in-community edge probability p and between-community edge probability q .

PPM

P

Doubt



Ans

Assuming $p < q$, find the maximum likelihood community assignments under a PPM.

Give your solution as two sets of node labels making up the two discovered communities. Justify your answer.

$$\text{Likelihood} \propto \frac{q^{|\text{out}|}}{p^{|\text{in}|}} = \left(\frac{q}{p}\right)^{|\text{out}|} \left(\frac{1-p}{1-q}\right)^{|\text{in}|}$$

for $p < q$ Likelihood > 1

\therefore likelihood is maximized by max out size

~~Problem 8~~ PageRank in a Wheel (6 credits)

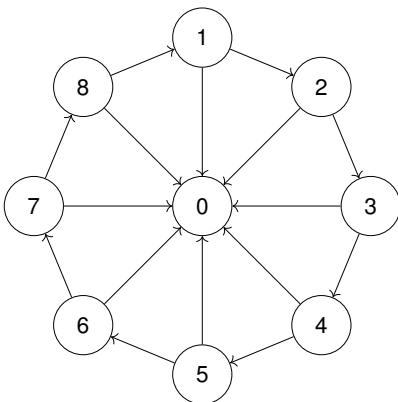


Figure 8.1: Example of a directed wheel graph with $n + 1 = 9$ nodes

Consider a directed graph of size $n + 1$ with a cycle of n nodes and an additional central node that every other node connects to (see figure). So we have a graph with node set $\mathcal{V} = \{0, 1, \dots, n\}$ and edge set

$$\mathcal{E} = \{(i, i+1) \mid i \in \{1, \dots, n-1\}\} \cup \{(n, 1)\} \cup \{(i, 0) \mid i \in \{1, \dots, n\}\}.$$

We want to compute the PageRank scores with a link-follow probability of β (a teleport probability of $1 - \beta$) and some arbitrary teleport vector π , $\sum_{i=0}^n \pi_i = 1$. Note that we index π from 0 to n .

We define the predecessor function pa as the index of the predecessor of a node in the directed cycle, i.e.

$$\text{pa}(1) = n \quad \text{and} \quad \text{pa}(i) = i - 1 \quad \forall i \in \{2, \dots, n\}.$$

You can write $\text{pa}^k(i)$ for the k -th predecessor of node i , i.e. $\text{pa}^3(i) = \text{pa}(\text{pa}(\text{pa}(i)))$ and $\text{pa}^0(i) = i$.

a) Set up the PageRank equations for all nodes in scalar form, i.e. each r_i separately instead of matrix form.

node 1, 2, .. n have two outgoing edges

$$r_i = \beta \frac{\pi_{\text{pa}(i)}}{2} + (1-\beta) r_i^* \quad \forall i \in \{1, \dots, n\}$$

central node only has incoming nodes from all the outer nodes

$$r_0 = \beta \sum_{i=1}^n \left(\frac{r_i}{2} \right) + (1-\beta) \pi_0$$

b) Why is this graph problematic for PageRank without random teleportation ($\beta = 1$)?

dead end at center. Random walk will be stuck there

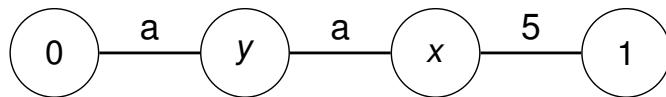
c) Show that the PageRank for node $i \in \{1, \dots, n\}$ in the outer cycle is given by

$$r_i = \frac{(1-\beta)}{1 - \left(\frac{\beta}{2}\right)^n} \sum_{j=0}^{n-1} \left(\frac{\beta}{2}\right)^j \pi_{\text{pa}^j(i)}.$$

$$\begin{aligned}
 r_i &= \frac{\beta}{2} \pi_{\text{pa}(i)} + (1-\beta) \pi_i^0 \\
 &= \frac{\beta}{2} \left(\frac{\beta}{2} \pi_{\text{pa}(\text{pa}(i))} + (1-\beta) \pi_{\text{pa}(i)} \right) + (1-\beta) \pi_i^0 \\
 &= \left(\frac{\beta}{2}\right)^2 \pi_{\text{pa}(\text{pa}(i))} + \frac{\beta}{2} (1-\beta) \pi_{\text{pa}(i)} + (1-\beta) \pi_i^0 \\
 &= \left(\frac{\beta}{2}\right)^2 \left[\frac{\beta}{2} \pi_{\text{pa}(\text{pa}(\text{pa}(i)))} + (1-\beta) \pi_{\text{pa}(\text{pa}(i))} \right] \\
 &\quad + \frac{\beta}{2} (1-\beta) \pi_{\text{pa}(i)} + (1-\beta) \pi_i^0 \\
 &= \left(\frac{\beta}{2}\right)^3 \pi_{\text{pa}(\text{pa}(\text{pa}(i)))} + \left(\frac{\beta}{2}\right)^2 (1-\beta) \pi_{\text{pa}(\text{pa}(i))} \\
 &\quad + \frac{\beta}{2} (1-\beta) \pi_{\text{pa}(i)} + (1-\beta) \pi_i^0 \\
 &\quad + \dots \\
 \pi_i^0 &= \left(\frac{\beta}{2}\right)^n \pi_{\text{pa}(i)} + (1-\beta) \sum_{j=0}^{n-1} \left(\frac{\beta}{2}\right)^j \pi_{\text{pa}^j(i)}
 \end{aligned}$$

Problem 9 Label Propagation (4 credits)

Consider the following graph



	0
	1
	2
	3
	4

The nodes labeled 0 and 1 are observed and from class 0 and 1, respectively. One edge has a fixed weight, the other two have a variable edge weight of $a \geq 0$. The two center nodes are unobserved and we call their labels x and y .

We want to predict classes for the two center nodes that minimize the Label Propagation objective exactly,

$$\frac{1}{2} \sum_{ij} W_{ij} (y_i - y_j)^2$$

where W is the weighted adjacency matrix and y_i, y_j are the labels of the nodes.

Find the set of all possible edge weights a that guarantee that node x is assigned to class 0. Justify your answer.

$$5(x-1)^2 + a^2(x-y)^2 + a(y-0)^2 = (a+5)x^2 - 10x + 5 - 2axy + 2ay^2 = E(a, x, y)$$

x	y	$E(a, 0, 0)$
0	0	5
1	0	a
0	1	$5 + 2a$
1	1	a

We set
 $x \leq 0$ when

$$a \geq 5$$

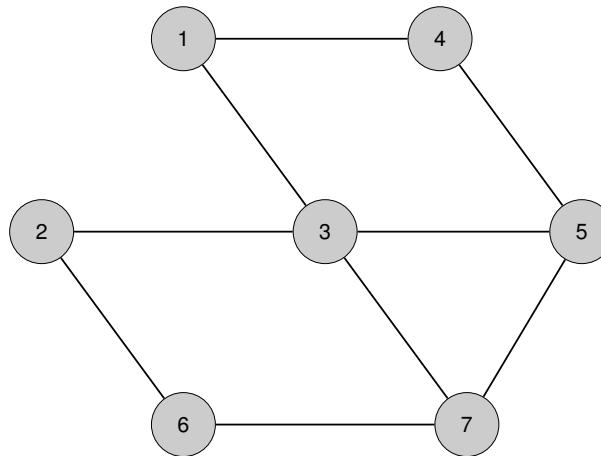
Problem 10 Adversarial Attacks on Graph Neural Networks (2 credits)

Suppose you are given the following two-layer graph neural network.

$$f(\mathbf{A}, \mathbf{X}) = \mathbf{Z} = \text{Softmax} (\hat{\mathbf{A}} \text{ReLU} (\hat{\mathbf{A}} \mathbf{X} \mathbf{W}_1) \mathbf{W}_2)$$

$\mathbf{X} \in \mathbb{R}^{N \times D}$ are the node features, \mathbf{Z} are the node predictions, \mathbf{W}_x are weight matrices of appropriate dimensions and $\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$ is the propagation matrix as defined for GCNs. Here, $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}$, where \mathbf{A} is the adjacency matrix and \mathbf{I} is the identity matrix, and $\tilde{\mathbf{D}}$ is a diagonal matrix of node degrees $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$.

The model was trained for the task of semi-supervised node classification, and we want to predict a class c for node 6 in the following graph \mathbf{A} :



- 0 1 a) An adversary with complete knowledge about the graph \mathbf{A} and the trained model $f(\mathbf{A}, \mathbf{X})$ may delete one edge to perturb the prediction for node 6. Deleting which of the following edges would lead to a greater change to the prediction for node 6? Justify your answer.

1. The edge connecting node 5 and 7
2. The edge connecting node 3 and 5
3. There is not enough information to determine which deletion leads to a greater change.

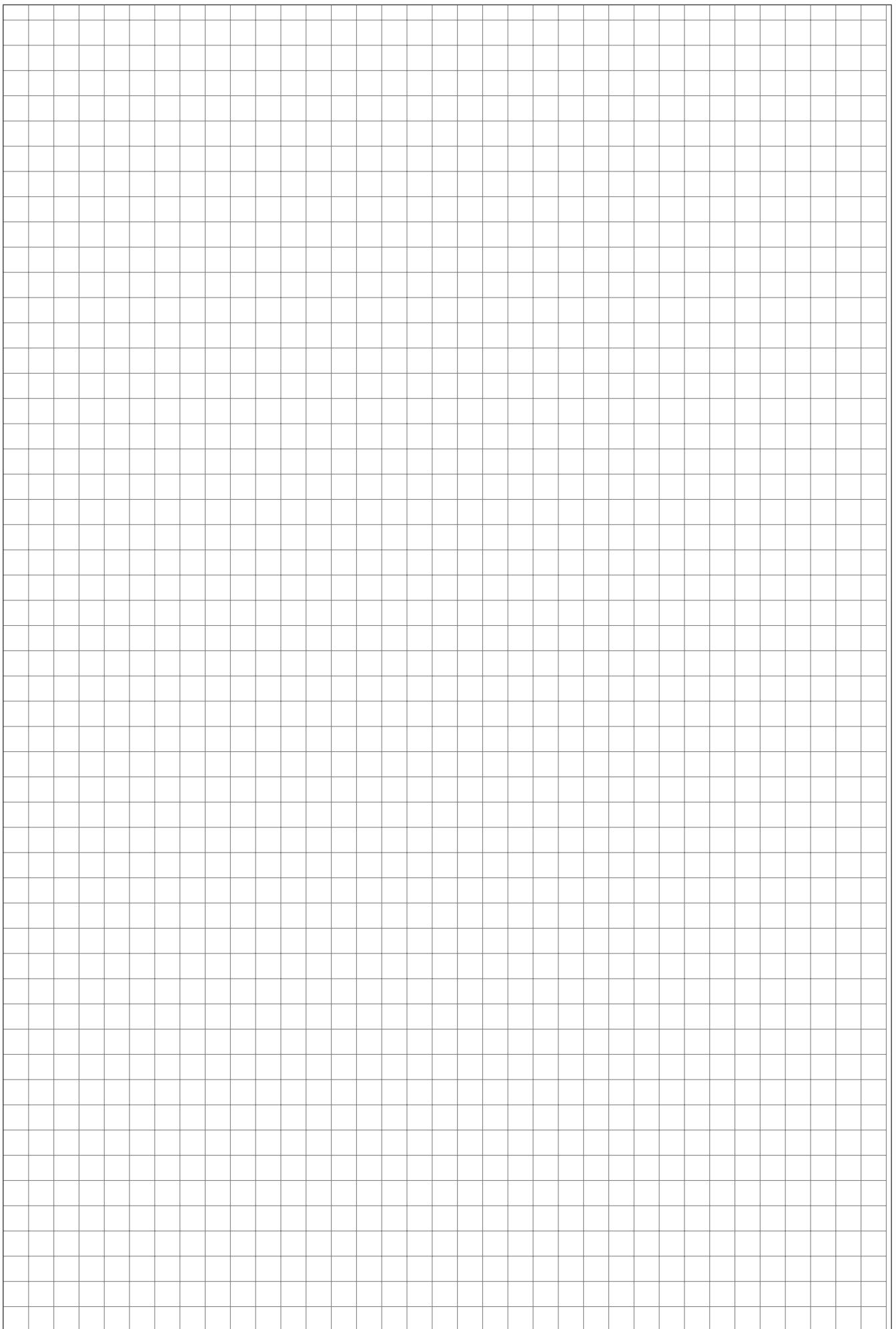
5 → 7 as it is two hops away
3 → 5 won't affect

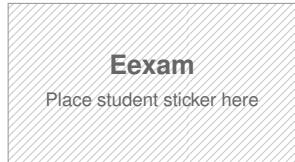
- 0 1 b) Assume we instead have a Personalized Propagation of Neural Predictions (PPNP) model instead of the two-layer GCN. How does this affect your choice? Justify your answer.

PPNP is infinite propagation with teleport
 not enough info

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

A large grid of squares, approximately 20 columns by 30 rows, intended for students to write their solutions. The grid is composed of thin black lines on a white background.





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 / Retake

Date: Wednesday 7th October, 2020

Examiner: Prof. Dr. Stephan Günnemann

Time: 14:15 – 15:30

Working instructions

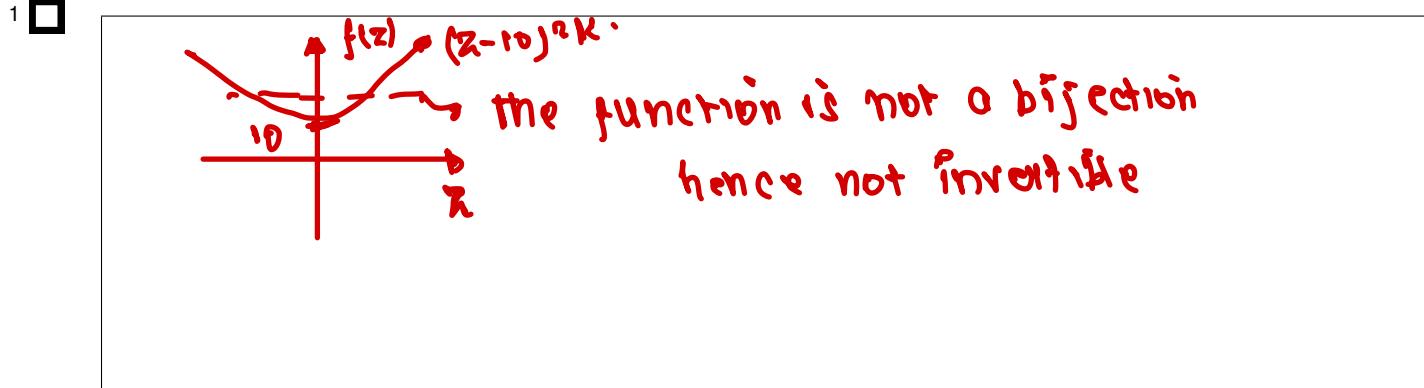
- This exam consists of **14 pages** with a total of **10 problems**.
Please make sure now that you received a complete copy of the exam.
- The total amount of achievable credits in this exam is 38 credits.
- Detaching pages from the exam is prohibited.
- Allowed resources:
 - all materials that you will use on your own (lecture slides, calculator etc.)
 - **not allowed are any forms of collaboration between examinees and plagiarism**
- You have to sign the code of conduct.
- Make sure that the **QR codes are visible** on every uploaded page. Otherwise, we cannot grade your exam.
- Only write on the provided sheets, **submitting your own additional sheets is not possible**.
- Last two pages can be used as scratch paper.
- All sheets (including scratch paper) have to be submitted to the upload queue. Missing pages will be considered empty.
- **Only use a black or blue color (no red or green)!**
- Write your answers only in the provided solution boxes or the scratch paper.
- **For problems that say "Justify your answer" you only get points if you provide a valid explanation.**
- **For problems that say "Prove" you only get points if you provide a valid mathematical proof.**
- If a problem does not say "Justify your answer" or "Prove" it's sufficient to only provide the correct answer.
- Instructor announcements and clarifications will be posted **on Piazza** with email notifications
- Exam duration - 75 minutes.

Left room from _____ to _____ / Early submission at _____



Problem 1 Normalizing Flows (5 credits)

- 0 a) Let $k \in \mathbb{N}$. We consider the transformation $f(z) = (z - 10)^{2k}$ from \mathbb{R} to \mathbb{R} . Is this transformation invertible? Justify your answer. If yes, compute the determinant of its Jacobian and its inverse.



- 0 b) We consider the transformation $f(z) = \begin{bmatrix} z_n \\ z_1 \\ \vdots \\ z_{n-1} \end{bmatrix}$ from \mathbb{R}^n to \mathbb{R}^n . Is this transformation invertible? Justify your answer. If yes, compute the determinant of its Jacobian and its inverse. Hint: The determinant of one elementary permutation (i.e. a permutation which interchanges any two rows) is -1

Yes, it can be seen as $(n-1)$ stacked permutations
determinant $(-1)^{n-1}$
inverse $f^{-1}(z) = \begin{bmatrix} z_2 \\ z_3 \\ \vdots \\ z_1 \end{bmatrix}$

how is this
 $(n-1)$ stacked
permutation?

- 0 c) We consider the transformation $f(z) = \begin{bmatrix} z_1 + 2z_2 \\ 3z_1 + 4z_2 \\ z_3 - 2z_4 \\ 3z_3 - 4z_4 \end{bmatrix}$ from \mathbb{R}^4 to \mathbb{R}^4 . Is this transformation invertible? Justify your answer. If yes, compute the determinant of its Jacobian and its inverse.

$$J = \begin{bmatrix} 1 & 2 & 0 & 0 \\ 3 & 4 & 0 & 0 \\ 0 & 0 & 1 & -2 \\ 0 & 0 & 3 & -4 \end{bmatrix}$$

$z_2 = x_1 \cdot x_2$

$z_1 + 2z_2 = x_1$

$3z_1 + 4z_2 = x_2$

$z_1 = 2x_2 - x_1$

$|J| = (-2)(-12) = 24$
as determinant $\neq 0$
invertible

~~Problem 2~~ Variational Inference (3 credits)

We are performing variational inference in some latent variable model $p_\theta(x, z)$ using the following family of variational distributions $\mathcal{Q}_1 = \{\mathcal{N}(z|\phi, 1) : \phi \in \mathbb{R}\}$.

Assume that

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

$$p(x|z) = \text{Bernoulli}(x | \sigma(z)) \propto \exp(xz - \log(1 + \exp(z)))$$

and x is known and fixed and $\sigma(z) = \frac{\exp(z)}{1+\exp(z)}$ is the sigmoid function. Does there exist a $q^* \in \mathcal{Q}_1$, such that $\text{KL}(q^*(z) \parallel p(z|x)) = 0$? Justify your answer.

0
1
2
3

for $\text{KL}(q^*(z) \parallel p(z|x)) = 0$

$$q^*(z) = p(z|x)$$

$$p(z|x) \propto p(x|z) \cdot p(z)$$

$$\propto \exp(xz - \log(1 + e^z)) \cdot \exp(-\frac{z^2}{2})$$

$$\propto \exp(xz - \log(1 + e^z) - \frac{z^2}{2})$$

This is not normal in z .

as $\log(1 + e^z) \rightarrow$ exponential term is present

~~DO~~

Problem 3 Robustness of Machine Learning Models (3 credits)

Consider a trained binary logistic regression model with weight vector $\mathbf{w} \in \mathbb{R}^d$ and bias $b \in \mathbb{R}$, where d is the data dimensionality. That is, the predicted probability of a sample $\mathbf{x} \in \mathbb{R}^d$ belonging to class 1 is:

$$p(y = 1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b),$$

where $\sigma(z) = \frac{1}{1+e^{-z}}$ is the logistic sigmoid function. An input sample is assigned to class 1 if $p(y = 1 | \mathbf{x}) > 0.5$, else it is assigned to class 0.

We would like to perform **robustness certification** of the logistic regression model using its **Lipschitz constant**. That is, we want to certify that the predicted class of the input sample does not change w.r.t. some radius in the input space.

- 0 a) Briefly explain why it is sufficient to consider $F(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, i.e. the model without the sigmoid activation function $\sigma(\cdot)$ for this purpose.

1 Because we just need to know which side of boundary is the input mapped to. (Classification)

- 0 b) Derive the (smallest possible) Lipschitz constant of $F(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ w.r.t. the L_2 norm. Justify your answer.

1 $\|\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \tilde{\mathbf{x}}\| \leq L \|\mathbf{x} - \tilde{\mathbf{x}}\|$
 let $a = \mathbf{x} - \tilde{\mathbf{x}}$
 $\|\mathbf{w}^T a\| \leq L \|a\|$
 $L \geq \left(\frac{\|\mathbf{w}^T a\|}{\|a\|} \right) \|\mathbf{w}\|$
 $\leq \|\mathbf{w}\|$
 (smallest value)

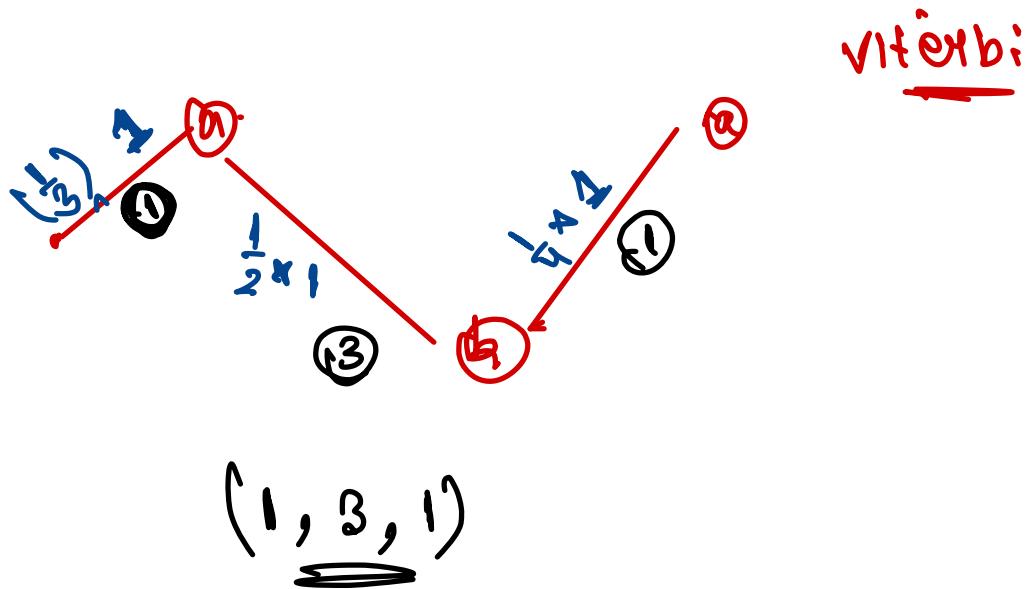
~~Problem 4~~ Hidden Markov Model (4 credits)

Consider an HMM where hidden variables are in $\{1, 2, 3\}$ and observed variables are in $\{a, b\}$. Let the model parameters be as follows:

$$\pi = \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix}, \quad A = \begin{bmatrix} 1 & 2 & 3 \\ 1 & \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \\ 2 & \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ 3 & \frac{1}{4} & \frac{1}{4} & \frac{1}{2} \end{bmatrix}, \quad B = \begin{bmatrix} a & b \\ 1 & \begin{bmatrix} 1 & 0 \\ \frac{2}{3} & \frac{1}{3} \end{bmatrix} \\ 2 & \begin{bmatrix} 0 & 1 \end{bmatrix} \\ 3 & \end{bmatrix}$$

Assume that the sequence $X_{1:3} = [aba]$ is observed. Find the most probable sequence $[Z_1, Z_2, Z_3]$ and compute its likelihood. Justify your answer.

0
1
2
3
4



~~Problem 5~~ Temporal Point Process (3 credits)

0
1
2
3

Consider a temporal point process defined on the interval $[0, 10\pi]$ with the conditional intensity function

$$\lambda^*(t) = \sin(t) + 2$$

What is the expected number of events that will be generated from this TPP on the interval $[0, 2\pi]$? Justify your answer.

$$\begin{aligned} \mathbb{E}[N[0, 2\pi]] &= \int_0^{2\pi} \lambda^*(t) dt = \int_0^{2\pi} (\sin(t) + 2) dt \\ &= -\cos t + 2t \Big|_0^{2\pi} = \boxed{4\pi} \end{aligned}$$

~~Problem~~ 6

Neural Network approaches for temporal data (4 credits)

We trained the following models to produce latent representations:

- M1: Skip-gram word2vec model
- M2: LSTM
- M3: Self-attention without positional encoding
- M4: Self-attention with positional encoding
- M5: A convolutional NN which produces the latent representation at time t based on the input at time t and $t - 1$.



At **inference time**, we use the following 6 sentences as inputs:

- S1: "I left"
- S2: "They left"
- S3: "I left yesterday"
- S4: "I go left"
- S5: "left I go"
- S6: "go left"

For each model (i.e. M1, M2, M3, M4, M5), which sets of sentences (e.g. (S1, S2), (S3, S4, S6)) are guaranteed to produce the same latent representation for the word "left"? Note that the answer may be zero, one or more sets of sentences. Justify your answer.

M1) skip gram : embeddings are the same, irrespective of the context

M2) LSTM : "I left" and "I left yesterday" produce same embedding \rightarrow left context

M3) self attention without positional encoding :
"I go left" and "left I go" are the same

M4) self attention with position encoding :
embeddings are different

M5) 1-D CNN with window 2

I go left , go left same

I left \Rightarrow I left yesterday.

Problem 7 PageRank Lollipop (7 credits)

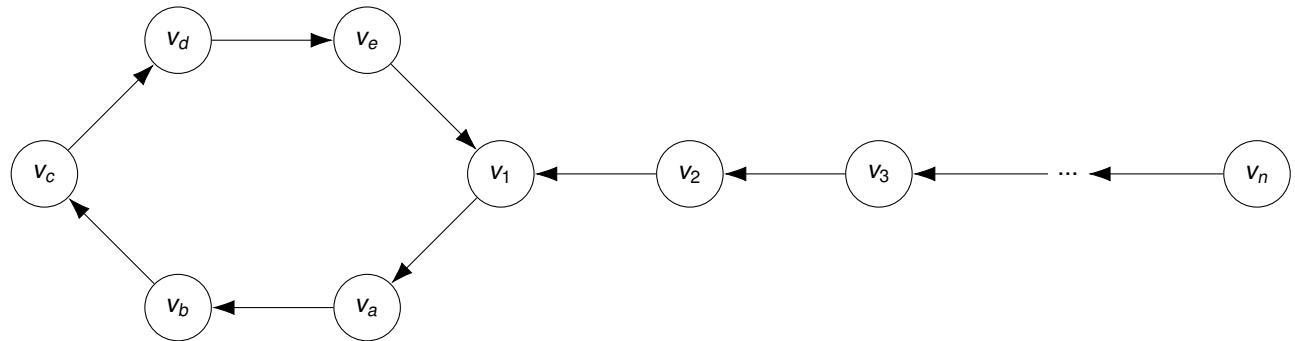


Figure 7.1: A directed “lollipop” graph with a tail of length n

Consider the directed, unweighted graph in Figure 7.1 with a “head” consisting of the six nodes v_1 , and v_a through v_e . Its “tail” consists of n nodes v_1, \dots, v_n where n is a parameter. Note, that we consider v_1 to be part of the head as well as the tail.

- 0 a) Which of the PageRank problems of *dead end*, *spider trap*, and *periodic states* apply here? Justify your answer. For each problem that applies, give a set of edges (at most 3 per problem) that would resolve that problem if they were inserted.
1
2

- ① This is a spider trap, as some nodes become unreachable once we enter the loop.
Connect and edge below the node in loop and open node
- ② Periodic problem exists in the loop (periodicity 6)
self loop to break periodicity.

In the following, we want to compute the topic-sensitive PageRank of the nodes. For that, we introduce the topic-sensitive teleport vector π where the topic is the "tail", i.e.

$$\pi_a = \dots = \pi_e = 0 \quad \text{and} \quad \pi_1 = \dots = \pi_n = \frac{1}{n}.$$

	0
	1
	2

State each node's PageRank equation for a teleport probability of $1 - \beta$.

$$\pi_n = \frac{1-\beta}{n} \quad \pi_i = \beta \pi_{i+1} + \frac{1-\beta}{n} \quad i \in \{2, \dots, n-1\}$$

$$\pi_a = \beta \pi_1 \quad \pi_b = \beta \pi_a \quad \pi_c = \beta \pi_b \quad \pi_d = \beta \pi_c$$

$$\pi_e = \beta \pi_d.$$

finally the node connecting the two.

$$\pi_1 = \beta(\pi_e + \pi_2) + \frac{1-\beta}{n}$$

c) Derive the sum of the PageRank of the "head" of the graph v_1, v_a, \dots, v_e as a function of n . Justify your answer.

$$\text{Reminder: } \sum_{i=0}^k \beta^i = \frac{1 - \beta^{k+1}}{1 - \beta} \quad \text{for } \beta \neq 1$$

	0
	1
	2
	3

$$\pi_a = \beta \pi_1, \quad \pi_b = \beta^2 \pi_1, \quad \pi_c = \beta^3 \pi_1, \quad \pi_d = \beta^4 \pi_1, \quad \pi_e = \beta^5 \pi_1,$$

$$\begin{aligned} \pi_2 &= \beta \pi_3 + \left(\frac{1-\beta}{n} \right) = \beta \left(\beta \pi_4 + \left(\frac{1-\beta}{n} \right) \right) + \left(\frac{1-\beta}{n} \right) = \beta^2 \pi_4 \\ &\quad + \beta \left(\frac{1-\beta}{n} \right) + \frac{1-\beta}{n} \\ &= \beta^{n-2} \pi_n + \left(\frac{1-\beta}{n} \right) \sum_{i=0}^{n-2} (\beta) \cdot \left(\frac{1-\beta}{n} \right)^i \left[\beta^{n-2} + \left(\frac{1-\beta}{n} \right) \right] \\ &\quad + \sum_{i=0}^{n-2} \beta^i \left(\frac{1-\beta}{n} \right)^{i+1} \end{aligned}$$

$$\pi_1 = \beta(\pi_e + \pi_2) + \frac{1-\beta}{n}$$

$$= \beta \left(\beta^5 \pi_1 + \frac{1-\beta}{n} \right)$$

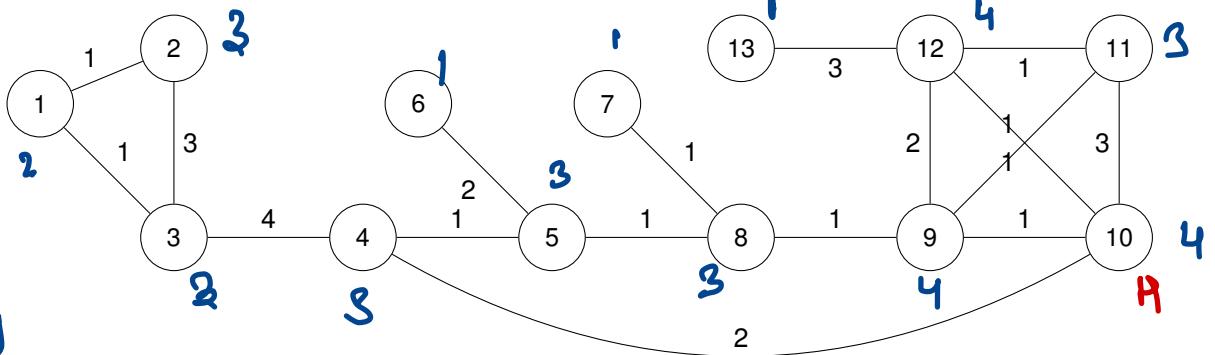
$$1 + \beta + \dots + \beta^{n-2}$$

$$= \left(\frac{1-\beta^{n-2}}{1-\beta} \right) \cdot \frac{1-\beta}{n}$$

~~Problem 8~~ Spectral Embeddings (3 credits)



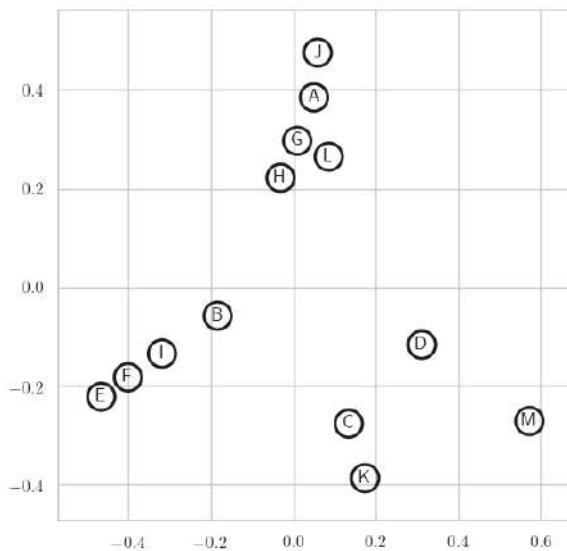
We use spectral embedding to map the following undirected, weighted graph into 2-dimensional space.



The following plot shows the embeddings as given by the eigenvectors belonging to the second and third smallest eigenvalues of the unnormalized Laplacian.

Second smallest e.v \rightarrow connectivity

$$\begin{bmatrix} f_2 \\ f_3 \end{bmatrix}$$



Fill out the following table that maps nodes in the graph to points in the plot.

Node	1	2	3	4	5	6	7	8	9	10	11	12	13
Point	E	F	I	B	C	K	M	D	L	H	G	A	J

Problem 9 Graph Neural Networks (4 credits)

Consider the following graph neural network.

4 hop

$$\mathbf{Z} = \sigma(\sigma(\hat{\mathbf{A}}\sigma(\hat{\mathbf{A}}\sigma(\hat{\mathbf{A}}\mathbf{X})\mathbf{W}_a)\mathbf{W}_b)\mathbf{W}_c)$$

$\mathbf{X} \in \mathbb{R}^{N \times D}$ are the node features, \mathbf{Z} are the node predictions, σ is the sigmoid function, \mathbf{W}_x are weight matrices of appropriate dimensions and $\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-\frac{1}{2}}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-\frac{1}{2}}$ is the propagation matrix as defined for GCNs, i.e. $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ is the adjacency matrix with added self-loops and $\tilde{\mathbf{D}}_{ii} = \sum_{j=1}^N \tilde{\mathbf{A}}_{ij}$ is the degree matrix of the amended adjacency matrix. Give the transformation and propagation depths of this model. Justify your answer.

Hint: We define

- propagation depth as the maximum distance that a model propagates information in the graph,
- transformation depth as the number of non-linear transformations that are applied to the features.

Each application of $\hat{\mathbf{A}}$ propagates info 1 level
so propagation depth = 4

transformation depth = no. of nonlinear
transforms = 4

b) Write down a formula for the node predictions \mathbf{Z} of a graph neural network with propagation depth 2 and transformation depth 2 where the transformations use fully connected layers with weight matrices \mathbf{W}_i and biases b_i .
Note: You can assume that operations such as matrix-vector summation broadcast as expected.

$$\hat{\mathbf{A}}\hat{\mathbf{A}}\sigma(\sigma(\mathbf{x}\mathbf{w}_1 + b_1)\mathbf{w}_2 + b_2)$$

	0
	1
	2

	0
	1
	2

Problem 10 Randomized Smoothing on Graph Neural Networks (2 credits)

Consider an arbitrary but fixed graph neural network $f(\mathbf{X})$ as base classifier. Here, \mathbf{X} denotes the adjacency matrix (i.e. the features are assumed to be constant).

For the smooth classifier $g(\mathbf{X})$ we use the randomization scheme $\phi(\mathbf{x})$:

$$g(\mathbf{x})_c = \mathcal{P}(f(\phi(\mathbf{x})) = c) = \sum_{\tilde{\mathbf{x}} \text{ s.t. } f(\tilde{\mathbf{x}})=c} \prod_{i=1}^{n^2} \mathcal{P}(\tilde{\mathbf{x}}_i | \mathbf{x}_i) \quad (1)$$

with

$$\mathcal{P}(\tilde{\mathbf{x}}_i | \mathbf{x}_i) = \begin{cases} p & \tilde{\mathbf{x}}_i = 1 - \mathbf{x}_i \\ 1 - p & \tilde{\mathbf{x}}_i = \mathbf{x}_i \end{cases} \quad (2)$$

Note that in contrast to the lecture here we do not distinguish between a probability for deleting p_d or adding p_a an edge. We simply use the "flip" probability p .

Furthermore, we consider the special case of $p = 0.5$.

- 0  a) What is the probability $\mathcal{P}(\tilde{\mathbf{x}}|\mathbf{x})$ for a directed graph or is there not enough information to determine it? Justify your answer.

1

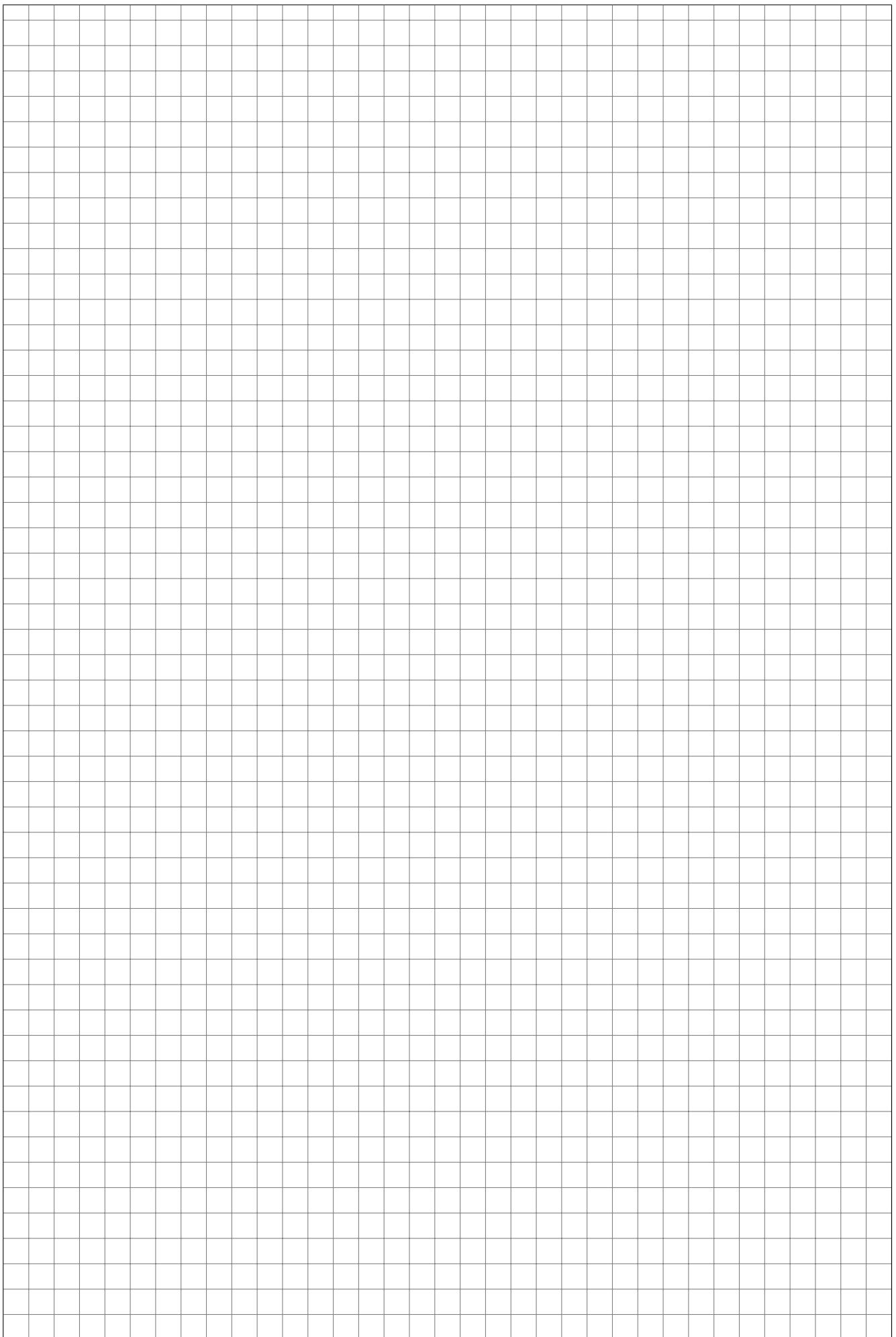
- 0  b) We are given the prediction for the original input $f(\mathbf{X}) = c$ and for a slightly different version $f(\mathbf{X}') \neq c$. What does the randomization scheme entail about the (exact) prediction of the smooth classifier $g(\mathbf{X})_c$ and $g(\mathbf{X}')_c$ (i.e. no Monte Carlo approximation)? Specifically, how do the probabilities for class c of the smooth classifier $g(\mathbf{X})_c$ and $g(\mathbf{X}')_c$ relate?

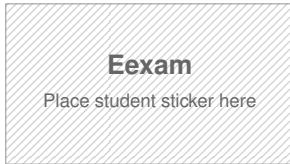
1. $g(\mathbf{X})_c = g(\mathbf{X}')_c$
2. $g(\mathbf{X})_c < g(\mathbf{X}')_c$
3. $g(\mathbf{X})_c > g(\mathbf{X}')_c$
4. There is not enough information to determine how $g(\mathbf{X})_c$ and $g(\mathbf{X}')_c$ relate

Justify your answer.

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

A large grid of squares, approximately 20 columns by 30 rows, intended for students to write their solutions. The grid is composed of thin black lines on a white background.





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.

Mining Massive Datasets

Exam: IN2323 / Endterm

Date: Friday 9th August, 2019

Examiner: Prof. Dr. Stephan Günnemann

Time: 13:30 – 15:00

P 1

P 2

P 3

P 4

P 5

P 6

P 7

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

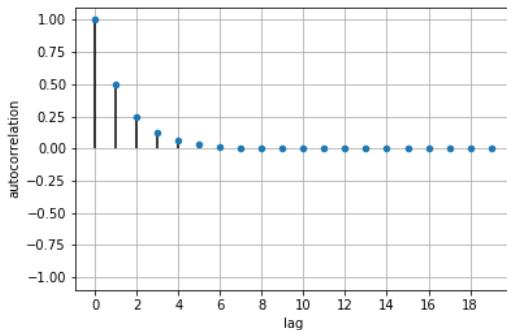
Working instructions

- This exam consists of **8 pages** with a total of **7 problems**.
Please make sure that you received a complete copy of the exam.
- You can earn 38 points.
- **Detaching pages from the exam is prohibited!**
- Allowed resources:
 - A4 sheet of handwritten notes (two sides)
 - **no other materials (e.g. books, cell phones, calculators) are allowed!**
- Only write on the sheets given to you by supervisors. If you need more paper, ask the supervisors.
- Last two pages can be used as scratch paper.
- All sheets (including scratch paper) have to be returned at the end.
- **Only use a black or a blue pen (no pencils, red or green pens)!**
- Write your answers only in the provided solution boxes or the scratch paper.
- **For problems that say "Justify your answer" or "Show your work" you only get points if you provide a valid explanation.** Otherwise it's sufficient to only provide the correct answer.
- Exam duration - 90 minutes.

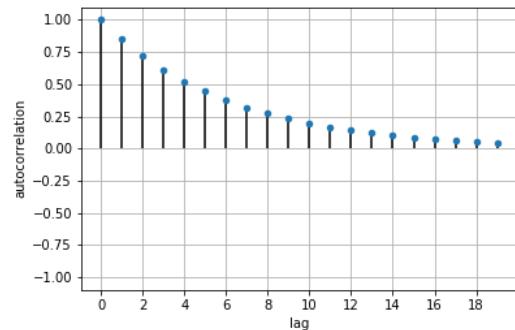
Left room from _____ to _____ / Early submission at _____

Problem 1 AR models: correlation function (4 points)

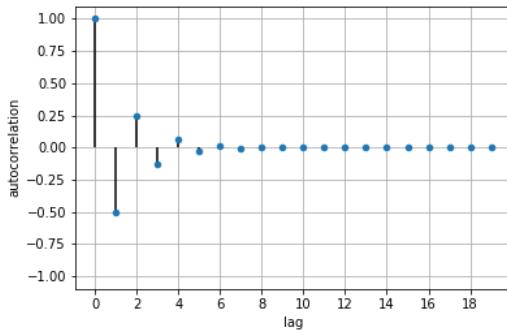
For each of the AR(1) and AR(2) models below find the corresponding autocorrelation function plot from Figure 1.1. Each plot was generated using one of the AR models from the list so that there is a one-to-one correspondence between the plots (1)-(4) and the AR processes (a)-(d). Everywhere $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$ with a positive σ .



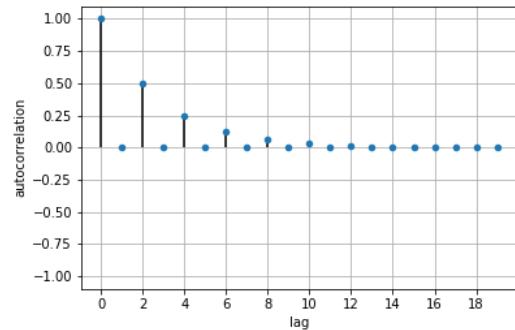
(1)



(2)



(3)



(4)

Figure 1.1: AR models: correlation function

a) $\mathcal{X}_t = -0.5\mathcal{X}_{t-1} + \epsilon_t$

- (1) (2) (3) (4)

b) $\mathcal{X}_t = 0.85\mathcal{X}_{t-1} + \epsilon_t$

- (1) (2) (3) (4)

c) $\mathcal{X}_t = 0.5\mathcal{X}_{t-2} + \epsilon_t$

- (1) (2) (3) (4)

d) $\mathcal{X}_t = 0.5\mathcal{X}_{t-1} + \epsilon_t$

- (1) (2) (3) (4)

Problem 2 Hidden Markov Models (6 points)

Consider the following Hidden Markov Model where Z_t are latent variables and X_t are continuous observed variables. We parametrize the prior and transition probabilities $P(Z_1 = i) = \pi_i$ and $P(Z_{t+1} = j | Z_t = i) = A_{ij}$ by:

$$\pi = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}, \quad A = \begin{bmatrix} 0.8 & 0.2 \\ 0.5 & 0.5 \end{bmatrix}$$

	0
1	
2	
3	
4	
5	
6	

We parametrize the emission probabilities with $X_t | Z_t = 0 \sim \text{Uniform}([8, 12])$ and $X_t | Z_t = 1 \sim \text{Uniform}([11, 13])$. We assume we observed $X = [8.5, 11, 13, 9.5]$.

Perform the Forward algorithm, compute α_t and compute the most probable state Z_t given X_1, \dots, X_t at every step t .

Hint 1: You can check your computations by comparing with $\alpha_4 = \begin{bmatrix} 3/3200 \\ 0 \end{bmatrix}$

Hint 2: Maintaining non-decimal fractions makes calculations easier (e.g. use 8/10 instead of 0.8).

$$\pi = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$

$$Pr(X=8.5 | Z_1=0) = \frac{1}{4}$$

$$x_1 = \underline{8.5}$$

$$Pr(X=8.5 | Z_1=1) = 0.$$

$$\alpha_1 = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix} \odot \begin{bmatrix} 1/4 \\ 0 \end{bmatrix} = \begin{bmatrix} 1/8 \\ 0 \end{bmatrix}$$

$$x_2 = 11$$

$$Pr(X=11 | Z_2=0) = 1/4$$

$$Pr(X=11 | Z_2=1) = 1/2$$

$$\alpha_2 = \begin{bmatrix} 1/4 \\ 1/2 \end{bmatrix} \odot \left(\begin{bmatrix} 0.8 & 0.2 \\ 0.5 & 0.5 \end{bmatrix} \begin{bmatrix} 1/8 \\ 0 \end{bmatrix} \right)$$

$$= \begin{bmatrix} 1/4 \\ 1/2 \end{bmatrix} \odot \begin{pmatrix} 1/10 \\ 1/40 \end{pmatrix} = \begin{bmatrix} 1/40 \\ 1/80 \end{bmatrix}$$

$$x_3 = 13$$

$$Pr(X=13 | Z_3=0) = 0$$

$$Pr(X=13 | Z_3=1) = 1/2$$

$$\alpha_3 = \begin{bmatrix} 0 \\ 1/2 \end{bmatrix} \odot \begin{bmatrix} 4/5 & 1/2 \\ 1/2 & 1/2 \end{bmatrix} \begin{bmatrix} 1/40 \\ 1/80 \end{bmatrix}$$

$$= \begin{bmatrix} 0 \\ 9/1600 \end{bmatrix}$$

$$x_4 = 11$$

probable states

$$\alpha_4 = \begin{bmatrix} 9/12800 \\ 0 \end{bmatrix}$$

$$0, 0, 1, 0$$

~~Problem 3~~ Recurrent neural networks (8 points)

Given a sequence of integers $\{x_1, \dots, x_n\}$, $x_t \in \{0, 1\}$, the goal is to output 1 at step t if the sum of all the inputs up to step t is even. That is, $h_t = 1$ if $\sum_{i=0}^t x_i$ is even, otherwise 0. For example, the outputs for a sequence 01101 are 10110. To model this you use an RNN with the following update equations:

$$\begin{aligned}\tilde{h}_t &= a \cdot h_{t-1} + b \cdot x_t \\ z_t &= f(c \cdot h_{t-1} + d \cdot x_t) \\ h_t &= f((1 - z_t) \cdot h_{t-1} + z_t \cdot \tilde{h}_t)\end{aligned}$$

where $h_0 = 1$ and $a, b, c, d \in \mathbb{R}$ and $f(x) = 1$ if $x > 0.5$, otherwise 0.

- 0 a) Fill out the table with all the possible input values x_t and input hidden state values h_{t-1} , and calculate the corresponding output h_t according to the given requirements.
1 2

sum upto t-1 → value added

h_{t-1}	x_t	h_t
0	0	0
0	1	1
1	0	1
1	1	0

- 0 1 b) Find the values for a, b, c and d such that this model achieves 100% accuracy. Show your work.
Hint: think about what can happen with a new input and what should gate z_t do in those cases.
2 3 4 5 6

$$z_t = f(c \cdot h_{t-1} + d \cdot x_t).$$

we know $x_t = 0$ retains the input
 $x_t = 1$ flips the input

$$c = 0 \quad d = 1$$

$h_t = f(\underline{x_t})$ just care whether
 $x_t = 1$ or 0

$h_{t-1} \Rightarrow$ past setting is forgotten

$$h_t = f(\underline{\tilde{h}_t})$$

$$f(b) = 1$$

$$f(a) = 1$$

$$f(a+b) = 0$$

Problem 4 Graph laws (4 points)

You are given an Erdős-Rényi graph $G(n, p)$, where n is the number of nodes and p is the probability of an edge.

a) What is the expected number of cliques of size m ?

$$\binom{n}{m} p^{\frac{m(m-1)}{2}}$$

0
 1
 2

b) We add a new node and connect it to every existing node with probability q . What is the expected number of newly created triangles?

$$\binom{n}{2} p q^2$$

0
 1
 2

Problem 5 Deep Generative Models (6 points)

The loss used in generative adversarial networks (GANs) can be written in the following form:

$$\min_{\theta} \max_{\phi} \mathcal{L}(\theta, \phi) = \min_{\theta} \max_{\phi} \mathbb{E}_{p^*(x)} [\log D_\phi(x)] + \mathbb{E}_{p(z)} [\log(1 - D_\phi(f_\theta(z)))]$$

distribution of noise

where $p^*(x)$ is the true data distribution, $p(z)$ is the distribution of the noise, f_θ is the generator, and D_ϕ is the discriminator.

From sheet

a) For a given generator (fixed parameters θ) assume there exists a discriminator $D_{\phi^*}(x)$ with parameters ϕ^* such that for all x :

$$D_{\phi^*}(x) = \frac{p^*(x)}{p^*(x) + p_\theta(x)}$$

where $p_\theta(x)$ is the distribution learned by the generator. Show that D_{ϕ^*} is **optimal**, i.e. $\phi^* = \arg \max_{\phi} \mathcal{L}(\theta, \phi)$. Hint: $\max_y [a \log(y) + b \log(1 - y)] = \frac{a}{a+b}$ for any $a, b \in \mathbb{R}_0^+, a + b > 0$.

$$\begin{aligned} \max_{\phi} \mathcal{L}(\theta, \phi) &= \max_{\phi} \mathbb{E}_{p^*(x)} [\log D_\phi(x)] + \mathbb{E}_{p(z)} [\log(1 - D_\phi(f_\theta(z)))] \\ &= \max_{\phi} \left\{ \int (p^*(x) \log D_\phi(x)) dx + \mathbb{E}_{p(z)} [\log(1 - D_\phi(f_\theta(z)))] \right\} \\ &= \max_{\phi} \left[\int \underbrace{p^*(x)}_a \underbrace{\log D_\phi(x)}_y + \underbrace{p_\theta(x)}_b \underbrace{\log(1 - D_\phi(x))}_{1-y} dx \right] \\ &= \int \left[\max_{\phi} \left\{ p^*(x) \log D_\phi(x) + p_\theta(x) \log(1 - D_\phi(x)) \right\} \right] dx \\ &= \int \left(\frac{p^*(x)}{p^*(x) + p_\theta(x)} \right) dx \end{aligned}$$

b) Show that $\mathcal{L}(\theta, \phi^*) = -\log(4) + 2 \cdot \text{JSD}(p^*(x) || p_\theta(x))$ where

$$\text{JSD}(p^*(x) || p_\theta(x)) = \frac{1}{2} \left[\text{KL}(p^*(x) || m(x)) + \text{KL}(p_\theta(x) || m(x)) \right]$$

is the Jensen–Shannon divergence, KL is the Kullback–Leibler divergence, and $m(x) = \frac{p_\theta(x) + p^*(x)}{2}$.

$$\begin{aligned} \mathcal{L}(\theta, \phi^*) &= \mathbb{E}_{p^*(x)} [\log D_{\phi^*}(x)] + \mathbb{E}_{p(z)} [\log(1 - D_{\phi^*}(f_\theta(z)))] \\ 2 \text{JSD}(p^*(x) || p_\theta(x)) &= \left[\text{KL}(p^*(x) || m(x)) + \text{KL}(p_\theta(x) || m(x)) \right] \\ &= \int p^*(x) \cdot \log \left(\frac{2p^*(x)}{p_\theta(x) + p^*(x)} \right) dx + p_\theta(x) \cdot \log \left(\frac{2p_\theta(x)}{p_\theta(x) + p^*(x)} \right) \\ &= \frac{\int 2p^*(x) \cdot dx}{2} + \frac{\int 2p_\theta(x) \cdot dx}{2} + \dots \\ &= 2(-\log(4)) + \dots \end{aligned}$$

Problem 6 Ranking (6 points)

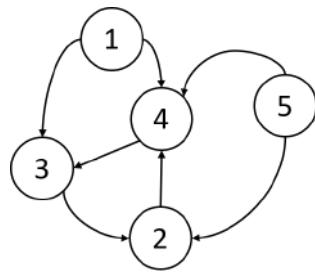
Given the following graph G let $e = (s, t)$ be a new edge, i.e. an edge which is not yet in the graph G . We denote with G_{new} the resulting graph when adding e to G . Your task is to identify the edge e based on the following observation:

The PageRank vector of G_{new} based on a teleport set $S = \{1, 5\}$ and $\beta = 0.85$ (i.e. the probability to teleport is 0.15) is given by

~~+~~

$$\pi = [0.0750, 0.1935, 0.2668, 0.2764, 0.1884]$$

What is the edge e (specify the source node s and target node t) that has been added? Justify your answer.



	0
	1
	2
	3
	4
	5
	6

$$r_1 = \frac{0.15}{1+1} = \frac{0.15}{2} = \underline{\underline{0.075}} \quad r_5 = \underline{\underline{0.075}}$$

now we need to determine the source node

$$\text{for } 1 \rightarrow 2 \quad M_2 = 0.075 + \beta \frac{\pi_1}{2} = 0.0962$$

$$\text{for } 2 \rightarrow 2 \quad M_2 = 0.075 + \beta \frac{\pi_2}{2}$$

for $3 \rightarrow 2$ ~~✓~~ matches

does not match for node 2, hence node 2 is involved as target node.

Problem 7 Spectral clustering (4 points)

SBM

You are given a graph G with $N = 1000$ nodes that is generated from a Stochastic Block Model with the following parameters

$$\pi = [0.5 \quad 0.5] \quad \eta = \begin{bmatrix} 0.05 & 0.5 \\ 0.5 & 0.05 \end{bmatrix}$$

	0
	1
	2
	3
	4

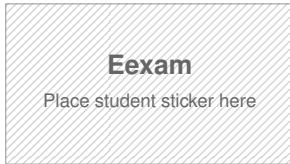
Assume that G is connected. Do you expect spectral clustering to recover the true communities z when applied to the graph G (yes or no)? Justify your answer.

as off diagonal in η is more than leading diagonals there are more edges connecting the different communities

spectral clustering minimizes the inter-community edges \rightarrow won't work.

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

A large grid of squares, approximately 20 columns by 30 rows, intended for students to write their solutions. The grid is composed of thin black lines on a white background.



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.

Mining Massive Datasets

Exam: IN2323 / Retake

Date: Monday 30th September, 2019

Examiner: Prof. Dr. Stephan Günnemann

Time: 08:00 – 09:30

P 1

P 2

P 3

P 4

P 5

P 6

P 7

P 8

I								
---	--	--	--	--	--	--	--	--

Working instructions

- This exam consists of **12 pages** with a total of **8 problems**.
Please make sure that you received a complete copy of the exam.
- You can earn 43 points.
- **Detaching pages from the exam is prohibited!**
- Allowed resources:
 - A4 sheet of handwritten notes (two sides)
 - **no other materials (e.g. books, cell phones, calculators) are allowed!**
- Only write on the sheets given to you by supervisors. If you need more paper, ask the supervisors.
- Last two pages can be used as scratch paper.
- All sheets (including scratch paper) have to be returned at the end.
- **Only use a black or a blue pen (no pencils, red or green pens)!**
- Write your answers only in the provided solution boxes or the scratch paper.
- **For problems that say "Justify your answer" or "Show your work" you only get points if you provide a valid explanation.** Otherwise it's sufficient to only provide the correct answer.
- Exam duration - 90 minutes.

Left room from _____ to _____ / Early submission at _____

Problem 1 AR models: stationarity (5 points)

Decide whether the following AR models are stationary or not. Everywhere $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$ with a positive σ .

Mark correct answers with a cross



To undo a cross, completely fill out the answer option



To re-mark an option, use a human-readable marking



a) $X_t = c + 0.1X_{t-1} + \epsilon_t$ with some $c \in \mathbb{R}$

$\Phi(L) = 1 + 0.1L$

stationary for $|c| > 0.1$, otherwise non-stationary

yes, always stationary

no, always non-stationary

b) $X_t = -3 + 0.2X_{t-1} - 0.01X_{t-2} + c\epsilon_t$ with some $c \in \mathbb{R} \setminus \{0\}$

stationary for $c > 0$, otherwise non-stationary

yes, always stationary

no, always non-stationary

$$\Phi(L) = 1 - \sum_{i=1}^2 \phi_i L^i = 1 - \underbrace{0.2L + 0.01L^2}_{(1 - 0.1L)^2} + 0.001L^2.$$

$$(1 - 0.1L)^2 \rightarrow L = 10 \text{ outside unit circle}$$

c) $X_t = 1 + 0.3X_{t-1} - 0.03X_{t-2} + 0.001X_{t-3} + \epsilon_t$

no, non-stationary

yes, stationary

$$\Phi(L) = 1 - 0.3L - 0.03L^2 + 0.001L^3 = 0$$

d) $X_t = -2 + 0.5X_{t-n} + \epsilon_t$ with some $n \in \mathbb{N}$

no, always non-stationary

stationary for $n \leq 2$, otherwise non-stationary

yes, always stationary

* const. can be ignored

e) $X_t = 2019 - \sum_{i=1}^n a^i X_{t-i} + \epsilon_t$ with some $n \in \mathbb{N}$, $a \in \mathbb{R} \setminus \{0\}$

stationary for $|a| < 1$, otherwise non-stationary

stationary for $|a|^n < 2019$, otherwise non-stationary

stationary for $n = 1$, $|a| < 1$, otherwise non-stationary

$$\Phi(L) = 2019 + \sum_{i=1}^n (aL)^i$$

$$2019 + \frac{1 - (aL)^{n+1}}{1 - aL}$$

$$\text{if } a \leq 1 \quad |aL| = \frac{1}{|a|}$$

Problem 2 Hidden Markov Models (7 points)

Consider the following Hidden Markov Model where $Z_t \in \{0, 1\}$ are latent variables and $X_t \in \{a, b\}$ are discrete observed variables. We parametrize the prior and transition probabilities $P(Z_1 = i) = \pi_i$, $P(Z_{t+1} = j|Z_t = i) = A_{ij}$ and $P(X_t = j|Z_t = i) = B_{ij}$ by:

$$\pi = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}, \quad A = \begin{bmatrix} \alpha & 1-\alpha \\ 1/2 & 1/2 \end{bmatrix}, \quad B = \begin{bmatrix} 3/4 & 1/4 \\ 1/4 & 3/4 \end{bmatrix}$$

We assume we observed $X = [a, b, a]$.

- a) Compute $P(Z_2|X_1, X_2)$ and $P(Z_2|X_1, X_2, X_3)$ as a function of α .

a) (i) Forward algorithm :

$$\alpha_1 = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix} \odot \begin{bmatrix} 3/4 \\ 1/4 \end{bmatrix} = \begin{bmatrix} 3/8 \\ 1/8 \end{bmatrix}$$

$$\alpha_2 = \begin{bmatrix} 1/4 \\ 3/4 \end{bmatrix} \odot \begin{bmatrix} \alpha & 1-\alpha \\ 1/2 & 1/2 \end{bmatrix} \begin{bmatrix} 3/8 \\ 1/8 \end{bmatrix} = \frac{1}{32} \begin{bmatrix} 3\alpha + 1 \\ 21/2 - 9\alpha \end{bmatrix}$$

(ii) Backward

$$\beta_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\begin{aligned} \beta_2 &= \begin{bmatrix} \alpha & 1-\alpha \\ 1/2 & 1/2 \end{bmatrix} \left(\begin{bmatrix} 3/4 \\ 1/4 \end{bmatrix} \odot \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right) = \begin{bmatrix} \alpha & 1-\alpha \\ 1/2 & 1/2 \end{bmatrix} \begin{bmatrix} 3/4 \\ 1/4 \end{bmatrix} \\ &= \frac{1}{4} \begin{bmatrix} 2\alpha + 1 \\ 2 \end{bmatrix} \end{aligned}$$

	0
	1
	2
	3
	4
	5

0 b) What is the most probable state Z_2 according to $P(Z_2|X_1, X_2)$ and $P(Z_2|X_1, X_2, X_3)$ for any value of $\alpha \in [0, 1]$.

1 Hint: $\sqrt{244} \approx 15.62$

2

Problem 3 RNNs & Word vectors (7 points)

You are solving a question-answering task. Given a context and a question, the goal is to find the answer inside the context. Below are two examples (1 and 2).

id	Context	Question	Answer
1	Mary was in the bathroom. Then she moved to the hallway.	Where is Mary?	hallway
2	John is in the hallway. Mary is there as well.	Where is Mary?	hallway

Assume that the question is represented with the vector \mathbf{q} . We want to know what is the probability that a word from the context is the answer. We decide to somehow represent every word w_i with an embedding \mathbf{h}_i and pass it together with \mathbf{q} through a neural network to get the probabilities. The only thing left to do is to decide how to get \mathbf{h}_i . We propose two approaches: sliding window and RNN.

- a) **Sliding window** — Every word w_i is represented with a pretrained word vector \mathbf{v}_i . A sliding window of size 2 takes the neighbouring words and constructs the embedding for w_i as a sum of vectors: $\mathbf{h}_i = \sum_{j=i-2}^{i+2} \mathbf{v}_j$. Is it possible for this model to find the right answer in example 1? What about example 2? Justify.

0
 1
 2

- b) **RNN** — As an alternative we use an RNN that takes pretrained word vectors \mathbf{v}_i from left to right and outputs \mathbf{h}_i as a word embedding. Why is this model able to output the right answer in example 1? Explain why it could fail answering example 2?

0
 1
 2

- c) Propose another model that overcomes the shortcomings of the sliding window and the RNN. Describe what the input would be and how would you calculate the embedding \mathbf{h}_i . Explain why this model could work on both examples.

0
 1
 2
 3

Problem 4 Deep Generative Model (5 points)

- 0 a) You are given a pseudo code implementation of 4 different variants of an AutoEncoder. Here, $\mathbf{x}_i \in \mathbb{R}^d$ is the input data, and $f_\theta : \mathbb{R}^d \mapsto \mathbb{R}^k$ and $g_\phi : \mathbb{R}^k \mapsto \mathbb{R}^d$ are fully connected feed-forward neural networks where θ and ϕ are learnable parameters. The final layers of f_θ and g_ϕ have no (i.e. have linear) activation functions. I_k is a $k \times k$ identity matrix, $\mathbf{0}_k$ is a k -dimensional vector of zeros, \mathcal{N} is the Normal distribution, σ is the softmax function, and $\text{diag}(\mathbf{x})$ takes a vector $\mathbf{x} \in \mathbb{R}^k$ and returns a $k \times k$ diagonal matrix with the vector values on the diagonal.

AutoEncoder 1

$$\begin{aligned}\epsilon_i &\sim \mathcal{N}(\mathbf{x}_i, I_d) \\ \mathbf{h}_i &= f_\theta(\epsilon_i) \\ \tilde{\mathbf{x}}_i &= g_\phi(\mathbf{h}_i) \\ \mathcal{L} &= \sum_i \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|_2\end{aligned}$$

AutoEncoder 2

$$\begin{aligned}\mathbf{h}_i &= f_\theta(\mathbf{x}_i) \\ \epsilon_i &\sim \mathcal{N}(\mathbf{h}_i, I_k) \\ \tilde{\mathbf{x}}_i &= g_\phi(\epsilon_i) \\ \mathcal{L} &= \sum_i \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|_2\end{aligned}$$

AutoEncoder 3

$$\begin{aligned}\mathbf{h}_i &= f_\theta(\mathbf{x}_i) \\ \epsilon_i &\sim \mathcal{N}(\mathbf{0}_k, I_k) \\ \tilde{\mathbf{x}}_i &= g_\phi(\mathbf{h}_i + \epsilon_i) \\ \mathcal{L} &= \sum_i \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|_2\end{aligned}$$

Is it **necessary** to use the reparametrization trick in order to compute the gradients of \mathcal{L} w.r.t. to **both** θ and ϕ in the above implementations? Answer with Yes or No and provide a justification. If the answer is Yes, modify the pseudo code to implement the reparametrization trick.

AutoEncoder 1

reparametrization not required as sampling does not depend on parameters.

AutoEncoder 2

sampling $\epsilon_i \sim \mathcal{N}(\underline{\mathbf{h}}_i, I_k) \rightarrow \mathbf{h}_i = f_\theta(\mathbf{x}_i) \therefore \text{reparametrization is required}$

$$\epsilon_i \sim \mathcal{N}(0, I_k)$$

$$\mathbf{h}_i = \mathbf{h}_i + \mathcal{N}(0, I_k)$$

AutoEncoder 3

Not required

- b) Assume the same setup as in a). The model specified by the following pseudo code is **not well defined**. Specify the reason why, and modify the pseudo code such that the model becomes well-defined. In addition, if you think it is **necessary** to use the reparametrization trick, please include it in your implementation.

0
 1
 2

$$\begin{aligned}\mathbf{h}_i &= f_{\theta}(\mathbf{x}_i) \\ \epsilon_i &\sim \mathcal{N}(\mathbf{0}_k, \text{diag}(\mathbf{h}_i)) \\ \tilde{\mathbf{x}}_i &= g_{\phi}(\epsilon_i) \\ \mathcal{L} &= \sum_i \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|_2\end{aligned}$$

Problem 5 Spectral clustering (3 points)

0
1
2
3

Given is the following matrix $M \in \mathbb{R}^{9 \times 9}$.

$$M = \begin{bmatrix} 1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ -1 & -1 & 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & -1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 2 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & -1 & 3 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & 1 \end{bmatrix}$$

Write down the **exact** value of the three smallest eigenvalues $\lambda_1, \lambda_2, \lambda_3$ of M and the respective eigenvectors x_1, x_2, x_3 .

2 connected components

$$\therefore \lambda_1 = \lambda_2 = \lambda_3 = 0.$$

~~Problem 6~~ Spectral clustering (3 points)

0
1
2
3

You are given an undirected graph $G = (V, E)$. It is known that the second smallest eigenvalue of the unnormalized Laplacian $L = D - W$ is equal to 10.

Let $\phi(G)$ denote the best possible ratio cut achievable on the graph G

$$\phi(G) = \min_{S \subset V} \text{ratio-cut}(S, \bar{S})$$

What are possible values of $\phi(G)$ for the given graph? Select all that apply. Justify your answer.

- a) 1
- b) 2
- c) 4
- d) 8
- e) 16

$\lambda_2 = \underline{\underline{10}}$.

relaxation of
min ratio cut
problem.

$$\frac{\text{ratio-cut}(S, \bar{S})}{|S|} \geq \lambda_2 \left(1 - \frac{|S|}{|V|}\right)$$

$$\text{ratio-cut}(S, \bar{S}) \geq \lambda_2 |S| \left(1 - \frac{|S|}{|V|}\right)$$

This should
be 10

$$10 \left(1 - \frac{1}{10}\right)$$

Problem 7 Ranking (5 points)

Given a graph G with 5 nodes, assume that you have access to several topic-sensitive PageRank vectors, each pre-computed using a **different** teleport set S , and the **same (fixed)** teleport parameter β , $0 < \beta < 1$.

- $\pi_{235} \in \mathbb{R}^5$, with teleport set $S = \{2, 3, 5\}$
- $\pi_{124} \in \mathbb{R}^5$, with teleport set $S = \{1, 2, 4\}$
- $\pi_{134} \in \mathbb{R}^5$, with teleport set $S = \{1, 3, 4\}$
- $\pi_3 \in \mathbb{R}^5$, with teleport set $S = \{3\}$

Assume that the random walker always teleports **uniformly at random** to each node in the teleport set.



Is it possible to compute each of the following PageRank vectors without access to the graph G , i.e. using only the above pre-computed vectors? If so, specify the exact equation as a function of $\pi_{235}, \pi_{124}, \pi_{134}$ and π_3 . If not, justify why not.

a) Is it possible to compute $\pi_{14} \in \mathbb{R}^5$ with teleport set $S = \{1, 4\}$?

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

$$\pi_{14} = \frac{3 * \underline{\pi_{134}} - \pi_3}{2}$$

b) Is it possible to compute $\pi_5 \in \mathbb{R}^5$ with teleport set $S = \{5\}$?

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

$$\pi_5 = 3 * \underline{\pi_{134}} - \pi_3 \quad \pi_2 = 3 * \pi_{124} - 3 * \pi_{134} + \pi_3$$

$$\pi_5 = 3 * \pi_{235} - 2 * \pi_{124} + 3 * \pi_{134} - 2 * \pi_3$$

c) Is it possible to compute $\pi_1 \in \mathbb{R}^5$ with teleport set $S = \{1\}$?

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

No we can't

d) Is it possible to compute $\pi_w \in \mathbb{R}^5$ with teleport set $S = \{1, 2, 3, 4, 5\}$, where we do not teleport to each node uniformly at random but rather with weights 0.2, 0.3, 0.1, 0.2, 0.2, respectively?

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

Problem 8 Graph Neural Networks (8 points)

Given an unweighted, undirected graph G with adjacency matrix $\mathbf{A} \in \{0, 1\}^{N \times N}$ and node attribute matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$, your task is to perform semi-supervised node classification with C classes using a graph convolutional network (GCN).

In matrix notation, a GCN with K layers is recursively defined as follows.

$$\begin{aligned}\mathbf{H}^{(0)} &= \mathbf{X} \\ \mathbf{H}^{(k)} &= \sigma^{(k)}(\tilde{\mathbf{A}} \mathbf{H}^{(k-1)} \mathbf{W}^{(k)}) \quad \text{for } k \in 1, \dots, K\end{aligned}$$

That is, $\mathbf{H}^{(K)} \in \mathbb{R}^{N \times C}$ contains the class predictions for **all nodes** stacked in a matrix. Here, $\sigma^{(k)}$ is the ReLU(\cdot) activation function for $k \in \{1, \dots, K-1\}$ and the softmax(\cdot) function for the final (output) layer $k = K$. $\mathbf{W}^{(k)}$ is the weight matrix of layer k .

$\tilde{\mathbf{A}} \in \mathbb{R}^{N \times N}$ is a degree-normalized version of the adjacency matrix whose entries are

$$\tilde{\mathbf{A}}_{uv} = \begin{cases} \frac{1}{\sqrt{d_u d_v}} & \text{if } \mathbf{A}_{uv} = 1 \\ \frac{1}{d_u} & \text{if } u = v \\ 0 & \text{else,} \end{cases}$$

where d_u is the degree of node u .

- 0 a) GCN is an instance of the differentiable message passing framework. Provide the corresponding message function M and update function U of GCN. For simplicity, you may write $\sigma^{(k)}$ for all layers.

1

2

3

- 0 b) Assume you have a GCN with 3 layers, i.e. $K = 3$. Provide the non-recursive (unrolled) expression for $\mathbf{H}^{(3)}$. That is, write down the single-line equation of $\mathbf{H}^{(3)}$ in matrix form.

1

c) We consider now the two following situations:

- (1) $\mathbf{A}_{uv} = 1$ for $u = v$ and 0 else.
- (2) $\sigma^{(k)}(x) = x$, i.e. identity activation function, for $k \in \{1, 2\}$

Apart from the additional information in each situation, the models are identical to the GCN defined above. Simplify the single-line, matrix-form expression of $\mathbf{H}^{(3)}$ given the additional information provided in each situation.

d) For each situation, find one equivalent model in the table below. **You may select each option only once.** Briefly justify your answer for each situation.

Recurrent neural network (RNN)	Linear regression
Feed-forward neural network (FFNN)	Label Propagation (LP)
Linear function	Deep Generative Model
Logistic regression on \mathbf{X}	Logistic regression on pre-processed features $\tilde{\mathbf{X}}$

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

A large grid of squares, approximately 20 columns by 30 rows, intended for students to write their solutions. The grid is composed of thin black lines on a white background.

Mining Massive Datasets — Final Exam

1	2	3	4	5	6	7	8	9	10	11	Σ
/8	/10	/11	/10	/8	/4	/8	/6	/4	/4	/12	/48

Do not write anything above this line

Name:

Student ID:

Signature:

- Only write on the sheets given to you by supervisors. If you need more paper, ask the supervisors.
- Pages 14-16 can be used as scratch paper.
- All sheets (including scratch paper) have to be returned at the end.
- **Do not unstaple the sheets!**
- Wherever answer boxes are provided, please write your answers in them.
- Please write your student ID (*Matrikelnummer*) on every sheet you hand in.
- **Only use a black or a blue pen (no pencils, red or green pens!).**
- You are allowed to use your A4 sheet of handwritten notes (two sides). **No other materials (e.g. books, cell phones, calculators) are allowed!**
- Exam duration - 90 minutes.
- This exam consists of 12 pages, 7 problems. You can earn 48 points.

1 Hidden Markov Models

Problem 1 [10 points] In this question, we will discuss hidden Markov models with **continuous** observations. We will use the notation from the lecture, where $Z_t \in \{1, \dots, K\}$ denotes the state at time t and $X_t \in \mathbb{R}$ denotes the observation at time t . The conditional probability of an observation at a state k is $\Pr(X_t | Z_t = k, \boldsymbol{\theta}) = \mathcal{N}(X_t | \mu_k, \sigma_k^2)$, i.e. a Gaussian distribution parametrized by mean μ_k and variance σ_k^2 . $\boldsymbol{\theta}$ is the set of parameters of the HMM, which includes the initial probabilities $\boldsymbol{\pi} \in \mathbb{R}^K$, transition probability matrix $\mathbf{A} \in \mathbb{R}^{K \times K}$ and the means and variances $\{\mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2\}$.

- a) Write down the log-likelihood $\log \Pr(Z_1, \dots, Z_T, X_1, \dots, X_T | \mathbf{A}, \boldsymbol{\pi}, \mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2)$ as a function of \mathbf{A} , $\boldsymbol{\pi}$, μ_k 's and σ_k^2 's.

- b) Write the forwards and backwards update equations for this HMM. Explain in a single line how they are different from the updates we studied in class (discrete observations).

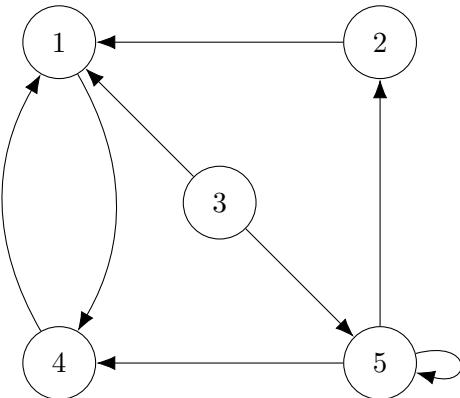
- c) You are given a sequence of observations $\{X_1, \dots, X_T\}$ and the corresponding states $\{Z_1, \dots, Z_T\}$.

Are the maximum likelihood estimates of A_{ij} and π_i for this model different from the ones for HMM with discrete observations (that we studied in class)? Explain why or why not.

- d) You are given a sequence of observations $\{X_1, \dots, X_T\}$ and the corresponding states $\{Z_1, \dots, Z_T\}$. Write down the closed-form maximum likelihood estimates for the parameters μ_k and σ_k^2 . You don't have to derive the closed-form MLE, stating it with an intuitive explanation is enough.

2 PageRank

Problem 2 [4 points] Given the graph above and the topic-specific PageRank vector $r = [0.3662, 0.0442, 0.0800, 0.3519, 0.1578]$ using the teleport set $S = \{3, 5\}$. What is the value of β (corresponding to the teleport probability $1 - \beta$)? Justify your answer.



3 Graph Clustering

Problem 3 [8 points] You are given a connected undirected unweighted graph with a set of nodes V . Let S_n^* be the partitioning minimizing the **normalized cut**, and let S_r^* be the partitioning minimizing the **ratio cut**. That is:

$$S_n^* = \arg \min_S \text{n-cut}(S) \quad S_r^* = \arg \min_S \text{r-cut}(S)$$

where $S = \{C_1, V \setminus C_1\}$ is a two-way partitioning of the nodes in the graph.

Prove or disprove that $\text{n-cut}(S_n^*) \leq \text{r-cut}(S_r^*)$.

Hint: Think about the values of the cuts when we plug in the same partitioning S .

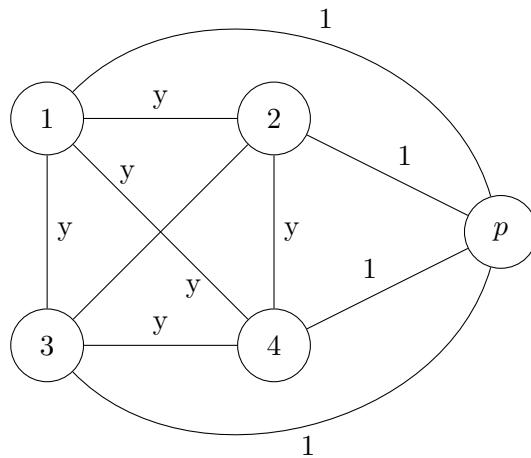
Furthermore, by definition we have that $\text{n-cut}(S_n^*) \leq \text{n-cut}(S)$. Combining these facts we obtain for any S :

$$\text{n-cut}(S_n^*) \leq \text{n-cut}(S) \leq \text{r-cut}(S)$$

Finally, replacing S with S_r^* above we obtain $\text{n-cut}(S_n^*) \leq \text{r-cut}(S_r^*)$.

You are given an undirected, weighted graph with a set of nodes $V = \{1, 2, \dots, n, p\}$. It contains one clique of n nodes that is fully interconnected, each edge having weight y , and one additional node p that is connected to all nodes in the clique with weight 1.

See the example figure below for $n = 4$ as an illustration. We want to partition the nodes into **two** clusters.

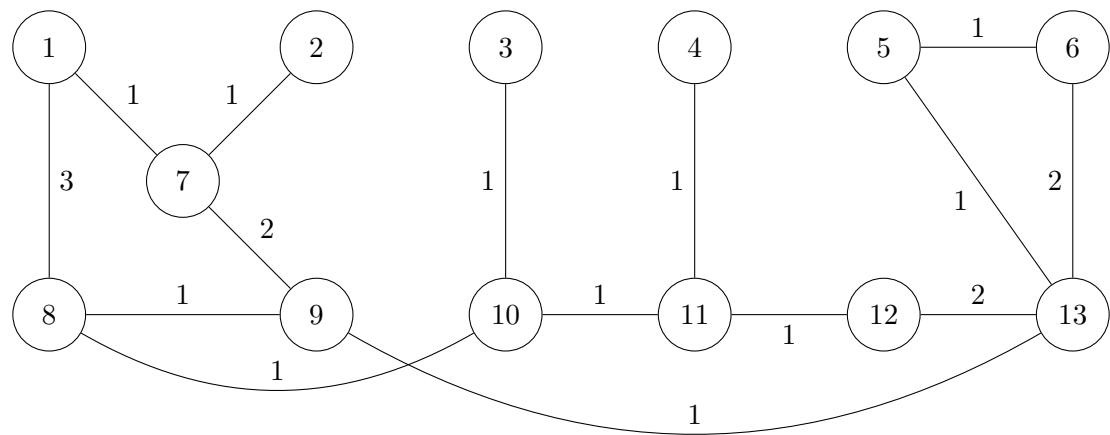


- a) What is the ratio cut when one of the clusters contains only the node p ? Provide the solution as a function of y and n .

- b) What is the ratio cut if one of the clusters corresponds to $S \subseteq \{1, 2, \dots, n\}$, with $|S| = k > 0$? Provide the solution as a function of y , n , and k .

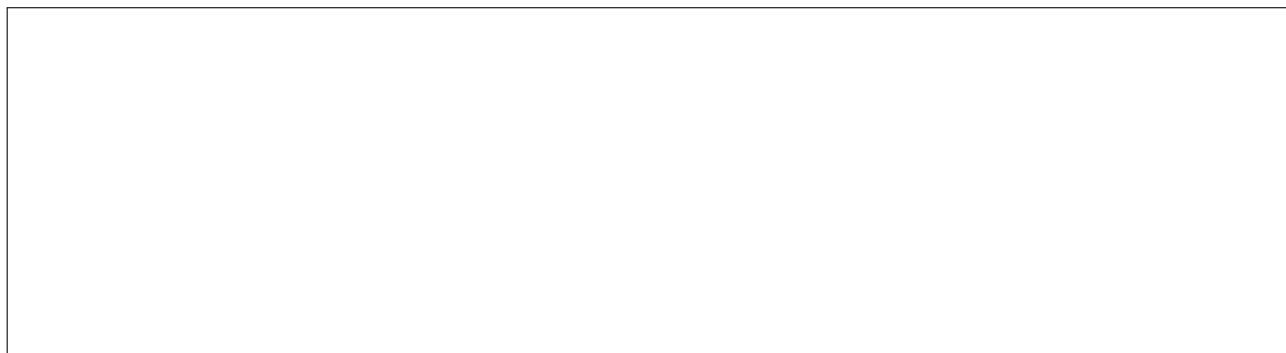
- c) When is the partitioning $\{p\}$ vs $\{1, 2, \dots, n\}$ the optimal solution for the minimum ratio cut? Specify all values of y and n , such that this holds.

Problem 4 [6 points] Spectral embedding has been applied to the following undirected weighted graph

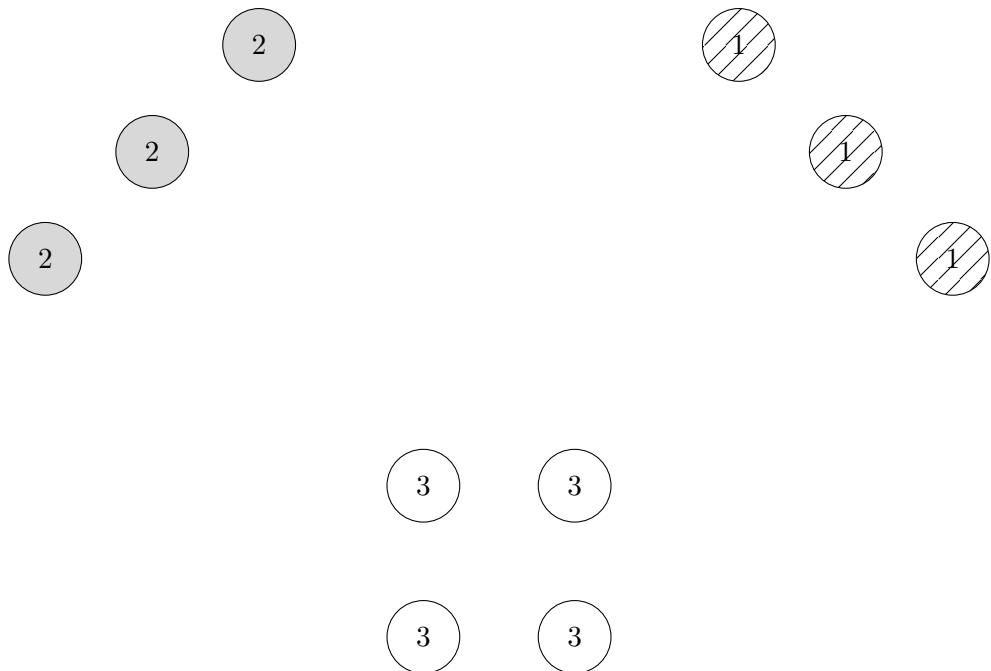


The plot below demonstrates the coordinates of the nodes in the embedding space (as defined by the second and the third eigenvectors of the unnormalized graph Laplacian).

Your task is to annotate the nodes in the plot below with their corresponding node IDs from the figure above.



Problem 5 [4 points] The graph below has been generated using the Stochastic block model with $K = 3$ communities (edges are hidden in the drawing). The community assignments of the 10 nodes in the graph are known (number inside each circle indicates the community ID z_i).



- (a) Draw the edges in the graph such that performing maximum likelihood estimation of the parameter η will yield

$$\boldsymbol{\eta} = \begin{bmatrix} 1 & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{4} & \frac{1}{3} & \frac{2}{3} \end{bmatrix}$$

In case you draw some edges incorrectly and need to correct your solution, please redraw the entire graph from scratch on pages 14-16. Note that there are multiple correct solutions and specifying any of them is enough.

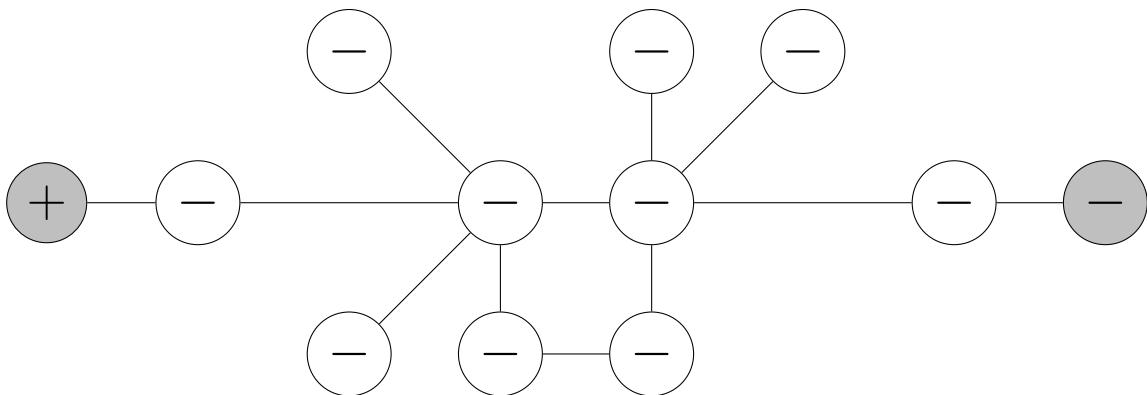
Provide a brief explanation below.

- (b) What is the maximum likelihood estimate of π for the given graph and communities?

4 Label Propagation

Problem 6 [4 points] You are given the graph below with two labeled seed nodes (shaded), belonging to two classes: + and -. Specify the labels of the remaining nodes (by writing + or - inside each node) that would be obtained with an **exact** solution of the standard binary label propagation problem (i.e. assuming label smoothness).

Justify your answer. Multiple correct solutions are possible, it is enough to provide one.



5 Deep Learning on Graphs

Problem 7 [12 points] You are given an undirected unweighted graph G with N nodes. Your task is to instantiate the differentiable message passing framework to compute PageRank. More specifically you have to specify:

- a) The input features $\mathbf{x}_v \in \mathbb{R}^2$ of each node v .
- b) The function M that computes the message from node u to its neighbor v .
- c) The function U that updates the hidden representation of a node v .

such that the hidden representations in the last layer approximately recover the PageRank score of the nodes in the graph G .

Hint 1: Increasing the number of layers in this formulation should improve the approximation of the PageRank scores. Hint 2: You do not need any trainable parameters.



