## Machine Learning for Graphs and Sequential Data Exercise Sheet 03

# VAE & GAN

**Problem 1:** Below we show pseudocode for implementing 3 autoencoder-like neural net architectures. The observed data is denoted as $\boldsymbol{x} \in \mathbb{R}^D$. Here, $g_{\boldsymbol{\lambda}} : \mathbb{R}^D \to \mathbb{R}^L$ and $f_{\boldsymbol{\psi}} : \mathbb{R}^L \to \mathbb{R}^D$ are fully connected feedforward neural networks with learnable parameters $\boldsymbol{\lambda}$ and $\boldsymbol{\psi}$. The output layers of $g_{\boldsymbol{\lambda}}$ and $f_{\boldsymbol{\psi}}$ have no (i.e. have linear) activation functions. $\mathcal{N}$ denotes the normal distribution, $\boldsymbol{I}_N$ is the $N \times N$ identity matrix, and $\boldsymbol{0}_N$ 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 $\boldsymbol{\lambda}$ and $\boldsymbol{\psi}$**. Answer "Yes" or "No" and provide a justification. If the answer is "Yes", modify the code to implement the reparametrization trick.

a) Model 1

$$\boldsymbol{z}_i \sim \mathcal{N}(\boldsymbol{x}_i, \boldsymbol{I}_D)$$
$$\boldsymbol{h}_i = g_{\boldsymbol{\lambda}}(\boldsymbol{z}_i)$$
$$\tilde{\boldsymbol{x}}_i = f_{\boldsymbol{\psi}}(\boldsymbol{h}_i)$$
$$\mathcal{L} = \|\boldsymbol{x}_i - \tilde{\boldsymbol{x}}_i\|_2^2$$

b) Model 2

$$\boldsymbol{h}_i = g_{\boldsymbol{\lambda}}(\boldsymbol{x}_i)$$
$$\boldsymbol{z}_i \sim \mathcal{N}(\boldsymbol{h}_i, \boldsymbol{I}_L)$$
$$\tilde{\boldsymbol{x}}_i = f_{\boldsymbol{\psi}}(\boldsymbol{z}_i)$$
$$\mathcal{L} = \|\boldsymbol{x}_i - \tilde{\boldsymbol{x}}_i\|_2^2$$

c) Model 3

$$\boldsymbol{h}_i = g_{\boldsymbol{\lambda}}(\boldsymbol{x}_i)$$
$$\boldsymbol{z}_i \sim \mathcal{N}(\boldsymbol{0}_L, \boldsymbol{I}_L)$$
$$\tilde{\boldsymbol{x}}_i = f_{\boldsymbol{\psi}}(\boldsymbol{h}_i + \boldsymbol{z}_i)$$
$$\mathcal{L} = \|\boldsymbol{x}_i - \tilde{\boldsymbol{x}}_i\|_2^2$$

**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.

$$\boldsymbol{h}_i = g_\lambda(\boldsymbol{x}_i)$$
$$\boldsymbol{z}_i \sim \mathcal{N}(\boldsymbol{0}_L, \text{diag}(\boldsymbol{h}_i))$$
$$\tilde{\boldsymbol{x}}_i = f_\psi(\boldsymbol{z}_i)$$
$$\mathcal{L} = \|\boldsymbol{x}_i - \tilde{\boldsymbol{x}}_i\|_2^2$$

**Problem 3:** The loss used in generative adversarial networks (GANs) can be written in the following form:
$$\min_{\boldsymbol{\theta}} \max_{\boldsymbol{\phi}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \min_{\boldsymbol{\theta}} \max_{\boldsymbol{\phi}} \mathbb{E}_{p^*(\boldsymbol{x})}[\log D_{\boldsymbol{\phi}}(\boldsymbol{x})] + \mathbb{E}_{p(\boldsymbol{z})}[\log(1 - D_{\boldsymbol{\phi}}(f_{\boldsymbol{\theta}}(\boldsymbol{z})))]$$

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

a) For a given generator (fixed parameters $\boldsymbol{\theta}$) assume there exists a discriminator $D_{\boldsymbol{\phi}^*}(\boldsymbol{x})$ with parameters $\boldsymbol{\phi}^*$ such that for all $\boldsymbol{x}$:
$$D_{\boldsymbol{\phi}^*}(\boldsymbol{x}) = \frac{p^*(\boldsymbol{x})}{p^*(\boldsymbol{x}) + p_{\boldsymbol{\theta}}(\boldsymbol{x})}$$

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

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

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

  – The generator is optimal i.e. $p_{\boldsymbol{\theta}}(\boldsymbol{x}) = p^*(\boldsymbol{x})$

  – The generator assigns a zero probability $p_{\boldsymbol{\theta}}(\boldsymbol{x}) = 0$ to a sample $\boldsymbol{x}$ whereas $p^*(\boldsymbol{x}) \neq 0$

  – The generator assigns a non-zero probability $p_{\boldsymbol{\theta}}(\boldsymbol{x}) \neq 0$ to a sample $\boldsymbol{x}$ whereas $p^*(\boldsymbol{x}) = 0$