

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 $\mathbf{x} \in \mathbb{R}^D$. 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, \mathbf{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

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

No, the samples do not depend on **learnable** parameters here.

b) Model 2

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

Yes, We need to replace line 2 with the following operations:

$$\begin{aligned} \boldsymbol{\epsilon}_i &\sim \mathcal{N}(\mathbf{0}_L, \mathbf{I}_L) \\ \mathbf{z}_i &= \mathbf{h}_i + \boldsymbol{\epsilon}_i \end{aligned}$$

c) Model 3

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

No, the samples do not depend on **learnable** parameters here.

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.

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

The covariance matrix of the normal distribution must be positive definite for the model to be well-defined. For the diagonal covariance, this means that all entries must be positive. We can ensure this by using the exponential function

$$\mathbf{h}_i = \exp(g_{\boldsymbol{\lambda}}(\mathbf{x}_i))$$

Since the samples \mathbf{z}_i depend on the parameters $\boldsymbol{\lambda}$ of the encoder, we need to use the reparametrization trick when sampling \mathbf{z}_i . For this, we replace line 2 with the following operations:

$$\begin{aligned} \boldsymbol{\epsilon}_i &\sim \mathcal{N}(\mathbf{0}_L, \mathbf{I}_L) \\ \mathbf{z}_i &= \boldsymbol{\epsilon}_i \odot \sqrt{\mathbf{h}_i} \end{aligned}$$

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

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

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

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

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

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

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

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

$$= \max_{D_\phi} \mathbb{E}_{p^*(\mathbf{x})}[\log D_\phi(\mathbf{x})] + \mathbb{E}_{p_\theta(\mathbf{x})}[\log(1 - D_\phi(\mathbf{x}))] \quad (2)$$

$$= \max_{D_\phi} \int [p^*(\mathbf{x}) \log D_\phi(\mathbf{x}) + p_\theta(\mathbf{x}) \log(1 - D_\phi(\mathbf{x}))] d\mathbf{x} \quad (3)$$

It's not obvious how we can find the discriminator D_{ϕ^*} that maximizes this integral. What if instead of maximizing the integral, we found the optimal discriminator **for every single \mathbf{x}** ? This would only allow us to achieve higher values of our optimization objective

$$\leq \int \max_{D_\phi} [p^*(\mathbf{x}) \log D_\phi(\mathbf{x}) + p_\theta(\mathbf{x}) \log(1 - D_\phi(\mathbf{x}))] d\mathbf{x} \quad (4)$$

Finding the optimal discriminator for specific fixed \mathbf{x} is easy. We use the formula from the problem statement, where we set $y = D_\phi(\mathbf{x})$, $a = p^*(\mathbf{x})$, and $b = p_\theta(\mathbf{x})$.

$$\begin{aligned} D_{\phi^*}(\mathbf{x}) &= \arg \max_{D_\phi(\mathbf{x})} (p^*(\mathbf{x}) \log D_\phi(\mathbf{x}) + p_\theta(\mathbf{x}) \log(1 - D_\phi(\mathbf{x}))) \\ &= \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + p_\theta(\mathbf{x})} \end{aligned} \quad (5)$$

That is, if our discriminator satisfies the property from Equation 5 for every \mathbf{x} , then it maximizes the expression under the integral sign in Equation 4 for every \mathbf{x} . Finally, if our discriminator D_{ϕ^*} maximizes the expression under the integral sign in Equation 4 for every \mathbf{x} , then it also maximizes the entire integral in Equation 3.

Putting everything together, we have shown, that the discriminator that satisfies Equation 5 for every point \mathbf{x} is the optimal discriminator for the GAN loss (assuming fixed generator parameters θ).

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

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

- If $p_\theta(\mathbf{x}) = p^*(\mathbf{x})$, then $D_{\phi^*}(\mathbf{x}) = \frac{1}{2}$
- If $p_\theta(\mathbf{x}) = 0$ and $p^*(\mathbf{x}) \neq 0$, then $D_{\phi^*}(\mathbf{x}) = 1$. The discriminator classifies all such samples as “Real”.
- If $p_\theta(\mathbf{x}) \neq 0$ and $p^*(\mathbf{x}) = 0$, then $D_{\phi^*}(\mathbf{x}) = 0$. The discriminator classifies all such samples as “Fake”.