



Homework 11  
 Due 5pm, Monday, June 10, 2024

**Problem 1: VLB for IWAE.** The standard variational lower bound (VLB) of VAE is

$$\log(p_\theta(x)) \geq \text{VLB}_{\theta,\phi}(x) = \mathbb{E}_{Z \sim q(z|x)} \left[ \log \left( \frac{p_\theta(x|Z)p_Z(Z)}{q_\phi(Z|x)} \right) \right],$$

where  $p_\theta(z|x)$  is the true posterior and  $q_\phi(z|x)$  is the approximate posterior. Define

$$\text{VLB}_{\theta,\phi}^{(K)}(x) = \mathbb{E}_{Z_1, \dots, Z_K \sim q_\phi(z|x)} \left[ \log \frac{1}{K} \sum_{k=1}^K \frac{p_\theta(x|Z_k)p_Z(Z_k)}{q_\phi(Z_k|x)} \right],$$

to be the VLB for importance weighted autoencoders (IWAE) [1]. To clarify,  $Z_1, \dots, Z_K$  are sampled independently from  $q_\phi(z|x)$ . Note that IWAE with  $K = 1$  coincides with the standard VAE, and  $\text{VLB}_{\theta,\phi}^{(1)} = \text{VLB}_{\theta,\phi}$ . Show:

- (a)  $\log p_\theta(x) \geq \text{VLB}_{\theta,\phi}^{(K)}(x)$  for all  $x$  and  $K \geq 1$ .
- (b) If  $K \geq M$ , then  $\text{VLB}_{\theta,\phi}^{(K)}(x) \geq \text{VLB}_{\theta,\phi}^{(M)}(x)$  for all  $x$ .
- (c) Let  $X_1, \dots, X_N$  be data for training the IWAE. Show that if  $q_\phi$  is “powerful enough”, then

$$\underset{\theta \in \Theta}{\text{maximize}} \sum_{i=1}^N \log p_\theta(X_i) = \underset{\theta \in \Theta, \phi \in \Phi}{\text{maximize}} \sum_{i=1}^N \text{VLB}_{\theta,\phi}^{(K)}(X_i).$$

What should be the precise meaning of “powerful enough”?

*Hint.* For (a), use the Jensen’s inequality. For (b), let  $I \subset \{1, \dots, K\}$  with  $|I| = M$  be a uniformly distributed subset of distinct indices from  $\{1, \dots, K\}$ . Then,  $\mathbb{E}_{I=\{i_1, \dots, i_M\}} \left[ \frac{a_{i_1} + \dots + a_{i_M}}{M} \right] = \frac{a_1 + \dots + a_K}{K}$  for any sequence of numbers  $a_1, \dots, a_K$ .

*Remark.* This analysis shows that  $\text{VLB}_{\theta,\phi}^{(K)}$  provides a tighter approximation of the log likelihood than  $\text{VLB}_{\theta,\phi}$ . However, using  $\text{VLB}_{\theta,\phi}^{(K)}$  requires more computation than  $\text{VLB}_{\theta,\phi}$ .

## References

- [1] Y. Burda, R. Grosse, and R. Salakhutdinov, Importance weighted autoencoders, *ICLR*, 2016.

**Problem 2: VAE with trainable prior.** In this problem, we consider the setup of training a VAE with a trainable prior. Specifically, we assume  $Z \sim r_\lambda(z)$ , where  $\lambda$  is a trainable parameter, and  $X \sim p_\theta(x | Z)$ . Let  $q_\phi(z | X)$  be the approximate posterior. Let

$$\text{VLB}_{\theta,\phi,\lambda}(X_i) = \mathbb{E}_{Z \sim q_\phi(z|X_i)} \left[ \log \left( \frac{p_\theta(X_i | Z) r_\lambda(Z)}{q_\phi(Z | X_i)} \right) \right].$$

- (a) Show that  $\log p_\theta(X_i) \geq \text{VLB}_{\theta,\phi,\lambda}(X_i)$ .
- (b) Describe how to evaluate stochastic gradients of  $\text{VLB}_{\theta,\phi,\lambda}(X_i)$  using the log-derivative trick.
- (c) Assume  $r_\lambda = \mathcal{N}(\lambda_1, \text{diag}(\lambda_2))$ , where  $\lambda_1, \lambda_2 \in \mathbb{R}^k$ ,  $q_\phi(z | X_i) = \mathcal{N}(\mu_\phi(X_i), \Sigma_\phi(X_i))$  with diagonal  $\Sigma_\phi$ , and  $p_\theta(X_i | z) = \mathcal{N}(f_\theta(z), \sigma^2 I)$ . Describe how to evaluate stochastic gradients of  $\text{VLB}_{\theta,\phi,\lambda}(X_i)$  using the reparameterization trick.

**Problem 3: Anomaly detection via flow models.** Assume we have a trained flow model that we use to evaluate the likelihood function  $p_\theta$ . In this problem, you will use this trained flow model to perform anomaly detection between the MNIST and KMNIST datasets. In step 1, load the MNIST and KMNIST datasets, and split the MNIST test dataset into “validation” and “test” sets. In step 2, define the flow model. In step 3, load the trained flow model. In step 4, calculate the mean and standard deviation of

$$\{\log p_\theta(Y_i)\}_{i=1}^M,$$

where  $Y_1, \dots, Y_M$  are the validation data. Define a threshold to be mean  $- 3$  standard deviations, and define inputs with log likelihood below this threshold to be anomalies. In step 5, check how many of the MNIST images within the test set are classified as anomalies and report the type I error rate. In step 6, check how many of the KMNIST images are classified as non-anomalies and report the type II error rate. Download the starter code `flow_anomaly.py`, which provides the implementation of steps 1–3. Complete the implementation of steps 4–6.

*Remark.* In this problem, we split the test data into validation and test sets because the entire training set was already used to train the flow model. If we were to train the flow model from scratch, it would be better to split the training set into the training and validation sets to set aside the validation data for step 3.

**Problem 4: Rock paper scissors and minimax optimization.** Consider a game of rock paper scissors between players  $A$  and  $B$ . Players  $A$  and  $B$  play randomized strategies with

$$p_A = \begin{bmatrix} \mathbb{P}(A \text{ plays rock}) \\ \mathbb{P}(A \text{ plays paper}) \\ \mathbb{P}(A \text{ plays scissors}) \end{bmatrix}, \quad p_B = \begin{bmatrix} \mathbb{P}(B \text{ plays rock}) \\ \mathbb{P}(B \text{ plays paper}) \\ \mathbb{P}(B \text{ plays scissors}) \end{bmatrix}.$$

Define

$$\Delta^3 = \{p = (p_1, p_2, p_3) \in \mathbb{R}^3 \mid p_1, p_2, p_3 \geq 0, p_1 + p_2 + p_3 = 1\}$$

so that  $p_A, p_B \in \Delta^3$ . In the game, a player receives 1 point for a win,  $-1$  points for a loss, and 0 points for a draw. Consider the minimax problem

$$\underset{p_A \in \Delta^3}{\text{minimize}} \quad \underset{p_B \in \Delta^3}{\text{maximize}} \quad \mathbb{E}_{p_A, p_B}[\text{points for } B].$$

(a) Show that

$$p_A^* = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}, \quad p_B^* = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}$$

is the unique solution to the minimax problem.

(b) Note that if  $p_B = (1/3, 1/3, 1/3)$ , then  $\mathbb{E}_{p_A, p_B}[\text{points for } B] = 0$  regardless of how  $A$  plays. Does this mean any strategy  $p_A \in \Delta^3$  is optimal for player  $A$ ? (Here, the word “optimal” is used informally. Think about whether any  $p_A \in \Delta^3$  is a best strategy for  $A$ .)

*Clarification.* We say  $(\theta^*, \phi^*)$  is a solution to the minimax problem

$$\underset{\theta \in \Theta}{\text{minimize}} \quad \underset{\phi \in \Phi}{\text{maximize}} \quad L(\theta, \phi)$$

if  $\theta^* \in \Theta$ ,  $\phi^* \in \Phi$ , and

$$L(\theta^*, \phi) \leq L(\theta^*, \phi^*) \leq L(\theta, \phi^*)$$

for all  $\theta \in \Theta$  and  $\phi \in \Phi$ , i.e., unilaterally deviating from  $\theta^*$  increases the value of  $L$  and unilaterally deviating from  $\phi^*$  decreases the value of  $L$ .

*Remark.* In the setup of GANs (which is what this problem is intended to prepare you for), if the generator is perfect, the discriminator cannot do better than a 50-50 guess in detecting fakes. However, the discriminator is still forced to learn to distinguish imperfect fakes, as otherwise, the generator can take advantage of the discriminator.