

In The Name of God

Sharif University of Technology  
Electrical Engineering Department

# Deep Generative Models

Assignment 2

Fall 2024

*Instructor: Dr. S. Amini*

Due on Aban 10, 1403 at 23:55



## 1 $D_{kl}$ intuitions

### 1.1 which one?

Consider you have a *target distribution*  $p(x | \theta)$  and an *approximate distribution*  $q(x | \phi)$ . Our objective is to optimize  $\phi$  to "fit"  $q(x | \phi)$  such that it becomes close to  $p(x | \theta)$ . A common dilemma arises: Should we solve

$$\arg \min_{\phi} \text{KL}(q(x | \phi), p(x | \theta)) \quad \text{or} \quad \arg \min_{\phi} \text{KL}(p(x | \theta), q(x | \phi))?$$

There is no one-size-fits-all answer to this question, as the optimal approach depends on the specific properties of  $p(x | \theta)$  and  $q(x | \phi)$ .

In this problem, you will be provided with several scenarios. For each one, choose one of the following options and provide a concise one-to-two-sentence justification for your choice, e.g. why a particular objective allows us to optimize  $\phi$  and how you plan to perform the optimization:

- A. We can minimize  $\text{KL}(q(x | \phi), p(x | \theta))$
- B. We can minimize  $\text{KL}(p(x | \theta), q(x | \phi))$
- C. We can minimize both  $\text{KL}(q(x | \phi), p(x | \theta))$  and  $\text{KL}(p(x | \theta), q(x | \phi))$
- D. None of the above, i.e. neither direction of KL divergence can be minimized.

The scenarios are listed below:

- **Scenario 1:** We can access samples and evaluate the exact density for both  $p(x | \theta)$  and  $q(x | \phi)$ .
- **Scenario 2:** We can access samples from  $p(x | \theta)$  but we do not have access to the density function. For  $q(x | \phi)$ , we do not have access to samples but we can evaluate the density.
- **Scenario 3:** We do not have access samples from  $p(x | \theta)$  and we only know its *unnormalized* density. However, we have access to both samples and (normalized) density function from  $q(x | \phi)$ .

**Hint 1:** You may not assume any closed-form solution for integral or expectation therefore you would need Monte Carlo estimation to estimate any expectations you come across.

**Hint 2:** You may use the log-derivative trick derived in Problem 2:

$$\nabla_{\phi} \mathbb{E}_{x \sim q(x|\phi)}[f(x; \phi)] = \mathbb{E}_{x \sim q(x|\phi)}[\nabla_{\phi} \log q(x | \phi) f(x)]$$

or you may use reparameterization trick.

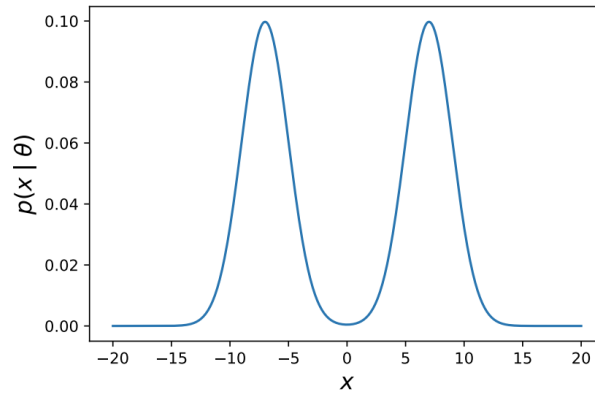


Figure 1: A two-mode Gaussian mixture

## 1.2 Mode-seeking or Mode-coverage

Recall the setting in previous problems, where you have a target distribution  $p(x | \theta)$  and an approximate distribution  $q(x | \phi)$ . In this problem, we will study the "mode-seeking" behavior of  $\text{KL}(q(x | \phi) \| p(x | \theta))$  and the "mode-coverage" behavior of  $\text{KL}(p(x | \theta) \| q(x | \phi))$ .

Now, assume we have  $p(x | \theta)$  as a bi-mode symmetric Gaussian mixture, e.g., the one illustrated in Fig. 1. And we have  $q(x | \phi)$  as a uni-variate Gaussian with learnable mean and variance. Now, please think about

- If we choose to minimize  $\text{KL}(q(x | \phi), p(x | \theta))$ , is there a global minima? If not, how many local minima would there exist? Where would the minima(s) be located?
- If we choose to minimize  $\text{KL}(p(x | \theta), q(x | \phi))$ , is there a global minima? If not, how many local minima would there exist? Where would the minima(s) be located?
- Based on your arguments from the (A) and (B), could you try to explain the meaning behind "mode-seeking" and "mode-coverage"?

## 2 General covariance matrix for VAE

In class, the VAE we considered having  $q(\mathcal{Z}|\phi)$  as a multivariate Gaussian with diagonal covariance, this is also known as mean-field Gaussian. Now, we will consider extending it to general covariance matrix. In particular, consider a random vector drawn from a Gaussian

$$\epsilon \sim \mathcal{N}(0, I),$$

which is then linearly transformed using the relation

$$z = \mu + L\epsilon$$

where  $L$  is a lower or (upper) triangular matrix. The off-diagonal elements define the correlations (covariances) of the elements in  $z$ . Show that  $z$  has a distribution  $\mathcal{N}(\mu, \Sigma)$  and write down an expression for  $\Sigma$  in terms of  $L$ . Explain why the diagonal elements of  $L$  must be positive. Describe how  $\mu$  and  $L$  can be expressed as the outputs of a neural network and discuss suitable choices of output-unit activation functions.

## 3 Gaussian Mixture latent model

In this question we will investigate how can we adapt VAE to have a Gaussian mixture embedding. the architecture of the model is in fig 2

This model generalizes the VAE model by using a Gaussian Mixture Model prior to replace the single Gaussian prior. This makes this network much more suitable for clustering tasks.

Below is an overview of how the network operates. Refer to the numbers in the figure 2.

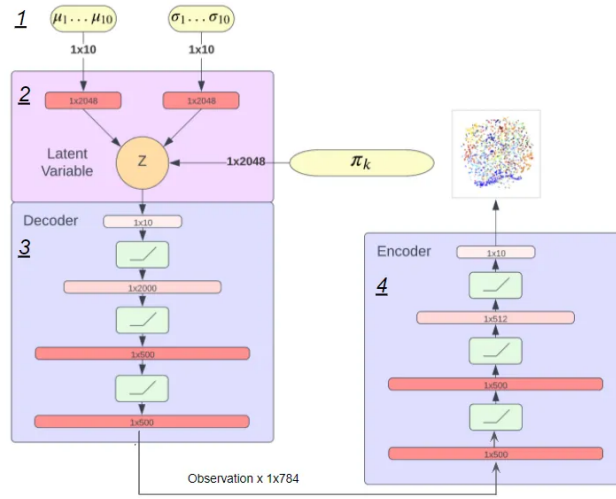


Figure 2: model architecture

1. A Cluster is selected from a Gaussian Mixture Model, and its mean and log variance is fed into the network
2. A latent embedding is generated based on the picked cluster
3. A DNN decodes the latent embedding into an observable  $x$
4. An Encoder network is used to maximize the ELBO

more precisely:

1. Choose a cluster  $c \sim \text{Cat}(\pi)$
2. Choose a latent vector  $z \sim \mathcal{N}(\mu_c, \sigma_c^2 \mathbf{I})$
3. Choose a sample  $x$ :
  - (a) If  $x$  is binary:
    - i. Compute the expectation vector  $\mu_x$ :

$$\mu_x = f(z; \theta)$$

- ii. Choose a sample  $x \sim \text{Ber}(\mu_x)$

- (b) If  $x$  is real-valued:

- i. Compute  $\mu_x$  and  $\sigma_x^2$ :

$$[\mu_x; \log \sigma_x^2] = f(z; \theta)$$

- ii. Choose a sample  $x \sim \mathcal{N}(\mu_x, \sigma_x^2 \mathbf{I})$

- $\pi_k$  is the prior probability for cluster  $k$
- $\text{Cat}(\pi)$  is the categorical distribution parameterized by  $\pi$
- $\mu_c$  and  $\sigma_c^2$  are the mean and variance of the Gaussian Distribution corresponding to cluster  $c$
- $f(z; \theta)$  is a neural network whose input is  $z$  and is parameterized by  $\theta$
- $\text{Ber}(\mu_x)$  and  $\mathcal{N}(\mu_x, \sigma_c^2)$  are multivariate Bernoulli and Gaussian Distributions.

### 3.1 Factorization

Using the processes above factorize the joint probability  $P(x, z, c)$  and explain how the network can compute each part of the factorized probability

### 3.2 ELBO

By maximizing  $p(x)$  prove that ELBO for this model will be

$$\mathcal{L}_{\text{ELBO}}(\mathbf{x}) = \mathbb{E}_{q(z, c|x)} [\log p(\mathbf{x}|z) + \log p(z|c) + \log p(c) - \log q(z|x) - \log q(c|x)]$$

note that we assumed a meanfield distribution for  $q(c, z|x)$

### 3.3 On derivation of $q(c_i | x)$

Prove that ELBO can be re expressed as below

$$\mathcal{L}_{\text{ELBO}}(\mathbf{x}) = \int_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \log \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q(\mathbf{z}|\mathbf{x})} d\mathbf{z} - \int_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) D_{KL}(q(\mathbf{c}|\mathbf{x}) \| p(\mathbf{c}|\mathbf{z})) d\mathbf{z}$$

Now explain how can we estimate  $q(c|x)$ , using the equation you derived?

### 3.4 A useful lemma for the next part

Prove that Given two multivariate Gaussian distributions  $q(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\sigma}}^2 \mathbf{I})$  and  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I})$ , we have:

$$\int q(\mathbf{z}) \log p(\mathbf{z}) d\mathbf{z} = \sum_{j=1}^J -\frac{1}{2} \log(2\pi\sigma_j^2) - \frac{\tilde{\sigma}_j^2}{2\sigma_j^2} - \frac{(\tilde{\mu}_j - \mu_j)^2}{2\sigma_j^2}$$

Where  $\mu_j, \sigma_j, \tilde{\mu}_j$  and  $\tilde{\sigma}_j$  simply denote the  $j^{\text{th}}$  element of  $\boldsymbol{\mu}, \boldsymbol{\sigma}, \tilde{\boldsymbol{\mu}}$  and  $\tilde{\boldsymbol{\sigma}}$ , respectively, and  $J$  is the dimensionality of  $\mathbf{z}$ .

### 3.5 Putting everything together

Prove that using the SGVB estimator and the reparameterization trick, along with the result of the previous parts the ELBO(x) for a bernouli  $x$  can be rewritten as

$$\begin{aligned} \mathcal{L}_{\text{ELBO}}(\mathbf{x}) = & \frac{1}{L} \sum_{l=1}^L \sum_{i=1}^D \left( x_i \log \mu_x^{(l)}|_i + (1 - x_i) \log(1 - \mu_x^{(l)}|_i) \right) \\ & - \frac{1}{2} \sum_{c=1}^K \gamma_c \sum_{j=1}^J \left( \log \sigma_c^2|_j + \frac{\tilde{\sigma}_j^2}{\sigma_c^2|_j} + \frac{(\tilde{\mu}_j - \mu_c|_j)^2}{\sigma_c^2|_j} \right) \\ & + \sum_{c=1}^K \gamma_c \log \frac{\pi_c}{\gamma_c} + \frac{1}{2} \sum_{j=1}^J (1 + \log \tilde{\sigma}_j^2) \end{aligned}$$

- $L$  is the number of Monte Carlo Samples in SGVB estimator
- $D$  is the dimensionality of  $\mathbf{x}$  and  $\mu_x$
- $x_i$  is the  $i^{\text{th}}$  element of  $\mathbf{x}$
- $J$  is the dimensionality of  $\mu_c$
- $\sigma_j^2, \mu_j, \tilde{\mu}_j, \tilde{\sigma}_j^2$  denote the  $j^{\text{th}}$  element of each variable
- $K$  is the number of clusters
- $\pi_c$  is the prior probability of cluster  $c$
- $\gamma_c$  denotes  $q(c|\mathbf{x})$

## 4 Vector Quantized Variational Autoencoders (VQ-VAE)

### \*(Optional)

Answer the below question.

1. Why VQ-VAE uses a discrete latent model instead of a continues one? give an example of a senorio where discrete latent is better that continues one?
2. How a discrete latent variable is modeled in VQ-VAE?
3. How gradient can be propagated to the input through the discrete latent variable?
4. How the dictionary of VQ-VAE updates?