

**Due Date: April 18th at 23:00**

Instructions

- For all questions that are not graded only on the answer, show your work! Any problem without work shown will get no marks regardless of the correctness of the final answer.
- Please try to use a document preparation system such as LaTeX. If you write your answers by hand, note that you risk losing marks if your writing is illegible without any possibility of regrade, at the discretion of the grader.
- Submit your answers electronically via the course GradeScope. Incorrectly assigned answers can be given 0 automatically at the discretion of the grader. To assign answers properly, please:
  - Make sure that the top of the first assigned page is the question being graded.
  - Do not include any part of the answer to any other questions within the assigned pages.
  - Assigned pages need to be placed in order.
  - For questions with multiple parts, the answers should be written in order of the parts within the question.
- Questions requiring written responses should be short and concise when necessary. Unnecessary wordiness will be penalized at the grader's discretion.
- Please sign the agreement below.
- It is your responsibility to follow updates to the assignment after release. All changes will be visible on Overleaf and Piazza.
- Any questions should be directed towards the TAs for this assignment (theoretical part): *Vitória Barin Pacela, Philippe Martin.*

**I acknowledge I have read the above instructions and will abide by them throughout this assignment. I further acknowledge that any assignment submitted without the following form completed will result in no marks being given for this portion of the assignment.**

Signature: \_\_\_\_\_

Name: \_\_\_\_\_

UdeM Student ID: \_\_\_\_\_

This assignment covers mathematical and algorithmic techniques in the families of deep generative models, such as VAEs, GANs, and contrastive learning.

**Question 1** (5-5-6). Consider a latent variable model  $p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$ , where  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$  and  $\mathbf{z} \in \mathbb{R}^K$ . The encoder network (aka “recognition model”) of a variational autoencoder,  $q_\phi(\mathbf{z}|\mathbf{x})$ , is used to produce an approximate (variational) posterior distribution over latent variables  $\mathbf{z}$  for any input data point  $\mathbf{x}$ .<sup>1</sup> This distribution is trained to match the true posterior by maximizing the evidence lower bound (ELBO):

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

Let  $\mathcal{Q}$  be the family of variational distributions with a feasible set of parameters  $\mathcal{P}$ ; i.e.  $\mathcal{Q} = \{q(\mathbf{z}; \pi) : \pi \in \mathcal{P}\}$ ; for example,  $\pi$  can be the mean and standard deviation of a normal distribution. We assume  $q_\phi$  is parameterized by a neural network (with parameters  $\phi$ ) that outputs the parameters,  $\pi_\phi(\mathbf{x})$ , of the distribution  $q \in \mathcal{Q}$ , i.e.  $q_\phi(\mathbf{z}|\mathbf{x}) := q(\mathbf{z}; \pi_\phi(\mathbf{x}))$ .

1.1 Show that maximizing the expected complete data log-likelihood (ECLL)

$$\mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})]$$

for a fixed  $q(\mathbf{z}|\mathbf{x})$ , w.r.t. the model parameter  $\theta$ , is equivalent to maximizing

$$\log p_\theta(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$$

This means the maximizer of the ECLL coincides with that of the marginal likelihood only if  $q(\mathbf{z}|\mathbf{x})$  perfectly matches  $p(\mathbf{z}|\mathbf{x})$ .

1.2 In this and the following task, the goal is to compare amortized variational inference (when  $q_\phi$  is optimized for the whole dataset) with the traditional variational inference (when  $q_\phi$  is optimized individually for each  $\mathbf{x}$ ). Consider a finite training set  $\{\mathbf{x}_i : i \in \{1, \dots, n\}\}$ ,  $n$  being the size the training data. Let  $\phi^*$  be the maximizer  $\arg \max_\phi \sum_{i=1}^n \mathcal{L}(\theta, \phi; \mathbf{x}_i)$  with  $\theta$  fixed. In addition, for each  $\mathbf{x}_i$  let  $q_i \in \mathcal{Q}$  be an “instance-dependent” variational distribution, and denote by  $q_i^*$  the maximizer of the corresponding ELBO. Compare  $D_{\text{KL}}(q_{\phi^*}(\mathbf{z}|\mathbf{x}_i) || p_\theta(\mathbf{z}|\mathbf{x}_i))$  and  $D_{\text{KL}}(q_i^*(\mathbf{z}) || p_\theta(\mathbf{z}|\mathbf{x}_i))$ . Which one is bigger?

1.3 Following the previous question, compare the two approaches in the second subquestion

- (a) in terms of bias of estimating the marginal likelihood via the ELBO, in the best-case scenario (i.e. when both approaches are optimal within the respective families)
- (b) from the computational point of view (efficiency)
- (c) in terms of memory (storage of parameters)

**Question 2** (3-3-7-3-2-2). The  $\beta$ -VAE [1] has been proposed as a follow-up from the VAE, and still is a relevant approach for learning a disentangled latent representation.

<sup>1</sup>Using a recognition model in this way is known as “amortized inference”; this can be contrasted with traditional variational inference approaches (see, e.g., Chapter 10 of Bishop’s *Pattern Recognition and Machine Learning*), which fit a variational posterior independently for each new data point.

The Beta-VAE introduces a hyperparameter  $\beta \in \mathbb{R}$ ,  $\beta \geq 1$ , which gives a higher weight to the regularization term (KL divergence between the real and the estimated posterior of  $z$ ) in the VAE, as opposed to the reconstruction error. The loss of a  $\beta$ -VAE is given by:

$$\mathcal{L}_{\theta, \phi, \beta}^{\beta \text{VAE}}(\mathbf{x}) = -\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{\theta}(\mathbf{x}|\mathbf{z}) + \beta \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z})) \quad (1)$$

where  $p_{\theta}(\mathbf{z})$  is an isotropic Gaussian.

- 2.1 This question is about the role of the reconstruction error. What would be the optimal solution when  $\beta = 0$ ?
- 2.2 This question is about the role of the KL divergence term. What would be the optimal solution if the regularization (KL divergence term) is optimized without any reconstruction error?
- 2.3 In a general autoencoder, let  $X \sim p(X)$  be the true generating process for the input data  $X$ ,  $Y = f_{\theta}(X)$  be the hidden variable, and  $p(X, Y)$  their associated joint.  
Prove that  $\max_{\theta} \mathbf{I}(X; Y) \geq \max_{\theta} \mathbf{H}(X) - l(X|Y)$ , where  $\mathbf{I}(X; Y)$  is the mutual information between  $X$  and  $Y$ ,  $\mathbf{H}(X)$  is the entropy of  $X$ , and  $l(X|Y)$  is the autoencoder loss function for binary inputs.
- 2.4 Propose a term for the loss that could replace the reconstruction error, considering the tradeoffs from the previous questions. (Hint: SSL-related)
- 2.5 Let  $z_a = f_a(x)$  and  $z_b = f_b(x)$  be two different representations of a latent factor under different parameterizations from the same observed data  $x$ .
  - (a) Give an example of  $z_a, z_b, f_a, f_b$  that would satisfy  $z_a \neq z_b$ . Consider the simplest case possible, where  $z_a, z_b \in \mathbb{R}^2$  and  $f_a, f_b$  are linear transformations.
  - (b) What does this suggest about the uniqueness of learning a latent variable from data  $x$ ? How would you mathematically characterize the entire set of  $z \in Z$  that could be reached from the same observed data  $x$ ?

**Question 3** (2-2-3). In this question, we are concerned with analyzing the training dynamics of GANs. Consider the following value function

$$V(d, g) = dg \quad (2)$$

with  $g \in \mathbb{R}$  and  $d \in \mathbb{R}$ . We will use this simple example to study the training dynamics of GANs.

1. Consider gradient descent/ascent with learning rate  $\alpha$  as the optimization procedure to iteratively minimize  $V(d, g)$  w.r.t.  $g$  and maximize  $V(d, g)$  w.r.t.  $d$ . We will apply the gradient descent/ascent to update  $g$  and  $d$  simultaneously. What is the update rule of  $g$  and  $d$ ? Write your answer in the following form

$$[d_{k+1}, g_{k+1}]^{\top} = A[d_k, g_k]^{\top}$$

where  $A$  is a  $2 \times 2$  matrix; i.e. specify the value of  $A$ .

- The optimization procedure you found in 6.1 characterizes a map which has a stationary point <sup>2</sup>, what are the coordinates of the stationary points?
- Analyze the eigenvalues of  $A$  and predict what will happen to  $d$  and  $g$  as you update them jointly. In other word, predict the behaviour of  $d_k$  and  $g_k$  as  $k \rightarrow \infty$ .

**Question 4** (3-5-4-3). Assume that we want to fit a data distribution  $p_d(\mathbf{x})$  with a parameterized density model defined as

$$p_\theta(\mathbf{x}) = \frac{\exp(-E_\theta(\mathbf{x}))}{Z_\theta},$$

where  $E_\theta(\mathbf{x})$  is the *energy* (i.e. an arbitrary function that maps  $\mathbf{x}$  to a scalar in  $\mathbb{R}$ ) parameterized with parameters  $\theta$  and  $Z_\theta$  is the normalizing constant:

$$Z_\theta = \int \exp(-E_\theta(x)) dx.$$

The issue with learning  $p_\theta(\mathbf{x})$  is the estimation of the normalizing constant  $Z_\theta$  which may involve an intractable integral. However, several methods exist for learning unnormalized probabilistic models, also known as Energy Based Models (EBMs). The questions in this problem will explore some methods related to learning EBMs.

- Noise Contrastive Estimation (NCE) is a method for learning the parameters of an EBM by contrasting it with another distribution with known density  $p_n(\mathbf{x})$ .  $p_n(\mathbf{x})$  is usually chosen to be simple with a tractable PDF that we can easily sample from (e.g. a Gaussian distribution). NCE defines the learning objective as a binary classification problem:

$$\mathcal{L}_{\text{NCE}}(E_\theta, p_d, p_n) := \mathbb{E}_{\mathbf{x} \sim p_d} [\log \sigma(E_\theta(\mathbf{x}))] + \mathbb{E}_{\mathbf{z} \sim p_n} [\log(1 - \sigma(E_\theta(\mathbf{z})))] ,$$

where  $\sigma$  is the sigmoid function. Assuming that  $E_\theta$  is expressive enough, give a close form solution of the optimal energy  $E_{\theta^*}$  as a function of  $p_n(\mathbf{x})$  and  $p_d(\mathbf{x})$ :

$$E_{\theta^*} = \arg \min_{E_\theta} \mathcal{L}_{\text{NCE}}(E_\theta, p_d, p_n). \quad (3)$$

- Contrastive learning also learns EBMs via a classification objective. It does so by contrasting *positive* and *negative* pairs of samples. In this question, we will assume that we are given a datasets of tuples representing a positive pair of samples  $(\mathbf{x}_i, \mathbf{x}'_i) \in \mathcal{D}$ . The negative pairs of samples are two elements  $(\mathbf{x}_i, \mathbf{x}'_j)$  sampled independently from  $\mathcal{D}$  (i.e.  $(\mathbf{x}_i, \cdot) \sim \mathcal{D}, (\cdot, \mathbf{x}'_j) \sim \mathcal{D}$ ). Let  $f_\theta : \mathcal{X} \rightarrow S^d$  and  $g_\phi : \mathcal{X}' \rightarrow S^d$  some mappings to the unit hyper-sphere, with  $\mathbf{z} = f_\theta(\mathbf{x})$  and  $\mathbf{y} = g_\phi(\mathbf{x}')$ , we define the contrastive objective as follows:

$$\mathcal{L}_{\text{contr.}}(f_\theta, g_\phi, \mathcal{D}) := \frac{1}{N} \sum_{i=1}^N \left[ -\log \frac{e^{\mathbf{z}_i^\top \cdot \mathbf{y}_i}}{\sum_{j=1}^M e^{\mathbf{z}_i^\top \cdot \mathbf{y}_j}} \right], \quad (4)$$

where  $M$  is the number of negative samples. Show that as  $M \rightarrow \infty$ , equation 4 can be decomposed into three terms:

$$\mathcal{L}_{\text{align}}(f_\theta, g_\phi, \mathcal{D}) := \frac{1}{N} \sum_{i=1}^N \mathbf{z}_i^\top \cdot \mathbf{y}_i,$$

<sup>2</sup>A stationary point is a point on the surface of the graph (of the function) where all its partial derivatives are zero (equivalently, the gradient is zero). Source: [https://en.wikipedia.org/wiki/Stationary\\_point](https://en.wikipedia.org/wiki/Stationary_point)

$$\mathcal{L}_{\text{unif}}(f_\theta, g_\theta, \mathcal{D}) := \frac{1}{N} \sum_{i=1}^N \log \frac{1}{M} \sum_{j=1}^M e^{\mathbf{z}_i^\top \mathbf{y}_j}$$

and a constant  $\log M$ . What does each term represent? What happens as  $M \rightarrow \infty$ ? Assume that the dimensionality of  $\mathcal{X}$  is greater than  $d$  the dimensionality of the hyper-sphere and the underlying density of the distribution generating  $\mathcal{D}$  is bounded<sup>3</sup>. Relate your answer to the information maximization principle.

3. The contrastive objective of equation 4 encodes positive and negative pairs of samples independently. Assume that we defined another network that jointly encode  $\mathbf{x}$  and  $\mathbf{x}'$ :  $h_\psi : \mathcal{X} \times \mathcal{X}' \rightarrow \mathcal{S}^d$ . For this question, assume that  $\mathbf{z} = h_\psi(\mathbf{x}, \mathbf{x}')$  and  $\mathbf{y} = g_\phi(\mathbf{x}')$ , where  $g_\phi$  is the same mapping as defined in the previous question. Derive an objective function  $\mathcal{L}_*(h_\psi, g_\phi, \mathcal{D})$  that have the same property as  $\mathcal{L}_{\text{contr.}}$  that we discussed in the previous question.
4. Assume that  $f_\theta$ ,  $g_\phi$  and  $h_\psi$  are all deep neural networks. Also, assume that we are training a vision-language encoders like CLIP. We can think of  $\mathbf{x}$  as image tokens and  $\mathbf{x}'$  as text tokens. In this case,  $f_\theta$  would be an image encoder,  $g_\phi$  would be a text encoder and  $h_\psi$  would be a joint image-text encoder. Discuss the trade-offs of learning an EBM with  $\mathcal{L}_{\text{contr.}}$  and  $\mathcal{L}_*$ .

## References

- [1] Higgins, Irina and Matthey, Loïc and Pal, Arka and Burgess, Christopher P. and Glorot, Xavier and Botvinick, Matthew M. and Mohamed, Shaker and Lerchner, Alexander, *beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework*, ICLR, 2017.

---

<sup>3</sup>These are details and are not necessary needed for a complete answer.