

Instructions

- For all questions, show your work!
- Use a document preparation system such as LaTeX.
- Submit your answers electronically via the course gradescope
- TA for this assignment is (theoretical part) : **Alexandra Volokhova** (IFT6135B) and **Ghait Boukachab** (IFT6135A).

Question 1 (2-2-4-2). 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 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 datapoint \mathbf{x} .¹ This distribution is trained to match the true posterior by maximizing the evidence lower bound (ELBO) :

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

We assume $q_\phi \in \mathcal{Q}$ where \mathcal{Q} is a parametric family, where we use ϕ to specify which member of the family we are using.

- 1.1 Show that data log likelihood $\log p_\theta(\mathbf{x})$ can be decomposed as a sum of ELBO and KL-divergence between variational and true posteriors over \mathbf{z} : $D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$
- 1.2 Show that maximizing ELBO w.r.t. ϕ is equivalent to minimizing KL-divergence between variational and true posteriors over \mathbf{z} w.r.t. ϕ .
- 1.3 In this and following task, the goal is to compare amortized variational inference (when q_ϕ is optimised for the whole dataset) with the traditional variational inference (when q_ϕ is optimised 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’s fix θ for simplicity. Let $q^* = \arg \max_{q_\phi \in \mathcal{Q}} \sum_{i=1}^n \mathcal{L}(\theta, \phi; \mathbf{x}_i)$ (i.e. q^* is the optimal variational distribution in the family \mathcal{Q} for given θ and training set). In addition, for each \mathbf{x}_i let $q_i^* = \arg \max_{q_\phi \in \mathcal{Q}} \mathcal{L}(\theta, \phi; \mathbf{x}_i)$. Compare $D_{\text{KL}}(q^*(\mathbf{z}|\mathbf{x}_i) || p_\theta(\mathbf{z}|\mathbf{x}_i))$ and $D_{\text{KL}}(q_i^*(\mathbf{z}|\mathbf{x}_i) || p_\theta(\mathbf{z}|\mathbf{x}_i))$. Which one is bigger?
- 1.4 Following the previous question, compare the two approaches in the second subquestion (justify the answers).
 - (a) which approach is better for estimating marginal likelihood via empirical ELBO
 - (b) which one is more computationally efficient
 - (c) which one is more memory efficient (storage of parameters)

Question 2 (5-2-7-2-2-5-2). In this task, we will go deeper into mathematics of diffusion models. Consider a denoising diffusion probabilistic model (DDPM) with the encoder process given by a linear Gaussian model : $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t|\sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t I)$, where $\beta_t \in (0, 1)$ is a fixed noise schedule. The forward diffusion process starts from initial image \mathbf{x}_0 from the dataset and ends at \mathbf{x}_T (T is a fixed number of steps). We assume $\mathbf{x}_T \sim p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T|0, I)$. The goal of the training is to learn a reversed (denoising process) $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$, which will allow to generate an image \mathbf{x}_0 starting from Gaussian noise \mathbf{x}_T .

1. 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 datapoint.

- 2.1 Given the equation for linear Gaussian encoder process $q(\mathbf{x}_t|\mathbf{x}_{t-1})$, show that the ground truth denoising process is

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}|\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t I) \quad (1)$$

where

$$\begin{aligned} \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) &= \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}x_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t \\ \tilde{\beta}_t &= \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t \\ \alpha_t &= 1 - \beta_t \\ \bar{\alpha}_t &= \prod_{s=1}^t \alpha_s \end{aligned} \quad (2)$$

If needed, you can use the following equation without proving it :

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t|\sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)I) \quad (3)$$

Hint : use Bayes rule and Markovian property of the encoder process

- 2.2 As we saw in task 2.1, it is possible to reverse the diffusion process analytically, without training anything. Explain, why we still need to train the reverse process $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ to generate images.
- 2.3 Now, let's derive the objective function for DDPM. Essentially, DDPM is a hierarchical variational autoencoder (with latent variables $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$) and its objective is an evidence lower bound (ELBO) for $\log p(\mathbf{x}_0)$. Show that

$$\log p(\mathbf{x}_0) \geq \mathcal{L}_{DDPM}(\theta; \mathbf{x}_0) = -L_0(\mathbf{x}_0) - \sum_{t=2}^T L_{t-1}(\mathbf{x}_0) - L_T(\mathbf{x}_0)$$

where

- (reconstruction term) $L_0(\mathbf{x}_0) = -\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0|\mathbf{x}_1)$
- (denoising matching term) $L_{t-1}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))$
- (prior matching term) $L_T(\mathbf{x}_0) = D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) || p(\mathbf{x}_T))$

The following equations might be useful for derivations :

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)q(\mathbf{x}_t|\mathbf{x}_0)}{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}$$

- 2.4 Which term in \mathcal{L}_{DDPM} doesn't affect optimisation over parameters and therefore can be excluded from the objective function ?
- 2.5 Compare ELBO for vanilla VAE (see previous task) and ELBO for DDPM. What is the key difference between them (in terms of trainable parameters) ?
- 2.6 Let's consider $L_{t-1}(\mathbf{x}_0)$ and $L_0(\mathbf{x}_0)$ in more detail.

- Using eq. 2 and 3, show that $\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon)$, where $\epsilon \sim \mathcal{N}(\epsilon|0, I)$

- A common parametrisation for denoising process is $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}|\mu_\theta(\mathbf{x}_t, t), \sigma_t^2 I)$, where mean of the Gaussian $\mu_\theta(\mathbf{x}_t, t)$ is trainable (here we consider σ_t^2 to be fixed for simplicity, while in practice it is trainable). However, instead of training a model to predict the $\mu_\theta(\mathbf{x}_t, t)$ directly, a common choice is to train a neural network ϵ_θ (aka "denoiser") to predict only the noise term : $\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}}\epsilon_\theta(\mathbf{x}_t, t))$. Show that

$$\mathbb{E}_{q(\mathbf{x}_0)} L_{t-1}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(0, I)} \left[\lambda_t \|\epsilon - \epsilon_\theta(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1-\alpha_t} \epsilon, t)\|^2 \right] + const \quad (4)$$

where $\lambda_t = \frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1-\alpha_t)}$ and $q(\mathbf{x}_0)$ is the groundtruth data distribution

Hint : you can use the equation for KL divergence between multivariate normal distributions without deriving it.

- Show that $\mathbb{E}_{q(\mathbf{x}_0)} L_0(\mathbf{x}_0)$ can be written in the same way as eq. 4

2.7 Finally, put together the equations for ELBO terms and get the DDPM loss function.

Question 3 (3-7). Let p_0 and p_1 be two probability distributions with densities f_0 and f_1 (respectively). We want to explore what we can do with a trained GAN discriminator. A trained discriminator is thought to be one which is "close" to the optimal one :

$$D^* := \arg \max_D \mathbb{E}_{\mathbf{x} \sim p_1} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_0} [\log(1 - D(\mathbf{x}))].$$

- 3.1 For the first part of this problem, derive an expression we can use to estimate the Jensen-Shannon divergence using a trained discriminator. We remind that the definition of JSD is $\text{JSD}(p_0, p_1) = \frac{1}{2} (KL(p_0 \parallel \mu) + KL(p_1 \parallel \mu))$, where $\mu = \frac{1}{2}(p_0 + p_1)$.
- 3.2 For the second part, we want to demonstrate that a optimal GAN Discriminator (i.e. one which is able to distinguish between examples from p_0 and p_1 with minimal NLL loss) can be used to express the probability density of a datapoint \mathbf{x} under f_1 , $f_1(\mathbf{x})$ in terms of $f_0(\mathbf{x})$ ². Assume f_0 and f_1 have the same support. Show that $f_1(\mathbf{x})$ can be estimated by $f_0(\mathbf{x})D(\mathbf{x})/(1 - D(\mathbf{x}))$ by establishing the identity $f_1(\mathbf{x}) = f_0(\mathbf{x})D^*(\mathbf{x})/(1 - D^*(\mathbf{x}))$.

Hint : Find the closed form solution for D^ .*

Question 4 (4-2-8-4-2). In this question, we will see why stop-gradient is critical for non-contrastive SSL methods like SimSiam and BYOL. We will show that removing stop-gradient results in collapsed representations, using the dynamics of SimSiam as our running example.

Consider a two-layer linear SimSiam model with the time-varying weight matrices given by $W(t) \in \mathbb{R}^{n_2 \times n_1}$ and $W_p(t) \in \mathbb{R}^{n_2 \times n_2}$. Note that $W(t)$ corresponds to the weights of the online **and** the target network, while $W_p(t)$ denotes the weights of the predictor. Let $\mathbf{x} \in \mathbb{R}^{n_1}$ be an input datapoint and $\mathbf{x}_1, \mathbf{x}_2$ be the two augmented versions of the input \mathbf{x} . Also note that in some instances, the dependence on time (t) is omitted for notational simplicity, and the weight matrices are referred to as W and W_p .

Let $\mathbf{f}_1 = W\mathbf{x}_1$ be the online representation of \mathbf{x}_1 and $\mathbf{f}_2 = W\mathbf{x}_2$ be the target representation of \mathbf{x}_2 . The learning dynamics of W and W_p can be obtained by minimizing SimSiam's objective function as shown below :

$$J(W, W_p) = \frac{1}{2} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [\|W_p \mathbf{f}_1 - \text{Stop-Grad}(\mathbf{f}_2)\|_2^2]. \quad (5)$$

2. You might need to use the "functional derivative" to solve this problem. See "19.4.2 Calculus of Variations" of the Deep Learning book or "Appendix D Calculus of Variations" of Bishop's Pattern Recognition and Machine Learning for more information.

4.1 Show (with proof) that the above objective can be simplified to :

$$J(W, W_p) = \frac{1}{2} [\text{tr}(W_p^T W_p F_1) - \text{tr}(W_p F_{12}) - \text{tr}(F_{12} W_p) + \text{tr}(F_2)], \quad (6)$$

where $F_1 = \mathbb{E} [\mathbf{f}_1 \mathbf{f}_1^T] = W(X + X')W^T$, $F_2 = \mathbb{E} [\mathbf{f}_2 \mathbf{f}_2^T] = W(X + X')W^T$ and $F_{12} = F_{21} = \mathbb{E} [\mathbf{f}_1 \mathbf{f}_2^T] = WXW^T$. Here, $X = \mathbb{E} [\bar{\mathbf{x}} \bar{\mathbf{x}}^T]$, where $\bar{\mathbf{x}}$ is the average augmented view of a datapoint \mathbf{x} and X' is the covariance matrix of augmented views \mathbf{x}' conditioned on \mathbf{x} and then averaged over the data \mathbf{x} , and tr is the Trace operation³.

4.2 Based on the above expression for $J(W, W_p)$, find the gradient update for W_p (the predictor network), denoting it as \dot{W}_p . In other words, obtain an expression for $\dot{W}_p = -\frac{\partial J}{\partial W_p}$ (the derivative of the objective function w.r.t the parameters W_p).

4.3 Consider the case when the Stop-Grad is removed. The gradient of the objective function $J(W, W_p)$ w.r.t the parameters W i.e. $\dot{W}(t) = -\frac{\partial J}{\partial W(t)}$, is given by :

$$\dot{W}(t) = \frac{d}{dt} \text{vec}(W(t)) = -H(t) \text{vec}(W(t)),$$

where $H(t)$ is a time-varying positive semi-definite matrix defined as

$$H(t) = X' \otimes (W_p(t)^T W_p(t) + I_{n_2}) + X \otimes (\tilde{W}_p(t)^T \tilde{W}_p(t)).$$

Here, \otimes is the Kronecker product⁴, $\tilde{W}_p(t) = (W_p(t) - I_{n_2})$, and "vec(W)" refers to the *vectorization* of a matrix W⁵. For simplicity, we are not taking weight decay into account here⁶.

If the minimal eigenvalue $\lambda_{\min}(H(t))$ is bounded away from zero, i.e. $\inf_{t \geq 0} \lambda_{\min}(H(t)) \geq \lambda_0 > 0$, then **prove that** $W(t) \rightarrow 0$.

Note : In order to prove the above question, the following property must be used :

For a time-varying positive definite matrix $H(t)$ whose minimal eigenvalues are bounded away from 0, the dynamics shown below :

$$\frac{d}{dt} \mathbf{w}(t) = -H(t) \mathbf{w}(t),$$

satisfies the constraint $\|\mathbf{w}(t)\|_2 = e^{-\lambda_0 t} \|\mathbf{w}(0)\|_2$, implying that $\mathbf{w}(t) \rightarrow 0$.

4.4 Consider the case when both the Stop-Grad **and** the predictor are removed. Show that the representations collapse i.e. $W(t) \rightarrow 0$. You may assume that X' is a positive definite matrix.

4.5 Speculate (in 1-2 sentences) as to why the stop-gradient and the predictor are necessary for avoiding representational collapse.

3. [https://en.wikipedia.org/wiki/Trace_\(linear_algebra\)](https://en.wikipedia.org/wiki/Trace_(linear_algebra)) [https://en.wikipedia.org/wiki/Trace_\(linear_algebra\)](https://en.wikipedia.org/wiki/Trace_(linear_algebra)).

4. For more information, see https://en.wikipedia.org/wiki/Kronecker_product#Matrix_equations https://en.wikipedia.org/wiki/Kronecker_product#Matrix_equations.

5. Also known as the "vec trick", it is obtained by stacking all the columns of a matrix A into a single vector.

6. Although omitted here, it must be noted that having weight decay is important. It has also been shown that, in practice, weight decay leads to stable learning.