

### Instructions

- Pour toutes les questions, montrez votre travail !
- Utiliser un système de rédaction de documents tel que *LaTeX*.
- Soumettez vos réponses par voie électronique via le système de notation du cours
- les TAs pour ce devoir sont (theoretical part) : **Alexandra Volokhova** (IFT6135B) and **Ghail Boukachab** (IFT6135A).

**Question 1** (2-2-4-2). Considérons un modèle de variable latente  $p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$ , où  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$  et  $\mathbf{z} \in \mathbb{R}^K$ . Le réseau encodeur (également appelé "modèle de reconnaissance") de l'autoencodeur variationnel,  $q_\phi(\mathbf{z}|\mathbf{x})$ , est utilisé pour produire une distribution postérieure approximative (variationnelle) sur les variables latentes  $\mathbf{z}$  pour tout point de données d'entrée  $\mathbf{x}$ .<sup>1</sup> Cette distribution est entraînée pour correspondre à la vraie postériorité en maximisant la limite inférieure de l'évidence (ELBO) :

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

On suppose que  $q_\phi \in \mathcal{Q}$  où  $\mathcal{Q}$  est une famille paramétrique, où nous indiquons  $\phi$  pour spécifier quel membre de la famille nous utilisons.

- 1.1 Montrer que la log-vraisemblance des données  $\log p_\theta(\mathbf{x})$  peut être décomposée en une somme d'ELBO et de divergence de KL entre les postérités variationnelles et réelles sur  $\mathbf{z}$  :  $D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$ .
- 1.2 Montrer que la maximisation de ELBO par rapport à  $\phi$  est équivalente à la minimisation de la divergence de KL entre les postérités variationnelles et réelles sur  $\mathbf{z}$  par rapport à  $\phi$ .
- 1.3 Dans cette sous-question et dans la suivante, l'objectif est de comparer l'inférence variationnelle armoisée (lorsque  $q_\phi$  est optimisé pour l'ensemble des données) avec l'inférence variationnelle traditionnelle (lorsque  $q_\phi$  est optimisé individuellement pour chaque  $\mathbf{x}$ ). Considérons un ensemble d'apprentissage fini  $\{\mathbf{x}_i : i \in \{1, \dots, n\}\}$ ,  $n$  étant la taille des données d'apprentissage. Fixons  $\theta$  pour plus de simplicité. Soit  $q^* = \arg \max_{q_\phi \in \mathcal{Q}} \sum_{i=1}^n \mathcal{L}(\theta, \phi; \mathbf{x}_i)$  (c'est-à-dire que  $q^*$  est la distribution variationnelle optimale dans la famille  $\mathcal{Q}$  pour un  $\theta$  et un ensemble d'apprentissage donnés). En outre, pour chaque  $\mathbf{x}_i$ , soit  $q_i^* = \arg \max_{q_\phi \in \mathcal{Q}} \mathcal{L}(\theta, \phi; \mathbf{x}_i)$ . Comparez  $D_{\text{KL}}(q^*(\mathbf{z}|\mathbf{x}_i) || p_\theta(\mathbf{z}|\mathbf{x}_i))$  et  $D_{\text{KL}}(q_i^*(\mathbf{z}|\mathbf{x}_i) || p_\theta(\mathbf{z}|\mathbf{x}_i))$ . Quel est le plus grand ?
- 1.4 Suite à la question précédente, comparez les deux approches dans la deuxième sous-question (justifiez les réponses).
  - (a) Quelle est la meilleure approche pour estimer la vraisemblance marginale via ELBO empirique ?
  - (b) Laquelle est la plus efficace en temps de calcul
  - (c) Laquelle est la plus efficace en termes de capacité de mémoire (stockage des paramètres)

**Question 2** (3\*-2-7-2-2-5-2). Dans cette question, nous allons creuser plus profondément dans les mathématiques des modèles de diffusion. Considérons un modèle probabiliste de diffusion à dispersion (DDPM) avec le processus encodeur donné par un modèle gaussien linéaire :  $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t | \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t I)$ , où  $\beta_t \in (0, 1)$  est un schéma fixe de bruits. Le processus de diffusion

1. L'utilisation d'un modèle de reconnaissance de cette manière est connue sous le nom "d'inférence amortie"; on peut l'opposer aux approches traditionnelles d'inférence variationnelle (voir, par exemple, le chapitre 10 de l'ouvrage de Bishop *Pattern Recognition and Machine Learning*), qui ajustent une postériorité variationnelle de manière indépendante pour chaque nouveau point de données.

"forward" commence par l'image initiale  $\mathbf{x}_0$  de l'ensemble de données et se termine à  $\mathbf{x}_T$  ( $T$  est un nombre fixe d'étapes). Nous supposons que  $\mathbf{x}_T \sim p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T|0, I)$ . L'objectif de l'entraînement est d'apprendre un processus inversé (processus de débruitage)  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ , qui permettra de générer une image  $\mathbf{x}_0$  à partir d'un bruit gaussien  $\mathbf{x}_T$ .

2.1 Etant donné l'équation du processus encodeur linéaire gaussien  $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ , montrez que le processus de débruitage de la "ground truth" est

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

où

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

Si nécessaire, vous pouvez utiliser l'équation suivante sans la prouver :

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

Indice : utiliser la règle de Bayes et la propriété markovienne du processus encodeur.

2.2 Comme nous l'avons vu dans la tâche 2.1, il est possible d'inverser le processus de diffusion analytiquement, sans entraînement. Expliquer pourquoi nous avons toujours besoin d'entraîner le processus inverse  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  pour générer des images.

2.3 Montrons maintenant la fonction objective du DDPM. Essentiellement, le DDPM est un auto-encodeur variationnel hiérarchique (avec des variables latentes  $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ ) et son objectif est une limite d'évidence (ELBO) pour  $\log p(\mathbf{x}_0)$ . Montrer que

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

où

- (terme de reconstruction)  $L_0(\mathbf{x}_0) = -\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0|\mathbf{x}_1)$
- (terme de correspondance pour le débruitage)  $L_{t-1}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))$
- (terme de correspondance préalable)  $L_T(\mathbf{x}_0) = D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) || p(\mathbf{x}_T))$

Les équations suivantes peuvent être utiles pour les dérivations :

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

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

2.4 Quel terme de  $\mathcal{L}_{DDPM}$  n'affecte pas l'optimisation des paramètres et peut donc être exclu de la fonction objective ?

2.5 Comparez ELBO pour vanilla VAE (voir question précédente) et ELBO pour DDPM. Quelle est la principale différence entre eux (en termes de paramètres entraînables) ?

2.6 Analysons  $L_{t-1}(\mathbf{x}_0)$  et  $L_0(x_0)$  plus en détail.

- En utilisant l'éq. 2 et 3, montrer que  $\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}}\boldsymbol{\epsilon})$ , où  $\boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon}|0, I)$
- Une paramétrisation courante du processus de débruitage est la suivante  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}|\mu_\theta(\mathbf{x}_t, t), \sigma_t^2 I)$ , où la moyenne de la gaussienne  $\mu_\theta(\mathbf{x}_t, t)$  est entraînable (nous considérons ici que  $\sigma_t^2$  est fixe pour des raisons de simplicité, alors qu'en pratique il est possible de l'entraîner). Cependant, au lieu d'entraîner un modèle pour prédire directement le  $\mu_\theta(\mathbf{x}_t, t)$ , un choix courant consiste à entraîner un réseau de neurones  $\boldsymbol{\epsilon}_\theta$  (également appelé "débruiteur") pour prédire uniquement le terme de bruit :  $\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t))$ . Montrer que

$$\mathbb{E}_{q(\mathbf{x}_0)} L_{t-1}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_0), \boldsymbol{\epsilon} \sim \mathcal{N}(0, I)} \left[ \lambda_t \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\alpha_t}x_0 + \sqrt{1-\alpha_t}\boldsymbol{\epsilon}, t)\|^2 \right] + const \quad (4)$$

où  $\lambda_t = \frac{\beta_t^2}{2\sigma_t^2\alpha_t(1-\alpha_t)}$  et  $q(\mathbf{x}_0)$  est la distribution des données "groundtruth". Astuce : vous pouvez utiliser l'équation de la divergence KL entre des distributions normales multivariées sans la déterminer.

- Montrer que  $\mathbb{E}_{q(\mathbf{x}_0)} L_0(\mathbf{x}_0)$  peut être écrit de la même manière que l'éq. 4

2.7 Enfin, rassemblez les équations pour les termes ELBO et obtenez la fonction de coût DDPM.

**Question 3** (3-7). Soit  $p_0$  et  $p_1$  deux distributions de probabilités avec les densités  $f_0$  et  $f_1$  (respectivement). Nous souhaitons explorer ce qu'on peut faire avec le discriminateur d'un GAN entraîné. Un discriminateur entraîné est considéré comme "proche" d'un discriminateur optimal :

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

3.1 Pour la première partie de ce problème, dérivez une expression permettant d'estimer la divergence de Jensen-Shannon (JSD) d'un discriminateur entraîné. Comme rappel, la JSD est  $\text{JSD}(p_0, p_1) = \frac{1}{2}(KL(p_0\|\mu) + KL(p_1\|\mu))$ , où  $\mu = \frac{1}{2}(p_0 + p_1)$ .

3.2 Pour la seconde partie, nous souhaitons démontrer qu'un discriminateur optimal d'un GAN (c.-à-d. un discriminateur qui peut distinguer des exemples provenant de  $p_0$  et de  $p_1$  avec une perte NLL minimale) peut être utilisé pour exprimer la densité de probabilité d'un exemple  $\mathbf{x}$  sous  $f_1$ ,  $f_1(\mathbf{x})$  en termes de  $f_0(\mathbf{x})$ <sup>2</sup>. Assumez que  $f_0$  et  $f_1$  ont le même support. Montrez que  $f_1(\mathbf{x})$  peut être estimé par  $f_0(\mathbf{x})D(\mathbf{x})/(1 - D(\mathbf{x}))$  en établissant l'identité  $f_1(\mathbf{x}) = f_0(\mathbf{x})D^*(\mathbf{x})/(1 - D^*(\mathbf{x}))$ .

*Astuce : Trouvez la solution analytique de  $D^*$ .*

**Question 4** (4-2-8-4-2). Dans cette question, nous allons voir que le *stop-gradient* est un élément critique pour les méthodes d'auto-supervision non contrastive comme SimSiam and BYOL. Nous allons démontrer que d'enlever le stop-gradient donne une représentation trivial en utilisant la dynamic de SimSiam comme exemple.

Considérez un réseau SimSiam à deux couches avec une matrice évoluant dans le temps donnée par  $W(t) \in \mathbb{R}^{n_2 \times n_1}$ . Notez que  $W(t)$  correspond au poids des réseaux online et target alors que

2. Il est possible que vous ayez à utiliser la *dérivée fonctionnelle* pour résoudre ce problème. Pour davantage d'information, voir "19.4.2 Calculus of Variations" du livre *Deep Learning* ou l'appendice "D. Calculus of Variations" du livre *Pattern Recognition and Machine Learning* de Bishop.

$W_p(t)$  dénote les pods du prédicteur. Soit  $\mathbf{x} \in \mathbb{R}^{n_1}$  une donnée d'entrée et  $\mathbf{x}_1, \mathbf{x}_2$  deux versions augmentées de  $\mathbf{x}$ . Notez aussi que dans certaines instances, la dépendances sur le temps (t) est omis pour simplifier la notation et que les matrices de poids sont référencées par  $W$  et  $W_p$ .

Soit  $\mathbf{f}_1 = W\mathbf{x}_1$  la représentation online de  $\mathbf{x}_1$  et  $\mathbf{f}_2 = W\mathbf{x}_2$  la représentation target de  $\mathbf{x}_2$ . Les dynamiques d'apprentissage de  $W$  et  $W_p$  sont obtenues en minimisant la fonction d'objectif de SimSiam comme montrée ci-dessous :

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

4.1 Montrez, avec preuve, que l'objectif ci-haut peut être simplifiée comme suit :

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

où  $F_1 = F_2 = \mathbb{E}[\mathbf{f}_1 \mathbf{f}_1^T] = W(X + X')W^T$  et  $F_{12} = F_{21} = \mathbb{E}[\mathbf{f}_1 \mathbf{f}_2^T] = WXW^T$ . Ici,  $X$  est la vue augmentée moyenne d'un point de donnée  $\mathbf{x}$  et  $X'$  est la matrice de covariance de vues augmentées  $\mathbf{x}'$  conditionnées sur  $\mathbf{x}$  et ensuite moyennée sur des données  $\mathbf{x}$  et  $\text{tr}$  est l'opération de trace<sup>3</sup>.

4.2 En vous basant sur l'expression ci-dessus pour  $J(W, W_p)$ , trouvez l'expression de mise-à-jour du gradient pour  $W_p$  (soit le réseau prédicteur). Autrement dit, obtenez une expression pour

$$\dot{W}_p = -\frac{\partial J}{\partial W_p} \quad (\text{la dérivée de la fonction d'objectif par rapport aux paramètres } W_p.)$$

4.3 Considérez le cas où le Stop-Grad est enlevé. Le gradient de la fonction d'objectif  $J(W, W_p)$  par rapport aux paramètres  $W$  (autrement dit  $\dot{W}(t) = -\frac{\partial J}{\partial W(t)}$  est donné par) :

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

où  $H(t)$  est une matrice semi-définie positive définie changeant dans le temps défini par

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

Ici,  $\otimes$  est le produit de Kronecker<sup>4</sup>,  $\tilde{W}_p = (W_p - I_{n_2})$ , et " $\text{vec}(W)$ " font références à la *vectorisation* de la matrice  $W$ <sup>5</sup>. Par simplicité, nous ne considérons pas le weight decay ici<sup>6</sup>.

Si la valeur propre minimal  $\lambda_{\min}(H(t))$  est bornée loin de zéro, i.e.  $\inf_{t \geq 0} \lambda_{\min}(H(t)) \geq \lambda_0 > 0$ , alors **prenez que**  $W(t) \rightarrow 0$ .

**Note :** Pour prouvez la question ci-dessus, la propriété suivante doit être utilisé :

Pour une matrice définitive positive variant dans le temps  $H(t)$  dont la valeur propre minimal est bornée loin de 0, la dynamique montrée ci-dessous :

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

satisfait la contrainte  $\|\mathbf{w}(t)\|_2 = e^{-\lambda_0 t} \|\mathbf{w}(0)\|_2$ , impliquant que  $\mathbf{w}(t) \rightarrow 0$ .

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. Pour plus d'information, voyez [https://en.wikipedia.org/wiki/Kronecker\\_product#Matrix\\_equations](https://en.wikipedia.org/wiki/Kronecker_product#Matrix_equations) [https://en.wikipedia.org/wiki/Kronecker\\_product#Matrix\\_equations](https://en.wikipedia.org/wiki/Kronecker_product#Matrix_equations)

5. Aussi connu comme le "vec trick", il est obtenue en empilant toutes les colonnes de la matrice  $A$  dans un seul vecteur.

6. Nous devons noter que l'utilisation du weight decay est importante. Il a été montré, un pratique, que weight decay mène à une apprentissage stable.

- 4.4 Considérez le cas où le Stop-Gradient **et** le prédicteur sont enlevés. Montrez que la représentation converge à la solution trivial, soit  $W(t) \rightarrow 0$ . Assumez que  $X'$  est une matrice définie positif.
- 4.5 Spéculez (en 1-2 phrases) pourquoi le stop-gradient et le prédicteur sont nécessaires pour éviter une représentation triviale.