

Due Date: April 22nd 23:59, 2019

Instructions

- *Lisez toutes les questions et instruction avant de commencer.*
- *Pour toutes questions, montrez votre travail!*
- *Soumettez votre code et votre rapport (pdf) électronique via la page du cours Gradescope. Si vous utilisez des notebook Jupyter, exportez les en pdf et soumettez les sur Gradescope.*
- *Vous pouvez mettre votre code dans un rpertoire Github et inclure le lien dans votre rapport!*

Problem 1

Les discriminateurs de modèles génératifs adversariaux (GAN) estiment une métrique entre deux distributions. Dans cette questions, vous allez implémenter un discriminateur qui estime le Jensen Shannon divergence (JSD) et un critic qui estime le Wasserstein distance (WD). Ensuite, en utilisant vos estimateurs, vous allez comparez les propriéts du JSD et de la WD.

Nous fournissons des samplers ¹ pour générer les différentes distributions que vous allez avoir besoin pour cette question. Vous pouvez utiliser n'importe quelle architecture que vous voulez, en autant que celles-ci ait assez de capacité pour estimer les distributions. Un MLP à trois couches devrait être suffisant. Vous pouvez utiliser SGD avec un learning rate de $1e-3$ et des mini-batch de taille 512.

1. Implémentez une fonction pour estimer le Jensen Shannon Divergence. Le JSD est défini comme $D_{JS}(p||q) = \frac{1}{2}D_{KL}(p||r) + \frac{1}{2}D_{KL}(q||r)$ où $r(x) = \frac{p+q}{2}$. Ainsi, la fonction d'objectif que votre réseau de neurone devrait optimiser est:

$$\arg \max_{\theta} \left\{ \log 2 + \frac{1}{2} \mathbf{E}_{x \sim p} [\log(D_{\theta}(x))] + \frac{1}{2} \mathbf{E}_{y \sim q} [\log(1 - D_{\theta}(y))] \right\}$$

2. Implémentez une fonction pour estimer le Wasserstein Distance. La forme dual du WD est défini comme: $W(p, q) = \sup_{\|T\|_L \leq 1} \mathbf{E}_{x \sim p} [T(x)] - \mathbf{E}_{y \sim q} [T(y)]$. Dans le but de contraindre votre fonction à être 1-lipschitz, nous recommandons d'utiliser la pénalité du gradient. Ainsi, l'objectif que votre réseau devra optimiser est

$$\arg \max_{\theta} \mathbf{E}_{x \sim p} [T_{\theta}(x)] - \mathbf{E}_{y \sim q} [T_{\theta}(y)] - \lambda \mathbf{E}_{z \sim r} [\|\nabla_z T_{\theta}(z) - 1\|_2].$$

r est la distributions sur $z = ax + (1-a)y$, où $x \sim p$, $y \sim q$ et $a \sim U[0, 1]$. Nous recommandons $\lambda \geq 10$.

3. Soit $Z \sim U[0, 1]$ une variable aléatoire suivant une distributions uniforme. Soit p une distribution défini comme $(0, Z)$ et q_{θ} la distributions défini comme (θ, Z) , où θ est un paramètre. Tracez une estimation du JSD et de la WD pour $\theta \in [-1, 1]$ avec des intervalles de 0.1 (c.-à-d. 21 points). L'axe des x devrait être la valeur de θ et l'axe des y devrait être votre estimation.

¹Voir le rpertoire du devoir IFT6135H19_assignment/assignment3/samplers.py

4. Soit f_0 la densité d'une distribution gaussienne standard une dimension, et f_1 une fonction de densité aléatoire donnée par la `distribution4`. Entraînez un discriminateur D_θ en maximisant la fonction de valeur

$$\mathbb{E}_{x \sim f_1} [\log D_\theta(x)] + \mathbb{E}_{y \sim f_0} [\log(1 - D_\theta(y))]$$

Estimez la densité f_1 en utilisant la procédure de la question 5 de la partie théorique. Tracez la sortie du discriminateur et l'estimation de la densité (en utilisant le script `density_estimation.py`).

Problem 2

Les Variational Autoencoders (VAEs) sont des modèles génératifs probabilistes. Ceci veut dire qu'ils sont utilisés pour estimer $p(\mathbf{x})$.

Entraînez une VAE sur le dataset *Binarised MNIST*, en utilisant le ELBO négatif comme vue en classe. La distribution a priori est $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$. Chaque pixel dans le dataset est binaire: Le pixel est soit noir, soit blanc; ceci veut dire que vous devez modéliser la probabilité $p(\mathbf{x}|\mathbf{z})$. En d'autres mots, le décodeur en sortie une distribution de Bernoulli (ceci veut dire que vous devez utiliser la perte d'entropie croisée binaire pour la reconstruction) sur les pixels. Ceci n'est pas compatible avec le dataset standard de MNIST disponible sur `torchvision`, donc soit (a) vous modifiez le code fournie avec le devoir ou, (b) vous téléchargez le dataset ici: <https://github.com/yburda/iwae/tree/master/datasets/BinaryMNIST>, et implémentez votre propre data loader.

Entraînez un VAE (10pts) Entraînez une VAE avec variable latente 100 dimensions. Ce qui suit est une suggestion d'hyperparamètres pour entraîner un VAE sur Binarized MNIST. Vous pouvez utiliser votre propre architecture:

Encoder (ϕ)

```
Conv2d(1, 32, kernel size=(3, 3))
ELU()
AvgPool2d(kernel size=2, stride=2)
Conv2d(32, 64, kernel size=(3, 3))
ELU()
AvgPool2d(kernel size=2, stride=2)
Conv2d(64, 256, kernel size=(5, 5))
ELU()
Linear(in features=256, out features=2 × 100)
(This final layer should output mean and log-variance)
```

Decoder (θ)

```
Linear(in features=100, out features=256)
ELU()
Conv2d(256, 64, kernel size=(5, 5), padding=(4, 4))
ELU()
UpsamplingBilinear2d(scale factor=2, mode=bilinear)
Conv2d(64, 32, kernel size=(3, 3), padding=(2, 2))
ELU()
UpsamplingBilinear2d(scale factor=2, mode=bilinear)
Conv2d(32, 16, kernel size=(3, 3), padding=(2, 2))
ELU()
Conv2d(16, 1, kernel size=(3, 3), padding=(2, 2))
```

Ceci n'est pas du code. Utilisez l'interface de votre framework favori pour implémenter ceci.

Utilisez ADAM avec un taux d'apprentissage de 3×10^{-4} , et entraînez pour 20 epochs. Évaluez le modèle sur un ensemble de validation en utilisant le **ELBO**. Aucun point ne sera donné ou enlevé si vous n'utilisez pas l'architecture donnée. Notez que pour cette question, vous devez:

1. Obtenir en moyenne un ELBO par instance de ≥ -96 sur l'ensemble de validation:

$$\frac{1}{|\mathcal{D}_{\text{valid}}|} \sum_{\mathbf{x}_i \in \mathcal{D}_{\text{valid}}} \mathcal{L}_{\text{ELBO}}(\mathbf{x}_i) \geq -96$$

2. Raportez le ELBO de votre modèle.

Vous êtes libre de modifier les hyperparamètres plus haut (sauf la taille de la variable latente) pour vous assurer que ceci fonctionne.

valuez le log-likelihood avec un VAE (20 pts) Les modèles VAE sont *valués* avec log-likelihood, approximations par importance sampling, qui a été couvert durant les lectures. La formule est reproduite ci-dessous avec plus de détails:

$$\log p(\mathbf{x} = \mathbf{x}_i) \approx \log \frac{1}{K} \sum_{k=1}^K \frac{p_{\theta}(\mathbf{x} = \mathbf{x}_i | \mathbf{z}_i^{(k)}) p(\mathbf{z} = \mathbf{z}_i^{(k)})}{q_{\phi}(\mathbf{z} = \mathbf{z}_i^{(k)} | \mathbf{x}_i)}; \quad \text{for all } k: \mathbf{z}_i^{(k)} \sim q_{\phi}(\mathbf{z} | \mathbf{x}_i)$$

et $\mathbf{x}_i \in \mathcal{D}$.

1. Avec M la taille de votre *minibatch* value, $K = 200$ le nombre d'échantillon utilisés pour votre importance sampling, D la dimension de votre input ($D = 784$ dans le cas de MNIST), et $L = 100$ la taille de votre variable latente, implémentez cette procédure d'importance sampling comme une fonction de:

- **donne:**

- votre modèle entrant la première partie de la question.
- (M, D) un array de \mathbf{x}_i 's.
- (M, K, L) un array de \mathbf{z}_{ik} 's.

- **retourne:**

- $(\log p(\mathbf{x}_1), \dots, \log p(\mathbf{x}_M))$ estimation de log-likelihood (M).

Montrez le fragment de code dans votre rapport.

Tip: Faites attention au débordement numérique² (probabilités d'une image deviennent très petites). Utilisez le truc du `LogSumExp`³ pour régler ce problème.

2. Rapportez votre évaluation de votre modèle entrant sur l'ensemble de validation et de test en utilisant, (a) le ELBO, et (b) l'estimation du log-likelihood ($\frac{1}{N} \sum_{i=1}^N \log p(\mathbf{x}_i)$), où N est la taille du dataset.

Assurez-vous de fournir votre code. Des points sont accordés pour l'implémentation.

²la génération d'un nombre trop petit pour être représenté sur le dispositif design l'emmagasine.

³<http://blog.smola.org/post/987977550/log-probabilities-semirings-and-floating-point>

Problem 3

Il y a présentement un intérêt grandissant pour la recherche en deep learning et vision par ordinateur pour la génération d'images. Un défi important pour la communauté de génération est la question d'évaluation. En information, nous voulons normalement des métriques objectives pour comparer nos algorithmes. Or, pour les tâches d'esthétiques, ce n'est pas clair si une telle métrique existe. Ainsi, valuer la qualité d'un modèle génératif demande souvent une combinaison d'évaluations qualitatives et quantitatives. Ainsi, l'évaluation de modèles génératifs est un domaine de recherche actif.

Dans cette question finale, nous explorons les avantages et les désavantages de deux modèles génératifs populaires en deep learning: GANs et VAEs, en utilisant des méthodes d'évaluations populaires de la littérature.

Street View House Numbers (SVHN) Le dataset SVHN⁴ est fait d'images venant de numéros de maisons house numbers de Google Street View. Comme MNIST, le but est de classifier (0-9), mais avec un dataset plus grand avec des images en couleurs.

Le dataset est disponible ici: http://ufldl.stanford.edu/housenumbers/train_32x32.mat (train set), et http://ufldl.stanford.edu/housenumbers/test_32x32.mat (test set). Splittez le training set en train et validation set comme vous voulez ⁵.

Alternativement, vous pouvez télécharger, splitter et itérateurs fournis avec ce devoir.

Generative models Pour ce devoir, vous devez implémenter les **deux** modèles suivants:

1. Variational Auto-Encoders (VAE):
Consiste d'un encodeur, $f(x)$, et d'un decoder $g_{\text{VAE}}(z)$ qui retourne la moyenne d'une fonction de probabilité gaussienne $p(x|z) = \mathcal{N}(x; g_{\text{VAE}}(z), \mathbf{I})$.
2. Generative Adversarial Networks (GANs):
Consiste d'un générateur $g_{\text{GAN}}(z)$, et d'un critic $d(x)$.

Notez que les deux modèles ont un $g(z)$, prenant un échantillon d'une distribution a priori (dans ce cas une distribution gaussienne, $\mathcal{N}(\mathbf{0}, \mathbf{I})$), et produit une image. Ceci est principalement ce que l'évaluation qualitative utilisera.

Hyperparameters & Training Pointers Entraînez les deux modèles avec variables latentes de dimension 100. Les deux modèles ont un générateur prenant un code latent et produit une image. Utilisez la même architecture pour les deux modèles génératifs

⁴<http://ufldl.stanford.edu/housenumbers/>

⁵Vous pouvez utiliser le **extra** test set disponible si vous le voulez, mais il n'y aura pas de points données ou enlevés si vous le faites.

Outre ces restrictions, vous êtes libre d'expérimenter avec différentes architectures et utilisez ce qui fonctionne. Une façon de vérifier est de regarder aux échantillons de $g(z)$. Ce qui suit sont des trucs pour entraîner ces modèles :

- Pour GANs, utilisez l'objectif WGAN-GP avec pénalité du gradient. Vous allez peut-être devoir ajuster le nombre de fois que vous mettez à jour les paramètres du discriminateur avant de mettre à jour ceux du générateur.
- Pour VAEs, faites attention avec la fonction d'activation que vous choisissez pour vous assurer que l'cartype est positif (exponentiel peut fonctionner, mais n'est peut-être pas optimale pour la stabilité).
- Plus important, il y a des métriques que vous pouvez surveiller lors de l'entraînement, comme la loss du critic pour le GAN et le KL divergence entre la distribution *priori* et l'approximation de la distribution *a posteriori* pour le VAE.

Qualitative Evaluation (25pts) Nous voulons valuer le “goodness” des modèles génératifs et si l'espace latent fait du sens. Par contre, ces concepts sont subjectifs et difficiles à quantifier. Ce qui suit sont des signaux que les auteurs incluent dans leur papier pour démontrer que leur modèle fait quelque chose de bien.

Pour les **deux** modèles :

1. **Fournissez des échantillons visuels.** Commentez, comparez et contrastez la qualité des échantillons de chaque modèle (e.g. floue, diversité).
2. **Nous voulons voir si le modèle a appris une représentation dml dans l'espace latent.** Échantillonnez un z aléatoire de votre distribution *priori*. Faites une petite perturbation de votre échantillon z pour *chaque dimension* (e.g. pour une dimension i , $z'_i = z_i + \epsilon$). ϵ doit être assez large pour voir des différences visuelles. Pour chaque dimension, observez si le changement résulte dans des variations visuelles (cela veut dire dans $g(z)$). Vous n'avez pas besoin de montrer toutes les dimensions, juste une couple qui résulte dans des changements intéressants.
3. **Comparez entre l'interpolation dans l'espace des données et l'espace latent.** Choisissez deux points aléatoires z_0 et z_1 dans l'espace latent échantillon de la distribution *priori*.
 - (a) Pour $\alpha = 0, 0.1, 0.2 \dots 1$ calculez $z'_\alpha = \alpha z_0 + (1 - \alpha) z_1$ et tracez les échantillons résultants $x'_\alpha = g(z'_\alpha)$.
 - (b) Utilisant les échantillons des données $x_0 = g(z_0)$ and $x_1 = g(z_1)$ et pour $\alpha = 0, 0.1, 0.2 \dots 1$ tracez les échantillons $\hat{x}_\alpha = \alpha x_0 + (1 - \alpha) x_1$.

Expliquez la différence entre les deux façons d'interpoler entre les images.

valuation quantitative (25pts) Certain des meilleures modèles génératifs n'ont pas de métrique d'évaluation tractable (ils fournissent seulement des échantillons). Plusieurs métriques sont utilisées dans la communauté de recherche pour évaluer les échantillons générés.

Le Frechet Inception Distance (FID) est une score souvent utilisé dans la littérature pour évaluer la génération de GANs. Essentiellement, elle calcule la distance l_2 entre le moment de premier et de deuxième ordre entre les représentations générées et la représentation cible. En pratique, cette représentation est une couche cachée d'un réseau de neurone pré-entraîné sur une tâche de classification.

Plus formellement, soit p une distribution cible et q la distribution venant du générateur's/décodateur's. (μ_p, Σ_p) et (μ_q, Σ_q) sont la moyenne et la covariance de la "représentations" de p et q , respectivement. Le FID est défini comme

$$d^2((\mu_p, \Sigma_p), (\mu_q, \Sigma_q)) = \|\mu_p - \mu_q\|_2^2 + \text{Tr}(\Sigma_p + \Sigma_q - 2(\Sigma_p \Sigma_q)^{1/2}) \quad (1)$$

Ici, nous voulons évaluer les échantillons générés par les GANs et les VAEs en utilisant le FID.

Dans le répertoire `assignment3` du répertoire du devoir, trouvez le script `score_fid.py`. Notre représentation est extraite en utilisant un classificateur custom entraîné sur des images de SVHN et leurs étiquettes respectives. Les TAs vous fournissent les modèles entraînés (`svhn_classifier.pt`), et les fonctions nécessaires pour charger les images (nos loaders chargent les images dans le même format que torchvision).

1. En utilisant l'équation 1, implémentez la fonction `calculate_fid_score(.,.)`, acceptant deux itérateurs en argument: un itérateur extrait les features des échantillons, et un itérateur des features extrait venant du test set.

La fonction devrait estimer $\mu_p, \Sigma_p, \mu_q, \Sigma_q$ et calculer $d^2((\mu_p, \Sigma_p), (\mu_q, \Sigma_q))$.

2. échantillonnez 1000 images de $g(z)$ venant du GAN et du VAE et mettez les dans le répertoire e.g. `path/to/sample_directory/samples/`. Assurez-vous qu'il y ait seulement le répertoire des `samples` dans le répertoire `path/to/sample_directory/`.

Alors, exécutez `python score_fid.py path/to/sample_directory/`. Rapportez le score que ce script retourne.