

Laissez des traces de votre démarche pour toutes les questions !
 For all questions, show your work ! Be as brief and specific as possible !

1. (10 point) Convolutional Neural Nets.

Consider a CNN with 6 layers applied to an input image of size 256×256 . For each layer we perform convolutions with kernels of size 5×5 , a 2×2 zero-padding and a stride of 1.

- (a) What are the dimensions of the feature maps for each layer ?
- (b) What is the effective receptive field for the neurons in the 6th layer (i.e. what dimensions of the input does a middle neuron in the feature map “see”)?
- (c) How much zero-padding would you require to make this a “full” convolution ?
- (d) How much zero-padding would you require to make this a “valid” convolution ?
- (e) How much zero-padding would you require to make this a “same” convolution ?

[FRANÇAIS] Considérez un CNN à six couches appliqué à une image de taille 256×256 . Pour chacune des couches, nous performons des convolutions avec des noyaux de taille 5×5 , un zéro padding de 2×2 et un ‘stride’ de 1.

- (a) Quelles sont les dimensions des *features maps* pour chacune des couches ?
- (b) Quel est le *effective receptive field* pour les neurones de la sixième couche (i.e. quelles sont les dimensions de l’image d’entrée qu’un neurone peut “voir”) ?
- (c) Combien de zéro padding auriez-vous besoin pour effectuer une *full convolution* ?
- (d) De combien de zéro padding auriez-vous besoin pour effectuer une *valid convolution* ?
- (e) Combien de zéro padding auriez-vous besoin pour effectuer une *same convolution* ?

2. (15 point) Recurrent Neural Networks.

Consider the behavior of a linear RNN :

$$\mathbf{h}_t = W\mathbf{h}_{t-1} + U\mathbf{x}_t + \mathbf{b},$$

where $\mathbf{h}_t \in \mathbb{R}^N$ and $\mathbf{x}_t \in \mathbb{R}^M$.

- (a) Write \mathbf{h}_t as a function of \mathbf{h}_0 .
- (b) Derive an expression for the Jacobian J :

$$\begin{bmatrix} \frac{\partial \mathbf{h}_{t,1}}{\partial \mathbf{h}_{0,1}} & \cdots & \frac{\partial \mathbf{h}_{t,1}}{\partial \mathbf{h}_{0,N}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{h}_{t,N}}{\partial \mathbf{h}_{0,1}} & \cdots & \frac{\partial \mathbf{h}_{t,N}}{\partial \mathbf{h}_{0,N}} \end{bmatrix}$$

- (c) What happens when $t \rightarrow \infty$? Under what conditions?

[FRANÇAIS] Considérez le RNN linéaire suivant :

$$\mathbf{h}_t = W\mathbf{h}_{t-1} + U\mathbf{x}_t + \mathbf{b},$$

où $\mathbf{h}_t \in \mathbb{R}^N$ et $\mathbf{x}_t \in \mathbb{R}^M$.

- (a) Exprimez \mathbf{h}_t en fonction de \mathbf{h}_0 .

- (b) Trouvez une expression pour la Jacobienne J :

$$\begin{bmatrix} \frac{\partial \mathbf{h}_{t,1}}{\partial \mathbf{h}_{0,1}} & \cdots & \frac{\partial \mathbf{h}_{t,1}}{\partial \mathbf{h}_{0,N}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{h}_{t,N}}{\partial \mathbf{h}_{0,1}} & \cdots & \frac{\partial \mathbf{h}_{t,N}}{\partial \mathbf{h}_{0,N}} \end{bmatrix}$$

- (c) Que ce passe-t-il lorsque $t \rightarrow \infty$? Sous quelles conditions?

3. (10 point) Optimization.

Briefly compare and contrast the following pairs of optimization algorithms (be specific) :

- (a) Adagrad and RMSprop,
- (b) RMSprop and Adam.

[FRANÇAIS] Comparez et contrastez brièvement les paires d'algorithmes d'optimisation suivantes (soyez spécifique!) :

- (a) Adagrad et RMSprop,
- (b) RMSprop et Adam.

4. (10 point) Regularization.

- (a) Consider training a neural network with stochastic gradient descent. Explain why the capacity of a neural network grow as the number of training iterations increases?

[FRANÇAIS] Considérez un réseau de neurones entraîné avec la descente de gradient stochastique. Expliquez pourquoi la capacité du modèle augmente avec chaque itération d'entraînement.

- (b) Consider an alternative regularization method to dropout where, for a given example, instead of using a fixed probability of dropping out a neuron from the network, we *learn* the dropout probability. More specifically, we imagine neuron i as having its own dropout probability parameter p_i and we update these dropout parameters via gradient descent (assuming we have solved the problem of estimating gradients for this parameter) in order to minimize the training loss.

- (i) Is this likely that this method would be an effective regularizer? What pitfalls do you foresee, if any?
- (ii) Above we mentioned that we have solved the problem of estimating gradients for the dropout parameters. What is the problem?

[FRANÇAIS] Considérez une méthode de régularisation similaire à *dropout*, où, pour un exemple donné, nous *apprenons* la probabilité de *dropout* au lieu d'utiliser une valeur fixe. Plus précisément, un neurone i aurait sa propre probabilité de *dropout* p_i . Ces paramètres seraient mis-à-jour par descente de gradient (en assumant que nous savons comment calculer le gradient de ces paramètres) afin de minimiser le coût d'entraînement.

- (i) Est-il probable que cette méthode soit un régularisateur efficace? Le cas échéant, qu'est-ce qui pourrait mal se passer?
- (ii) Plus haut, nous avons assumé que nous savions comment calculer le gradient des paramètres *dropout* p_i . Quel problème devons-nous surmonter afin d'y arriver?

5. (10 point) CAE and Weight Decay

- (a) For the case of a linear autoencoder with a squared error loss function, prove that the contractive autoencoder (CAE) penalty is equivalent to weight decay. Specifically, for the linear auto-encoder of data \mathbf{x} , assume the reconstruction is $\tilde{\mathbf{x}} = W^\top \mathbf{h}(\mathbf{x})$, where $\mathbf{h}(\mathbf{x}) = W\mathbf{x}$ and the CAE loss function is : $L = (\mathbf{x} - \tilde{\mathbf{x}})^\top (\mathbf{x} - \tilde{\mathbf{x}}) + \lambda \|\nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x})\|_F^2$.

[FRANÇAIS] Prouvez que dans le cas d'un auto-encodeur linéaire accompagné d'une fonction de coût quadratique, la pénalité contractante est équivalente au *weight decay*. Dans votre réponse, présumez que la reconstruction de \mathbf{x} est donnée par $\tilde{\mathbf{x}} = W^\top \mathbf{h}(\mathbf{x})$, où $\mathbf{h}(\mathbf{x}) = W\mathbf{x}$, et que la fonction de coût du CAE est $L = (\mathbf{x} - \tilde{\mathbf{x}})^\top (\mathbf{x} - \tilde{\mathbf{x}}) + \lambda \|\nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x})\|_F^2$.

- (b) With the equivalence above thus established, describe if or how the CAE penalty is distinct from standard weight decay when they are both applied to an autoencoder with a single hidden layer of rectified linear units (ReLus).

[FRANÇAIS] En présument l'équivalence ci-haut, décrivez si ou comment la pénalité contractante diffère du *weight decay* conventionnel lorsqu'appliquée à un auto-encodeur avec une couche cachée dont la fonction d'activation utilisée est le *rectified linear* (ReLus).

6. (15 point) VAE.

- (a) Consider you have a latent variable model over observations \mathbf{x} parametrized by the conditional distribution $p(\mathbf{x} | \mathbf{z})$ and prior over the latent variables $\mathbf{z} : p(\mathbf{z})$. Imagine we want to approximate the intractable posterior distribution $p(\mathbf{z} | \mathbf{x})$ with a variational approximation $q(\mathbf{z} | \mathbf{x})$. Prove that the gap between the log probability (density) and the variational lower bound (ELBO) is the KL divergence between the true posterior and the approximation $q(\mathbf{z} | \mathbf{x})$. Specifically, prove :

$$\log p(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z} | \mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z} | \mathbf{x})} \right] - \mathbb{E}_{q(\mathbf{z} | \mathbf{x})} \left[\log \frac{p(\mathbf{z} | \mathbf{x})}{q(\mathbf{z} | \mathbf{x})} \right] \quad (1)$$

[FRANÇAIS] Considérez que vous avez un modèle à variable latente qui modélise les observations \mathbf{x} par une distribution conditionnel $p(\mathbf{x} | \mathbf{z})$ ainsi qu'un prior sur les variables latentes \mathbf{z} : $p(\mathbf{z})$. Imaginez que nous voulons estimer la distribution postérieure $p(\mathbf{z} | \mathbf{x})$ avec une approximation variationnelle. Prevez que la différence entre la log probabilité (densité) et la *variational lower bound* (ELBO) est la KL divergence entre la vrai distribution postérieure et l'approximation $q(\mathbf{z} | \mathbf{x})$. Plus précisément, prouvez que :

$$\log p(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z} | \mathbf{x})} \right] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{z} | \mathbf{x})}{q(\mathbf{z} | \mathbf{x})} \right] \quad (2)$$

- (b) As discussed in class, the traditional mean-field variational method corresponds to factorizing the variational approximation to the posterior distribution as a product of distributions : $q^{mf}(z_i) = \prod_j \mathcal{N}(z_{i,j} | m_{i,j}, \sigma_{i,j}^2)$ and maximizing the lower bound directly with respect to the variational parameters $m_{i,j}$ and $\sigma_{i,j}^2$ for each example separately (i.e. no encoder network is learned). Could the Inverse Auto-regressive Flow (IAF) approach to inference in the VAE outperform the mean-field method ? (Explain your answer.)

[FRANÇAIS] Comme discuté en classe, la méthode *mean-field* correspond à factoriser l'approximation variationnelle de la distribution postérieure comme étant un produit de distribution : $q^{mf}(z_i) = \prod_j \mathcal{N}(z_{i,j} | m_{i,j}, \sigma_{i,j}^2)$. La borne inférieure est ainsi maximiser par rapport aux paramètres variationnelles : $m_{i,j}$ et $\sigma_{i,j}^2$ pour chacun des exemples (il n'y a donc aucun encodeur !). Est-ce qu'utiliser l'approche *Inverse Auto-regressive Flow (IAF)* pour inférer nos valeurs nous permettrait d'avoir de meilleurs résultats que la méthode *mean-field* ? Expliquer votre réponse.

7. (15 points) GANs

For input $\mathbf{x} \in \mathbb{R}^N$ Consider we have a linear regression model as a very simple GAN discriminator :

$$D(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + \mathbf{b}. \quad (3)$$

- (a) Would it be advisable to use such a simple discriminator ? Justify your answer by describing what problem / benefits you might expect to encounter.
- (b) For the loss used in the WGAN-GP seen in class :

$$L = \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g}[D(\tilde{\mathbf{x}})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r}[D(\mathbf{x})] + \lambda \mathbb{E}_{\hat{\mathbf{x}} \sim \mathbb{P}_{\hat{\mathbf{x}}}[(\|\nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}})\|_2 - 1)^2]} \quad (4)$$

where $\hat{\mathbf{x}}$ is interpolated between the generated $\tilde{\mathbf{x}}$ and a true data sample \mathbf{x} . For this simple GAN discriminator above, please provide the discriminator gradient update for the model parameters \mathbf{w} : $\nabla_{\mathbf{w}} L$.

- (c) How does PacGAN help avoid the mode collapse problem that can plague standard GAN training ?

[FRANÇAIS] Pour $\mathbf{x} \in \mathbb{R}^N$, considérez que nous avons un modèle de régression linéaire comme discriminateur pour un GAN :

$$D(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + \mathbf{b}. \quad (5)$$

- (a) Est-ce une bonne idée d'avoir un modèle aussi simple comme discriminateur ? Justifiez votre réponse en décrivant les problèmes/avantages que vous pourriez rencontrer.

- (b) Voici la fonction de coût du WGAN-GP que nous avons vue en classe :

$$L = \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g}[D(\tilde{\mathbf{x}})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r}[D(\mathbf{x})] + \lambda \mathbb{E}_{\hat{\mathbf{x}} \sim \mathbb{P}_{\hat{x}}}[(\|\nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}})\|_2 - 1)^2] \quad (6)$$

où $\hat{\mathbf{x}}$ est une interpolation entre un exemple généré $\tilde{\mathbf{x}}$ et un vrai exemple \mathbf{x} . Pour le discriminateur linéaire décrit plus haut, donnez le gradient par rapport aux paramètres \mathbf{w} : $\nabla_{\mathbf{w}} L$.

- (c) Comment est-ce que *PacGAN* peut nous aider à éviter le problème de *mode collapse* qui arrive souvent lors de l'entraînement de GAN classique ?

8. (15 point) Generative models

Consider you are training a Boltzmann machine over observations $\mathbf{x} \in \{0,1\}^N$ with two sets of latent variables $\mathbf{y} \in \{0,1\}^M$ and $\mathbf{z} \in \{0,1\}^K$. The joint probability is parametrized as follows :

$$P(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \frac{1}{Z} \exp(-E(\mathbf{x}, \mathbf{y}, \mathbf{z})) \quad (7)$$

where

$$Z = \sum_{\mathbf{x}} \sum_{\mathbf{y}} \sum_{\mathbf{z}} \exp(-E(x, y, z)) \quad (8)$$

and

$$E(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \sum_{ijk} W_{ijk} x_i y_j z_k \quad (9)$$

- (a) Describe a block-Gibbs sampling strategy that will efficiently sample from sets of independent variables conditionally on the others. Derive and state clearly all required conditionals.
- (b) Is the computation of the unnormalized marginal probability in \mathbf{x} ($\log P(\mathbf{x}) + \log Z$) linear in the number of latent variables, as is the RBM ?
- (c) Is it possible to define an efficient training strategy based on a Contrastive Divergence-like learning algorithm ? Explain your answer.

[FRANÇAIS] Considérez que vous entraîné une *Boltzmann machine* sur des observations $\mathbf{x} \in \{0,1\}^N$ avec deux ensembles de variables latentes : $\mathbf{y} \in \{0,1\}^M$ and $\mathbf{z} \in \{0,1\}^K$. La probabilité jointe est paramétrisée comme suit :

$$P(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \frac{1}{Z} \exp(-E(\mathbf{x}, \mathbf{y}, \mathbf{z})) \quad (10)$$

où

$$Z = \sum_{\mathbf{x}} \sum_{\mathbf{y}} \sum_{\mathbf{z}} \exp(-E(\mathbf{x}, \mathbf{y}, \mathbf{z})) \quad (11)$$

et

$$E(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \sum_{ijk} W_{ijk} x_i y_j z_k \quad (12)$$

- (a) Décrivez une stratégie d'échantillonnage de *blocs-Gibbs* qui échantillonne à partir d'ensembles de variables conditionnellement indépendantes les unes des autres. Dérivez et énoncez clairement les distributions conditionnelles requises.
- (b) Est-ce que le calcul de la probabilité marginale (non-normalisée) dans \mathbf{x} ($\log P(\mathbf{x}) + \log Z$) est linéaire en fonction du nombre de variables latentes, comme pour les RBMs ?
- (c) Est-il possible de définir une stratégie d'entraînement efficace qui s'inspire de l'algorithme *Contrastive Divergence* ? Expliquez votre réponse.