

Date: 10/04/2019. Time/Temps: 12h30 - 14h30

Instructions

- For all long-form questions, show your work!
- Montrez les traces de votre démarche pour toutes les questions longues!

Question 1 [ENGLISH] True/False (30pts, 2 points each/chacun)

- (a) As the capacity of neural network increases, we expect the training error to increase.
- (b) In training a neural network, if training set size increases, we would expect the difference between the training and generalization error to decrease.
- (c) Convolutioning a feature map of size $(32, 32)$ with a 3×4 kernel, a stride of 2 without zero padding yields a feature map of size $(15, 15)$.
- (d) The RMSProp optimizer uses momentum to accelerate learning.
- (e) The Adam optimizer adapts the learning rate by using the running average of the elementwise square of gradient to estimate its second moment.
- (f) Weight decay has no impact on the training of a neural network with Batch Normalization.
- (g) Maximizing the ELBO wrt to the encoder $q(\mathbf{z}|\mathbf{x})$ is equivalent to minimizing the KL divergence $D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$.
- (h) PixelCNNs are efficient to train because they can exploit parallelism across pixels despite being autoregressive generative models.
- (i) Reparameterization trick is used to reduced the bias in approximating the true posterior $p(\mathbf{z}|\mathbf{x})$.
- (j) The softmax activation function is shift-invariant.
- (k) For some choice of the proposal distribution $q(\mathbf{h})$, the following inequality holds

$$\mathbb{E}_{\mathbf{h}} \left[\log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h})} \right] < \mathbb{E}_{\mathbf{h}_1, \dots, \mathbf{h}_K} \left[\frac{1}{K} \sum_{j=1}^K \log \frac{p(\mathbf{x}, \mathbf{h}_j)}{q(\mathbf{h}_j)} \right]$$

where $\mathbf{h}, \mathbf{h}_1, \dots, \mathbf{h}_K$ are independently and identically distributed by $q(\mathbf{h})$ and $K > 1$.

- (l) The optimal discriminator of a vanilla GAN is the density function of the real data distribution.
- (m) The Jensen Shannon divergence between two distributions is always $\log 2$ whenever they have disjoint support.
- (n) The Wasserstein-GAN requires regularizing the generator network to be 1-Lipschitz.
- (o) An "adversarial example" is the name given to GAN generated examples that successfully fool the GAN discriminator.

Question 1 [FRANÇAIS] Vrai / Faux: (30pts, 2 points chacun)

- (a) Lorsque la capacité d'un réseau de neurone augmente, nous nous attendons à ce que l'erreur d'entraînement augmente.
- (b) Dans l'entraînement d'un réseau de neurone, si la taille de l'ensemble d'entraînement augmente, la différence entre les erreurs d'entraînement et de généralisation devrait diminuer.
- (c) Effectuer une convolution sur un *feature map* de taille $(32, 32)$ avec un noyau de taille 3×4 , un *stride* de 2 et sans *zero padding* produit un *feature map* de taille $(15, 15)$.
- (d) L'optimiseur RMSProp utilise le momentum pour accélérer l'apprentissage.
- (e) L'optimiseur Adam adapte le taux d'apprentissage en utilisant la moyenne courante (*running average*) du gradient au carré par élément pour estimer le moment d'ordre deux.
- (f) Le *weight decay* n'a pas d'impact sur l'entraînement de réseau de neurone avec *Batch Normalization*.
- (g) Maximiser le *ELBO* par rapport à l'encodeur $q(\mathbf{z}|\mathbf{x})$ est équivalent à minimiser la divergence KL $D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$.
- (h) Il est efficace d'entraîner des *PixelCNNs*, car nous pouvons exploiter le parallélisme au travers des pixels, même si ce dernier est un modèle génératif autorégressif.
- (i) La technique de reparametrization (*Reparameterization trick*) est utilisée pour réduire le biais dans l'approximation de la vraie probabilité à postériori $p(\mathbf{z}|\mathbf{x})$.
- (j) La fonction d'activation Softmax est invariante au décalage.
- (k) Pour un choix de distribution $q(\mathbf{h})$, l'inégalité suivante est valable

$$\mathbb{E}_{\mathbf{h}} \left[\log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h})} \right] < \mathbb{E}_{\mathbf{h}_1, \dots, \mathbf{h}_K} \left[\frac{1}{K} \sum_{j=1}^K \log \frac{p(\mathbf{x}, \mathbf{h}_j)}{q(\mathbf{h}_j)} \right]$$

où $\mathbf{h}, \mathbf{h}_1, \dots, \mathbf{h}_K$ sont indépendamment et identiquement distribuée par $q(\mathbf{h})$ et $K > 1$.

- (l) Le discriminateur optimal du GAN originale est la fonction de densité des données de la distribution des vraies données.
- (m) La divergence de Jensen Shannon entre deux distributions est toujours $\log 2$ lorsque leur support est disjoint.
- (n) Le generateur du Wasserstein-GAN doit être régularisé pour être 1-Lipschitz.
- (o) Un "*adversarial example*" est le nom donné aux exemples générés par le générateur d'un GAN qui trompent le discriminateur d'un GAN.

Question 2 [ENGLISH] Short answer questions (*20pts, 2 points each*)

- (a) How does one diagnose overfitting?
- (b) Explain why the capacity of a neural network grows as the number of training iterations increases.
- (c) Very briefly, explain the difference between AdaGrad and RMSprop.
- (d) Applying early stopping to a linear model is equivalent to what form of regularization?
- (e) Compute the full convolution (with kernel flipping) for the following 1D matrices $[1, 2, 3] * [1, 0, 1]$ (Hint: in $f * g$, the kernel is g).
- (f) Compute the valid convolution (with kernel flipping) for the above (i.e. for $[1, 2, 3] * [1, 0, 1]$).
- (g) The standard ReLU hidden unit is given by the following

$$a = \mathbf{w}^T \mathbf{x} + b, \quad h = \text{ReLU}(a)$$

Specify the changes to the pre-activation function a when we apply Batch Normalization.

- (h) The BERT transformer model was trained on two self-supervised tasks, name or describe these.
- (i) Given an encoder $f : \mathcal{X} \rightarrow \mathcal{H}$ and a decoder $g : \mathcal{H} \rightarrow \mathcal{X}$, write out the loss function of a contractive autoencoder. Use L_2 reconstruction loss.
- (j) Describe the Meta-Learning evaluation setting, i.e. describe the composition of the meta-test set and specify how we measure the Meta-Learning model performance.

Question 2 [FRANÇAIS] Questions à réponse courte. (*20pts, 2 points chaque*)

- (a) Comment diagnostiquer le surapprentissage (*Overfitting*)?
- (b) Expliquez pourquoi la capacité d'un réseau de neurone augmente avec le nombre d'itération d'entraînement.
- (c) Expliquez très brièvement la différence entre AdaGrad et RMSprop.
- (d) Quel type de régularisation est équivalent à appliquer la technique d'arrêt précoce (*early stopping*) à un modèle linéaire?
- (e) Calculez la convolution complète (*full convolution*) avec retournement de noyau (*kernel flipping*) pour les matrices 1D suivante $[1, 2, 3] * [1, 0, 1]$
- (f) Calculez la convolution valid (*valid convolution*) avec retournement de noyau (*kernel flipping*) pour les matrices de la sous-question précédente (c.-à-d. pour $[1, 2, 3] * [1, 0, 1]$).
- (g) L'unité ReLU standard est donnée par ce qui suit

$$a = \mathbf{w}^T \mathbf{x} + b, \quad h = \text{ReLU}(a)$$

Spécifiez les changements à la fonction de pré-activation a lorsque nous appliquons *Batch Normalization*.

- (h) Le model transformeur BERT est entraîné sur deux tâches auto-supervisées. Nommez ou décrivez celles-ci.
- (i) Soit l'encodeur $f : \mathcal{X} \rightarrow \mathcal{H}$ et décodeur $g : \mathcal{H} \rightarrow \mathcal{X}$, écrivez la fonction de coût d'un auto-encodeur contractif (*contractive autoencoder*). Utilisez la fonction de coût de reconstruction L_2 .
- (j) Décrivez le cadre d'évaluation utilisé en Meta-Learning, c.-à-d. décrivez la composition du *meta-test set* et spécifiez comment nous mesurons la performance des modèles de Meta-Learning.

Question 3 [ENGLISH] Supervised neural nets (10pts)

Let \mathbf{x} be an n -dimensional vector. Recall that the softmax function $S : \mathbf{x} \in \mathbb{R}^n \mapsto S(\mathbf{x}) \in (0, 1)^n$ is defined as $S(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$.

- (a) Let \mathbf{x} be a function of vector \mathbf{u} . Show that the gradient $\nabla_{\mathbf{u}} \log S(\mathbf{x}(\mathbf{u}))_i$ is equal to

$$\nabla_{\mathbf{u}} \mathbf{x}(\mathbf{u})_i - \mathbb{E}_j[\nabla_{\mathbf{u}} \mathbf{x}(\mathbf{u})_j]$$

where j is a random index following a categorical distribution with probability $S(\mathbf{x}(\mathbf{u}))_j$.

- (b) Let \mathbf{y} and \mathbf{x} be K -dimensional vectors related by $\mathbf{y} = S(\mathbf{x})$. Use the fact you found in the previous question to derive the gradient of the cross-entropy loss (i.e. negative log likelihood) with respect to the input of the softmax, $\nabla_{\mathbf{u}} L(\mathbf{x}, \mathbf{c})$, where \mathbf{c} is a one-hot vector corresponding to the class label (i.e. a vector of all zeros except for a 1 in the position associated with the correct class):

$$L(\mathbf{x}, \mathbf{c}) = \sum_{i=1}^K -c_i \log y_i$$

Question 3 [FRANÇAIS] Réseaux de neurones supervisés (10pts)

Soit \mathbf{x} un vecteur à n -dimensions. Rappelez-vous que la fonction softmax $S : \mathbf{x} \in \mathbb{R}^n \mapsto S(\mathbf{x}) \in (0, 1)^n$ est définie comme $S(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$.

- (a) Soit \mathbf{x} fonction du vecteur \mathbf{u} . Montrez que le gradient $\nabla_{\mathbf{u}} \log S(\mathbf{x}(\mathbf{u}))_i$ est égal à

$$\nabla_{\mathbf{u}} \mathbf{x}(\mathbf{u})_i - \mathbb{E}_j[\nabla_{\mathbf{u}} \mathbf{x}(\mathbf{u})_j]$$

où j est un indice aléatoire suivant une distribution catégorique avec probabilité $S(\mathbf{x}(\mathbf{u}))_j$.

- (b) Soit \mathbf{y} et \mathbf{x} des vecteurs à n -dimensions en relation par $\mathbf{y} = S(\mathbf{x})$. Utilisez le résultat obtenu à la question précédente pour dériver le gradient de *cross-entropy loss* (c.-à.-d. le *negative log likelihood*) par rapport à l'entrée du softmax, $\nabla_{\mathbf{u}} L(\mathbf{x}, \mathbf{c})$, où \mathbf{c} est un vecteur *one-hot* correspondant à l'étiquette de la classe (c.-à.-d. un vecteur de 0 partout à l'exception de la position associée à la bonne classe qui est représentée par un 1):

$$L(\mathbf{x}, \mathbf{c}) = \sum_{i=1}^K -c_i \log y_i$$

Question 4 [ENGLISH] RNNs and gradient penalty regularization (20pts)

Denote by σ the logistic sigmoid function. Consider the following RNN:

$$\begin{aligned}\mathbf{h}_t &= \sigma(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1}) \\ y_t &= \mathbf{v}^\top \mathbf{h}_t\end{aligned}$$

Here, each $\mathbf{x}_t \in \mathbb{R}^d$, $\mathbf{h}_t \in \mathbb{R}^m$ and $y_t \in \mathbb{R}$.

Let z_t be the true target of the prediction y_t . Consider a regularized L_2^2 -based loss function (per time-step). The form of the regularization we will consider is inspired from the WGAN-GP and we will use it to encourage smoothness from \mathbf{x}_t to y_t . Specifically, $L = \sum_t L_t$ where $L_t = (z_t - y_t)^2 + \|\nabla_{\mathbf{x}_t} y_t\|_2^2$.

- (a) Draw the computational graph for forward propagation of this RNN, unrolled for 3 time steps (from $t = 1$ to $t = 3$). Include and label the initial hidden states for the RNN: h_0 .
- (b) Express the gradient $\nabla_{\mathbf{h}_t} L$ recursively in terms of $\nabla_{\mathbf{h}_{t+1}} L$.
- (c) Derive a simplified expression for the gradient $\nabla_{\mathbf{x}_t} y_t$.
- (d) Derive simplified expressions for $\nabla_{\mathbf{v}} L$, $\nabla_{\mathbf{W}} L$ and $\nabla_{\mathbf{U}} L$ in terms of known quantities.

Question 4 [FRANÇAIS] RNNs régularisation par pénalité de gradient (20pts)

Dénotez par σ la fonction logistique sigmoid. Considérez le RNN suivant:

$$\begin{aligned}\mathbf{h}_t &= \sigma(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1}) \\ y_t &= \mathbf{v}^\top \mathbf{h}_t\end{aligned}$$

Ici, chaque $\mathbf{x}_t \in \mathbb{R}^d$, $\mathbf{h}_t \in \mathbb{R}^m$ et $y_t \in \mathbb{R}$.

Soit z_t la vraie cible de la prédiction de y_t . Considérez une fonction de coût L_2^2 -based régularisée (par pas de temps). La forme de régularisation que nous allons considérée est inspirée par WGAN-GP. Nous allons l'utiliser pour encourager le lissage (*smoothness*) de \mathbf{x}_t à y_t . Spécifiquement, $L = \sum_t L_t$ où $L_t = (z_t - y_t)^2 + \|\nabla_{\mathbf{x}_t} y_t\|_2^2$.

- (a) Dessinez le graph de calcul de la propagation avant pour ce RNN, déroulé sur 3 pas de temps (de $t = 1$ à $t = 3$). Incluez et étiquettez l'état caché initial du RNN: h_0 .
- (b) Exprimez le gradient $\nabla_{\mathbf{h}_t} L$ récursivement en terme de $\nabla_{\mathbf{h}_{t+1}} L$.
- (c) Dérivez une expression simplifiée pour le gradient $\nabla_{\mathbf{x}_t} y_t$.
- (d) Dérivez une expression simplifiée pour $\nabla_{\mathbf{v}} L$, $\nabla_{\mathbf{W}} L$ and $\nabla_{\mathbf{U}} L$ en terme de quantités connues.

Question 5 [ENGLISH] VAEs (20pts)

Consider a latent variable model $\mathbf{z} \sim p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ where $\mathbf{z} \in \mathbb{R}^K$, and $\mathbf{x} \sim p_\theta(\mathbf{x}|\mathbf{z})$. The encoder network of variational autoencoder, $q_\phi(\mathbf{z}|\mathbf{x})$, is used to produce an approximate (variational) posterior distribution over the latent variables \mathbf{z} for any input datapoint \mathbf{x} .

- (a) Prove that the log-likelihood of the data $\log p_\theta(\mathbf{x})$ can be expressed as the sum of the evidence lower bound, $\mathcal{L}[q_\phi]$, and the KL divergence between $p_\theta(\mathbf{z} | \mathbf{x})$ and $q_\phi(\mathbf{z} | \mathbf{x})$.
- (b) Decompose the variational gap (i.e. the KL) into the approximation gap and amortization gap, and explain what they are.
- (c) How might one reduce this approximation gap?
- (d) Given K i.i.d. samples drawn from $q_\phi(\mathbf{z} | \mathbf{x})$, how would you estimate the variational gap? This method should be more accurate if a larger number of samples K is used.

Question 5 [FRANÇAIS] VAEs (20pts)

Considérez un modèle à variables latentes $\mathbf{z} \sim p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_K)$ où $\mathbf{z} \in \mathbb{R}^K$, et $\mathbf{x} \sim p_\theta(\mathbf{x}|\mathbf{z})$. Le réseau encodeur de l'encodeur variationnel (*variational autoencoder*), $q_\phi(\mathbf{z}|\mathbf{x})$, est utilisé pour produire une distribution à postériori approximée (*variational*) sur les variables latentes \mathbf{z} pour n'importe quel point de donnée en entrée \mathbf{x} .

- (a) Prouvez que le *log-likelihood* des données $\log p_\theta(\mathbf{x})$ peut être exprimé comme une somme du *evidence lower bound*, $\mathcal{L}[q_\phi]$, et de la divergence KL entre $p_\theta(\mathbf{z} | \mathbf{x})$ et $q_\phi(\mathbf{z} | \mathbf{x})$.
- (b) Décomposez le gap variationnel (*variational gap*) (c.-à-d. le KL) en gap d'approximation et en gap d'amortization. Expliquez ce qu'ils sont.
- (c) Comme est-il possible de réduire ce gap d'approximation?
- (d) Soit K échantillons tiré i.i.d. de $q_\phi(\mathbf{z} | \mathbf{x})$, comment estimeriez-vous le gap variationnel? Cette méthode est plus précise si un grand nombre K d'échantillons est utilisé.