

Deep Learning — Final Exam

Date: Thursday, September 1, 2022

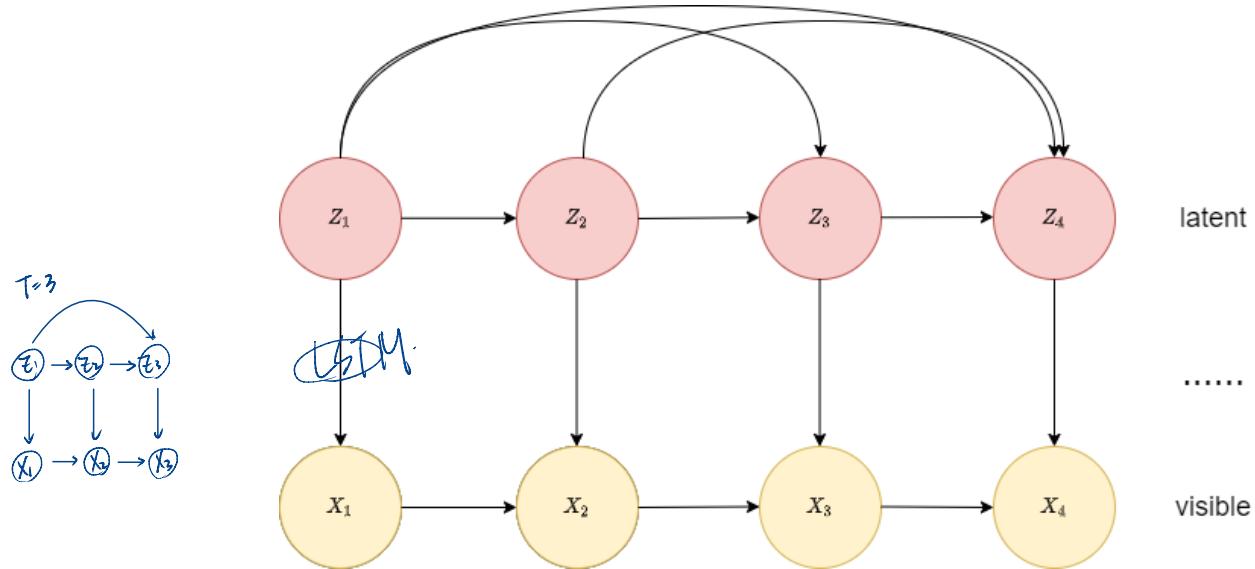
Time: 01:20pm – 04:20pm (180 minutes)

Format: Open book

Instructions:

- 1) You may give your answers in Chinese or English.
- 2) Please give your answers in succinct phrases or point form.
- 3) Please write your answers clearly (with explicit denotation of labels and symbols used).

9. (15 pts) Consider the following latent factor model, where $Z_i, i = 1, 2, \dots, T$ are latent variables and $X_i, i = 1, 2, \dots, T$ are visible variables.



- (a) (3 pts) Design an inevitable mapping from X_i 's to their latent representations Z_i 's using the flow model idea. Use $T = 3$ as an example. Think about how you should convert each of the X_i 's into the corresponding Z_i 's and what the conditioning variable should be utilized in each of these conversions. The conversions may involve coupling layers as fundamental building blocks.
- (b) (3 pts) According to the graphical model, factorize the joint distribution $p(Z_{1-T})$.
- (c) (3 pts) Describe how you would evaluate the probability $p(X_1, X_2, X_3)$ of one specific sample X_1, X_2, X_3 once the model is trained.
- (d) (3 pts) How would you train the flow model in (a)?
- (e) (3 pts) Describe how you would sample X_i 's from Z_i 's using the designed flow model.

(a)

$$(b) p(Z_{1-T}) = p(Z_1, Z_2, \dots, Z_T) = \prod_{k=1}^T p(Z_k | Z_{1:k-1}), \text{ where } k \in \{1, 2, \dots, T\}$$

$$(c) p(X_1, X_2, X_3) = p(X_1 | Z_1) p(X_2 | X_1, Z_2) p(X_3 | X_2, Z_3)$$

Q

2. (28 pts) Following the latent factor model given in Problem 1, design a variational autoencoder to learn the distribution $p(X_1, X_2, \dots, X_T)$.

- (a) (3 pts) According to the graphical model, factorize the encoding distribution $q(Z_{1-T}|X_{1-T}) = \prod_{t=1}^T q(Z_t|''\dots'')$. Apply the d-separation rule to identify the conditioning factors in ''...''.
- (b) (3 pts) By observing $q(Z_t|''\dots'')$, choose a neural network to implement it. You must specify (a) the network type, e.g. CNN or RNN, and the rationale; (b) indicate explicitly its input and output, and whether there is any hidden connection and/or output recurrence; and (c) specify the assumption about the distribution.
- (c) (3 pts) According to the graphical model, factorize the decoding distribution $p(X_{1-T}|Z_{1-T}) = \prod_{t=1}^T p(X_t|''\dots'')$. Apply the d-separation rule to identify the conditioning factors in ''...''.
- (d) (3 pts) By observing $p(X_t|''\dots'')$, choose a neural network to implement it. You must specify (a) the network type, e.g. CNN or RNN, and the rationale; (b) indicate explicitly its input and output, and whether there is any hidden connection and/or output recurrence; and (c) specify the assumption about the distribution.
- (e) (3 pts) Choose a neural network to implement the prior distribution $P(Z_1, Z_2, \dots, Z_T)$. You must specify (a) the network type, e.g. CNN or RNN, and the rationale; (b) indicate explicitly its input and output, and whether there is any hidden connection and/or output recurrence; and (c) specify the assumption about the distribution.
- (f) (10 pts) Draw a diagram to connect all these networks together for end-to-end training and specify explicitly the loss function. If re-parameterization is needed, indicate where it should be used.
- (g) (3 pts) Describe how you would sample X_i 's once the model is trained.

3. (20 pts) Convolution, dropout, activation function, and CNN

- (2 pts) Describe the reasons we don't use a full image for Transformer and how ViT addresses the issues.
- (3 pts) Explain the idea of dropout.
- (3 pts) What is the size of the output feature map for an 256 x 256 input image after convolution with kernel (3,3), padding (2,2), and stride (2,2)?
- (6 pts) What distinguishes CNN, RNN, and an attention module? What benefit does the attention module provide?
- (6 pts) What may cause gradient vanish problem and how to solve it? Explain your answer.

(a) When we using Transformer for image data,

using the full image as input will cost a lot of computation and memory.

Vision Transformer (ViT) solve this issue by dividing the image into smaller patches.

Each patches is treated as a token, similar to a word in NLP tasks.

↓
it use self-attention to capture relationships between different patches.

(b) Dropout 通过在训练过程随机丢弃一定比例的神经元，使神经网络不会过度依赖部分神经元
降低了 overfitting 的可能性。

(c) The output size = $\lfloor (\text{input size} + 2 \times \text{padding} - \text{kernel size}) / \text{stride} \rfloor + 1$

$$\begin{aligned} &= \lfloor (128 + 2 \times 2 - 3) / 2 \rfloor + 1 \\ &\quad \text{(dilation) } \times \text{kernel size} \\ &= 128 + 1 \\ &= 129 \end{aligned}$$

$$3 + 2n = 256 + 2 \times 2$$

↓ ↓ ↓ ↓
kernel stride input padding

↓

$$n = 128 \Rightarrow n = 128$$

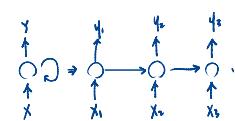
↓

$$\text{output size} = 129 \times 129$$



RNN: 用於處理序列資料 (ex. text, time-series data)

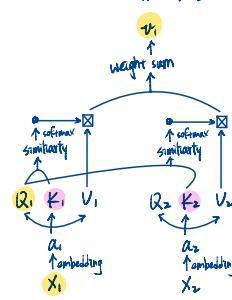
具有遞迴結構，各依循資料的順序依序依次輸入，並使用前一次的隱藏層影響這次的計算



Attention module: 利用 Q, K, V 三元組 代表注意力机制。

透過权重分配，使模型能專注於與前位移最相關的輸入部分

其目的是想取代 RNN 原本要做的事情，只是不用一個接一個地處理



→ 使用 attention module 的好處有以下

① 藉由 self-attention，使 CNN 在輸入資料維度擴大時能大幅減少所需層數

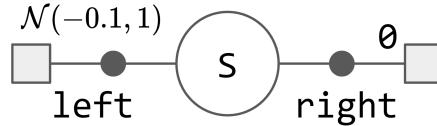
② 不用像 RNN 那麼處理資料，可平行化進行。

通过添加在某些層的跳躍連接，可以改善梯度傳播，有助於訓練更深層的神經網路

減少梯度的消失效

② 可以在反向傳遞過程中設置閾值將梯度限制在一範圍中，以防止梯度爆炸或消失的問題

- (b) The ϵ -greedy policy improvement can converge to the optimal policy by decreasing ϵ over time.
As ϵ decreases, the policy becomes more greedy, which means it selecting action with higher action values more frequently.
① When ϵ is initially high (close to 1), the policy is exploratory, as it selects random actions with equal probability. This exploration allows the agent to gather the information about the environment and learn the values.
② As the training proceeds, ϵ will be decreasing by decay mechanism. This reduction in ϵ let the policy from exploration to choosing higher action with respect to q_π is an improvement, that is $v_{\pi'}(s) \geq v_\pi(s)$.
4. (14 pts)
- (a) RL3-P.42 (a) (7 pts) Prove the following theorem: For any ϵ -greedy policy π , the ϵ -greedy policy π' values gradually with respect to q_π is an improvement, that is $v_{\pi'}(s) \geq v_\pi(s)$.
- (b) (7 pts) Show how ϵ -greedy policy improvement will converge to the optimal.
③ As ϵ close to zero, the policy becomes purely greedy, and always selecting the action with highest value.
At this point, the policy has converge to the optimal policy, as it consistently selects the actions that maximize the expected cumulative reward.
5. (18 pts) Consider the undiscounted MDP shown below.



Episodes start in S with a choice between two actions, left and right. The right action transitions immediately to the terminal state with a reward of zero. The left action also transitions to the terminal state but with a reward drawn from a normal distribution with mean -0.1 and variance 1.0 .

- (a) (2 pts) What is the optimal state-action value $Q^*(S, \text{left})$ and $Q^*(S, \text{right})$?
(b) (2 pts) What is the probability of taking the left action using the ϵ -greedy policy ($\epsilon = 0.1$) with respect to Q^* ?
(c) (14 pts) What problem will you encounter when training with Deep Q-learning (DQN)?
In order to avoid the problem, what method can we apply?

$$(a) Q^*(s, \text{left}) = E[R \sim N(-0.1, 1)] = -0.1$$

$$Q^*(s, \text{right}) = 0$$

- (b) The ϵ -greedy policy involves choosing the action with the highest state-action value with a prob. of $1-\epsilon$ (exploration), and choosing a random action with a prob. of ϵ .

$$\Rightarrow a = \begin{cases} \underset{a}{\operatorname{argmax}} Q^*(s, a) = Q^*(s, \text{right}), & \text{prob.} = 1-\epsilon = 0.9 \\ \text{random}, & \text{prob.} = \epsilon = 0.1 \end{cases}$$

Therefore, we obtain the prob. of taking the left action using the ϵ -greedy policy is 0.1 .

- (c) When training with DQN, one problem that can be encountered is the issue of overestimation of state-action value.

First, DQN uses an approximation function (typically a neural network) to estimate the state-action value, so the neural network may tends to overestimate the values for certain actions, and the overestimation will lead to suboptimal policies and slower convergence.

To address the issue, we can use DDQN, which uses two approximation functions that are trained on different samples, one for selecting the best action and the other for calculating the value of this action, since these two functions seen two different samples, it is unlikely that they overestimate the same action.