

MFDNN
Homework 11

24/05/30

(2/4)

$$p_\theta(z) = \int p(z, z) dz = \int p_\theta(z|z)p_z(z) dz = \int \frac{p_\theta(z|z)p_z(z)}{q_\phi(z|x)} q_\phi(z|x) dz$$

$$\begin{aligned} 1.(a) \quad \log p_\theta(z) &= \log(E_{Z \sim p_z} [p_\theta(z|Z)]) \\ &= \log(E_{Z \sim q_\phi(z|x)} [p_\theta(z|Z) \frac{p_z(z)}{q_\phi(z|x)}]) \\ &= \log\left(\frac{1}{K} \sum_{k=1}^K E_{z_k \sim q_\phi(z|x)} \left[p_\theta(z|z_k) \frac{p_z(z_k)}{q_\phi(z_k|x)} \right]\right) \end{aligned}$$

$f(t) = \log\left(\frac{1}{K} \sum_{i=1}^K t_i\right) = \log(t)$ is concave so, by Jensen's inequality,

$$\geq E_{z_1, \dots, z_K \sim q_\phi(z|x)} \left[\log\left(\frac{1}{K} \sum_{k=1}^K p_\theta(z|z_k) \frac{p_z(z_k)}{q_\phi(z_k|x)}\right) \right] = VLB_{\theta, \phi}^{(K)}(x) \quad \blacksquare$$

$$(b) \text{ Define } a_k = p_\theta(z|z_k) \frac{p_z(z_k)}{q_\phi(z_k|x)}$$

$$\text{Then } VLB_{\theta, \phi}^{(K)} = E_{z_1, \dots, z_K \sim q_\phi(z|x)} \left[\log\left(\frac{1}{K} \sum_{k=1}^K a_k\right) \right]$$

$$\text{and } VLB_{\theta, \phi}^{(M)} = E_{z'_1, \dots, z'_M \sim q_\phi(z|x)} \left[\log\left(\frac{1}{M} \sum_{k=1}^M a_k\right) \right]$$

w.l.o.g. (without loss of generality) say that the z'_1, \dots, z'_M sampled for $VLB^{(M)}$ are a subset of the z_1, \dots, z_K sampled for $VLB^{(K)}$ indicated by $I \subset \{1, \dots, K\}$ with $|I|=M \leq K$ as required. We can then rewrite $VLB^{(M)}$ as

$$VLB_{\theta, \phi}^{(M)} = E_{z_1, \dots, z_K \sim q_\phi} \left[E_{I \in \{z_1, \dots, z_K\}} \left[\log\left(\frac{1}{M} \sum_{k=1}^M a_{i_k}\right) \right] \right]$$

i.e. for the outer expectation we sample the same Z_k as for $VLB_{\theta, \phi}^{(K)}$ and $a_{i_k} = p_\theta(z|z_{i_k}) \frac{p_z(z_{i_k})}{q_\phi(z_{i_k}|x)}$

where $Z_{i_k} \in \{z_1, \dots, z_K\}$

$$\leq E_{z_{i_k} \sim q_\phi} \left[\log\left(E_I \left[\frac{1}{M} \sum_{k=1}^M a_{i_k} \right]\right) \right] \quad \text{by Jensen's inequality}$$

$$= E_{z_{i_k} \sim q_\phi} \left[\log\left(\frac{1}{K} \sum_{k=1}^K a_k\right) \right] \quad \text{by given hint}$$

$$= VLB_{\theta, \phi}^{(K)}(x)$$

So for $K \geq M$, $VLB_{\theta, \phi}^{(M)} \leq VLB_{\theta, \phi}^{(K)}$ \blacksquare

(c) Powerful enough refers to if the neural network, parameterized by ϕ , underlying q_ϕ can accurately represent the true posterior distribution. i.e. if $q_\phi(z|x) \approx p_\theta(z|x)$ sufficiently well
 $\exists \phi^* \text{ s.t. } \forall z_k \quad z_k \text{ for } k=1, \dots, K$

In this case,

$$\underset{\theta, \phi}{\text{maximize}} \sum_{i=1}^N \text{VLB}_{\theta, \phi}^{(k)}(X_i)$$

$$\begin{aligned} \text{close to equality by "more powerful"} &\approx \underset{\theta}{\text{maximize}} \sum_{i=1}^N E_{z_1, \dots, z_K \sim p_\theta(z|x)} \left[\log \frac{1}{K} \sum_{k=1}^K \frac{p_\theta(x|z_k) p_\theta(z_k)}{q_\phi(z_k|x)} \right] \\ &= \underset{\theta}{\text{maximize}} \sum_{i=1}^N \log \left(\frac{1}{K} \sum_{k=1}^K p_\theta(x) \right) = \underset{\theta}{\text{maximize}} \sum_{i=1}^N \log p_\theta(x) \quad \blacksquare \end{aligned}$$

$$\begin{aligned} 2.(a) \quad \log p_\theta(X_i) &= \log \left(E_{z \sim r_\lambda(z)} [p_\theta(X_i|z)] \right) \\ &= \log \left(E_{z \sim q_\phi(z|X_i)} \left[\frac{p_\theta(X_i|z) r_\lambda(z)}{q_\phi(z|X_i)} \right] \right) \\ &\geq E_{z \sim q_\phi(z|X_i)} \left[\log \left(\frac{p_\theta(X_i|z) r_\lambda(z)}{q_\phi(z|X_i)} \right) \right] \quad \text{by Jensen's inequality} \\ &= \text{VLB}_{\theta, \phi, \lambda}(X_i) \quad \blacksquare \end{aligned}$$

$$(b) \quad \underline{\nabla} \text{VLB}(X_i) = (\underline{\nabla}_\theta \text{VLB}(X_i), \underline{\nabla}_\phi \text{VLB}(X_i), \underline{\nabla}_\lambda \text{VLB}(X_i))$$

$$\begin{aligned} \underline{\nabla}_\theta \text{VLB}(X_i) &= \underline{\nabla}_\theta \int \log \left(\frac{p_\theta(X_i|z) r_\lambda(z)}{q_\phi(z|X_i)} \right) q_\phi(z|X_i) dz \\ &= \int \underline{\nabla}_\theta (p_\theta(X_i|z)) \frac{1}{p_\theta(X_i|z)} q_\phi(z|X_i) dz + \underline{0} \end{aligned}$$

$$= E_{z \sim q_\phi(z|X_i)} [\underline{\nabla}_\theta (\log(p_\theta(X_i|z)))]$$

$$\underline{\nabla}_\phi \text{VLB}(X_i) = E_{z \sim q_\phi(z|X_i)} \left[(\underline{\nabla}_\phi \log q_\phi(z)) \log \left(\frac{p_\theta(X_i|z) r_\lambda(z)}{q_\phi(z|X_i)} \right) \right] \quad \text{by log-differentiation rule for VLBs (HW10.1)}$$

$$\underline{\nabla}_\lambda \text{VLB}(X_i) = E_{z \sim q_\phi(z|X_i)} [\underline{\nabla}_\lambda (\log(r_\lambda(z)))] \quad \text{by same logic as } \underline{\nabla}_\theta$$

So stochastic gradients of $\text{VLB}_{\theta, \phi, \lambda}(X_i)$ can be computed by:

$$\begin{aligned} \underline{\nabla}_{\theta, \phi, \lambda} \text{VLB}_{\theta, \phi, \lambda}(X_i) &\approx \frac{1}{K} \sum_{k=1}^K \left(\underline{\nabla}_\theta (\log(p_\theta(X_i|z_k))), (\underline{\nabla}_\phi \log(z_k)) \log \left(\frac{p_\theta(X_i|z_k) r_\lambda(z_k)}{q_\phi(z_k|X_i)} \right), \right. \\ &\quad \left. \text{number of elements in a batch of SG calculation} \right) \quad \underline{\nabla}_\lambda (\log(r_\lambda(z_k))) \quad \text{where } z_k \sim q_\phi(z|X_i) \\ &\quad \text{z}_{ik} \text{ (sample for each } X_i \text{ once)} \end{aligned}$$

2.(c) ∇_{θ} and ∇_{λ} same as above ✓ (Expectation distribution q_{ϕ} doesn't depend on θ or λ so no need for the trick)

∇_{θ} evaluation changes if we use reparametrization vs. log-derivative trick:

$$\nabla_{\theta} VLB(X_i) = \nabla_{\theta} E_{z \sim q_{\phi}(z|X_i)} \left[\log \left(\frac{p_{\theta}(X_i|z)r_{\lambda}(z)}{q_{\phi}(z|X_i)} \right) \right] + E[\log \left(\frac{r_{\lambda}}{q_{\phi}} \right)]$$

Redo 2.(c) $VLB_{\theta, f, \lambda}(X_i) = E_{z \sim q_{\phi}(z|X_i)} \underbrace{[\log p_{\theta}(X_i|z)]}_{①} - D_{KL}(q_{\phi}(z|X_i) || r_{\lambda}(z))$ (slide 89)

- $① = E_{z \sim q_{\phi}(z|X_i)} \left[\log \left((2\pi\sigma^2)^{-d/2} \exp \left(-\frac{1}{2} (X_i - f_{\theta}(z))^T \left(\frac{1}{\sigma^2} I \right) (X_i - f_{\theta}(z)) \right) \right) \right]$

$$= -\frac{1}{2\sigma^2} E_{z \sim N(\mu_{\phi}(X_i), \Sigma_{\phi}(X_i))} [\|X_i - f_{\theta}(z)\|^2] - \frac{d}{2} \log(2\pi\sigma^2)$$

By reparametrization trick:

$$= -\frac{1}{2\sigma^2} E_{\varepsilon \sim N(0, I)} \left[\|X_i - f_{\theta}(\mu_{\phi}(X_i) + \sum_{\phi}^{ln} (X_i) \varepsilon)\|^2 \right] - \frac{d}{2} \log(2\pi\sigma^2)$$

→ Apply ∇_{ϕ} : $-\frac{1}{2\sigma^2} E_{\varepsilon \sim N(0, I)} [\nabla_{\phi} \| \dots \|^2]$ which can be evaluated + back-propagated

- $② = -\frac{1}{2} \left(\text{tr} \left([\text{diag}(\lambda_2)]^{-1} \sum_{\phi}(X_i) \right) + (\lambda_1 - \mu_{\phi}(X_i))^T [\text{diag}(\lambda_2)]^{-1} (\lambda_1 - \mu_{\phi}(X_i)) - k + \log \left(\frac{\det(\text{diag}(\lambda_2))}{\det(\sum_{\phi}(X_i))} \right) \right)$

by HW9.5

→ Gradient ∇_{θ} can then be calculated + back propagated

Redo 4.(a) $E_{P_A, P_B} [\text{points for } B] = P_A^T M P_B$

$$M := \begin{bmatrix} 0 & -1 & 1 \\ 1 & 0 & -1 \\ -1 & 1 & 0 \end{bmatrix} \quad P_C = \begin{bmatrix} P_{C,\text{rock}} \\ P_{C,\text{paper}} \\ P_{C,\text{scissors}} \end{bmatrix}$$

$$= P_{A1}(P_{B3} - P_{B2}) + P_{A2}(P_{B1} - P_{B3}) + P_{A3}(P_{B2} - P_{B1}) \quad (\text{as I did originally})$$

- Clearly, $P_A^T M P_B \geq \min_{P_B \in \Delta^3} P_A^T M P_B \Rightarrow \min_{P_B \in \Delta^3} P_A^T M P_B \geq \max_{P_A \in \Delta^3} P_A^T M P_B$

$$\Rightarrow \min_{P_A} \max_{P_B} P_A^T M P_B \geq \min_{P_B} \max_{P_A} P_A^T M P_B \text{ by above}$$

$$= \max_{P_B} \min_{P_A} P_A^T M P_B \text{ by above}$$

$$= \max_{P_B} \min_{P_A} P_A^T M P_B \quad (**)$$

$$\min_{P_A} \max_{P_B} P_A^T M P_B \leq \max_{P_B} P_A^T M P_B = 0 \quad (**)$$

$$0 \geq \min_{P_A} \max_{P_B} P_A^T M P_B$$

$$= \min_{P_A} \max_{P_B} -P_B^T M P_A$$

$$= -\max_{P_A} \min_{P_B} P_B^T M P_A$$

$$= -\max_{P_B} \min_{P_A} P_A^T M P_B$$

$$\geq -\min_{P_A} \min_{P_B} P_A^T M P_B \geq 0 \quad \text{by (*) and then (**)}$$

by (**)

$$\therefore P_A^T M P_B = (P_A^T M P_B)^T = P_B^T M^T P_A = -P_B^T M P_A$$

$$\Rightarrow \min_{P_A} \max_{P_B} P_A^T M P_B = \max_{P_B} \min_{P_A} P_A^T M P_B = 0 \Leftrightarrow \text{existence of saddle pt solution } P_A^*, P_B^*$$

- Show unique by proving $P_A^* M P_B \leq P_A^{*T} M P_B^* = 0 \leq P_A^T M P_B^* \quad \forall P_A, P_B \in \Delta^3$ only for given P_A^*, P_B^*

Suppose $P_B^* \neq [1/3, 1/3, 1/3]^T \Leftrightarrow M P_B^* = [P_{B_1} - P_{B_2}, P_{B_2} - P_{B_3}, P_{B_3} - P_{B_1}]^T \neq [0, 0, 0]^T$

\Rightarrow At least one entry of $M P_B^*$ is negative

So pick P_A such that $P_A^T M P_B^* < 0$ to obtain a contradiction with $P_A^T M P_B \geq 0 \quad \forall P_A \in \Delta^3$ above.

$$\Rightarrow P_B^* = [1/3, 1/3, 1/3]^T$$

Analogous argument with $P_A^* M P_B \leq 0$ to show $P_A^* = [1/3, 1/3, 1/3]^T$

So given P_A^*, P_B^* is indeed the unique solution.

(b) (I'd misunderstood the Q the first time - I assumed we had set that B plays p_B^*)

If A chooses $P_A \neq P_A^*$, then B can choose one p_B of $[1, 0, 0]^T, [0, 1, 0]^T, [0, 0, 1]^T$ to make

$E_{P_A, P_B} [\text{points for } B] > 0$! So no, not just any strategy is optimal for A. Only in the specific case where B plays p_B^* is A free to pick any strategy since $E_{P_A, P_B^*} [\text{points for } B] = 0$.

4.(a) In this model, each game is independent so $E[\text{points for } B] = \text{num games} \times E[\text{points won by } B \text{ in 1 game}]$

$$E[\text{points won by } B \text{ in 1 game}] = (P(B \text{ rock}) P(A \text{ scissors}) + P(B \text{ paper}) P(A \text{ rock}) + \\ P(B \text{ scissors}) P(A \text{ paper})) - (P(B \text{ rock}) P(A \text{ paper}) + \\ P(B \text{ paper}) P(A \text{ scissors}) + P(B \text{ scissors}) P(A \text{ rock})) + 0$$

Notation:

$$(p_C)_i =: p_{Ci} \\ = p_{B1} p_{A3} + p_{B2} p_{A1} + p_{B3} p_{A2} - p_{B1} p_{A2} - p_{B2} p_{A3} - p_{B3} p_{A1} \\ = p_{B1} (p_{A3} - p_{A2}) + p_{B2} (p_{A1} - p_{A3}) + p_{B3} (p_{A2} - p_{A1})$$

$$\Rightarrow E_{p_A^*, p_B^*} [\text{points for } B] = p_{B1} \left(\frac{1}{3} - \frac{1}{3}\right) + p_{B2} \left(\frac{1}{3} - \frac{1}{3}\right) + p_{B3} \left(\frac{1}{3} - \frac{1}{3}\right) = 0 \quad \text{for all } p_B \in \Delta^3$$

$$E_{p_A^*, p_B^*} [\text{points for } B] = \frac{1}{3} (p_{B2} - p_{B1}) + \frac{1}{3} (p_{B1} - p_{B3}) + \frac{1}{3} (p_{B3} - p_{B1}) = 0 \quad \text{for all } p_B \in \Delta^3$$

$$E_{p_A^*, p_B^*} [\text{points for } B] = 0 \Rightarrow E_{p_A^*, p_B^*} [\text{points for } B] \leq E_{p_A^*, p_B} [\text{points for } B] \leq E_{p_A, p_B^*} [\text{points for } B] \\ \text{for all } p_A, p_B \in \Delta^3$$

\Rightarrow so p_A^*, p_B^* is a solution to minimax problem

Suppose there exists another solution (p_A', p_B') , without lost one necessarily different from p_A^*, p_B^* .

• If only one is different (wlog assume $p_A' = p_A^*$ and $p_B' \neq p_B^*$),

$E_{p_A^*, p_B'} [\text{pts for } B] = E_{p_A^*, p_B^*} [\text{pts for } B] = 0$ still and then $E_{p_A^*, p_B'} [\text{pts for } B] \geq 0$ regard:

$$E_{p_A^*, p_B'} [\text{pts for } B] = p_B' \cdot (p_A^* \times 1)$$

... ?

• If both are different ($p_A' \neq p_A^*$ and $p_B' \neq p_B^*$)

... ?

4.(b) Yes, since this is a symmetric game in the sense that a loss for one player is a win for the other. Therefore $E_{p_A, p_B} [\text{pts for } B] = 0 \Rightarrow E_{p_A, p_B} [\text{pts for } A] = 0$. So if B plays with $p_B = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ then no particular strategy can be better for A in terms of winning more points in expectation; i.e. any $p_A \in \Delta^3$ is optimal for A.