



یادگیری عمیق

نیم سال دوم ۰۳-۰۲
مدرس: مهدیه سلیمانی

تمرین چهارم

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- در طول ترم امکان ارسال با تاخیر تمرین های نظری بدون کسر نمره تا سقف ۵ روز و تمرین های عملی تا سقف ۱۰ روز وجود دارد. محل بارگزاری جواب تمرین های نظری بعد از ۳ روز و تمرین های عملی بعد از ۵ روز بسته خواهد شد و پس از گذشت این مدت، پاسخ های ارسال شده پذیرفته نخواهند شد.
- هم فکری در انجام تمرین مانعی ندارد، فقط توجه داشته باشید که پاسخ تمرین حتما باید توسط خود شخص نوشته شده باشد. همچنین در صورت هم فکری در هر تمرین، در ابتدای جواب تمرین نام افرادی که با آن ها هم فکری کرده اید را حتما ذکر کنید.
- برای پاسخ به سوالات نظری در صورتی که از برگه خود عکس تهیه می کنید، حتما توجه داشته باشید که تصویر کاملا واضح و خوانا باشد. در صورتی که خوانایی کافی را نداشته باشد، تصحیح نخواهد شد.
- محل بارگذاری سوالات نظری و عملی در هر تمرین مجزا خواهد بود. به منظور بارگذاری بایستی تمارین تئوری در یک فایل pdf با نام `HW4_[First-Name]_[Last-Name]_[Student-Id].pdf` و تمارین عملی نیز در یک فایل مجزای زیپ با نام `HW4_[First-Name]_[Last-Name]_[Student-Id].zip` بارگذاری شوند.
- در صورت وجود هرگونه ابهام یا مشکل، در کوثرای درس آن مشکل را بیان کنید و از پیغام دادن مستقیم به دستیاران آموزشی خودداری کنید.
- طراحان این تمرین: آقایان جواهریان، حسینی، علیخانی، ثقفیان

بخش نظری (۶۵ نمره (+۱۵ امتیازی))

سوال اول: استنباط متغیر (۲۰ نمره)

در این تمرین پایه های تئوری مدل VAE و پیشرفت های حاصل شده بر روی آن مورد بررسی قرار می گیرند. یک مدل احتمالاتی که به صورت مجموعه ای از متغیر های قابل مشاهده X و مجموعه متغیر های پنهان Z می باشد را در نظر بگیرید؛ در استنباط بیزی ما علاقه مند به محاسبه توزیع پسین Z بعد از مشاهده دادگان هستیم. چالش اصلی این روش، intractable بودن توزیع پسین می باشد. استنباط متغیر (Variational Inference) یکی از روش های حل این چالش می باشد. در این روش، محاسبه توزیع پسین را به یک مسئله بهینه سازی روی خانواده ای از توزیع ها تبدیل می کنیم:

$$q = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} D_{\text{KL}}(q(z) \parallel p_{\theta}(z|x))$$

۱. با بسط دادن رابطه مطرح شده، به رابطه ELBO برسید و نشان دهید، ELBO کران پایینی برای $\log p_{\theta}(x)$ می باشد.

۲. فرض کنید مجموعه دادگان $D = \{x_i\}$ را در اختیار داریم؛ اگر بخواهیم تخمین MLE برای پارامترهای مدل ارائه دهیم؛ بایستی تابع زیر را بیشینه کنیم:

$$\log p(D|\theta) = \sum_{i=1}^N \log p(x = x_i|\theta)$$

اما محاسبه مستقیم $\log p(x = x_i|\theta)$ برایمان مقدور نیست. بنابراین به جای رابطه بالا، رابطه زیر را بیشینه می کنیم:

$$\sum_{i=1}^N \max_{q \in \mathbb{Q}} \text{ELBO}(q, x_i, \theta)$$

حال فرض کنید خانواده Q به کمک پارامتر ψ قابل پارامتریز کردن می باشد. به کمک الگوریتم کاهش گرادیان، یک الگوریتم برای انجام MLE ارائه دهید.

۳. رابطه قبلی از لحاظ محاسباتی سنگین می باشد. برای حل این چالش دو روش ارائه می شود:

Stochastic VI (\bar{I}): در این روش به کمک تخمین مونته کارلو، جمع موجود در تابع هزینه، با یک جمع روی batch mini تخمین زده می شود:

$$\sum_{i=1}^N \max_{q \in \mathbb{Q}} \text{ELBO}(q, x_i, \theta) \simeq \frac{N}{B} \sum_{i=1}^B \max_{q \in \mathbb{Q}} \text{ELBO}(q, x_i, \theta)$$

Amortized VI (ب): در این روش به جای پیدا کردن N پارامتر متغیر $(\psi_{1:N})$ ، یک شبکه عصبی با پارامتر ϕ بهینه می شود که با ورودی گرفتن x_i ، پارامتر بهینه ψ_i را خروجی دهد:

$$q(z_n|\psi_n) = q(z_n|g_\phi(x_n)) = q_\phi(z_n|x_n)$$

$$\sum_{i=1}^B \max_{q \in \mathbb{Q}} \text{ELBO}(q, x_i, \theta) = \max_{\phi} \sum_{i=1}^B \text{ELBO}(q_\phi(\cdot|x_i), x_i, \theta)$$

این دو روش را در الگوریتم طراحی شده برای قسمت قبل اعمال کنید و الگوریتم تغییر یافته را ارائه دهید.

۴. حال با در نظر گرفتن مفروضات زیر و به کمک reparameterization trick الگوریتم خود را برای این حالت خاص، مشخص کنید:

$$p_\theta(z) = \mathcal{N}(0, I) \quad (۱)$$

$$p_\theta(x|z) = \mathcal{N}(f_\theta(z), \sigma^2 I) \quad (۲)$$

$$q_\phi(\cdot|x) = \mathcal{N}(\mu_\phi(x), \text{diag}(\sigma_\phi^2(x))) \quad (۳)$$

* توجه کنید که بایستی تابع هزینه خود را به ساده ترین شکل ممکن در بیاورید. (بایستی KL دو توزیع نرمال را ساده کنید)

۵. در چند سوال قبل به بررسی حالت خاصی از VI به نام fixed-form VI پرداختیم. می توان به جای انتخاب یک خانواده پارامتری از توزیع های متغیر و انجام مسئله بهینه سازی به کمک کاهش گرادیان، انتخاب های دیگری نیز انجام داد. در این سوال به بررسی یک روش پر استفاده دیگر به نام Mean-Field VI می پردازیم. فرض کنید متغیر نهان Z یک بردار J بعدی باشد و توزیع متغیر را به صورت زیر بنویسیم:

$$q_{\psi}(z) = \prod_{j=1}^J q_j(z_j)$$

حال فرض کنید تمامی q_i ها بغیر از q_j ثابت باشند. نشان دهید q_j که مقدار ELBO را بیشینه می کند به صورت زیر می باشد:

$$q_j(z_j) \propto \exp [\mathbb{E}_{z_{-j}} [\log p_{\theta}(x, z_j, z_{-j})]]$$

که در آن z_{-j} مجموعه تمامی z_i ها بغیر از z_j می باشد.

سوال دوم: Diffusion Models (۲۵ نمره)

در کلاس درس با خانواده ای از مدل های دیفیوژنی به نام DDPMs (Denoising Diffusion Probabilistic Models) آشنا شدیم و در این سوال قصد داریم با خانواده دیگری از این مدل های دیفیوژنی به نام VDMs (Variational Diffusion Models) آشنا بشویم. برعکس DDPMs که فضای نهان را به صورت گسسته در زمان در نظر می گرفت، VDMs فضای نهان به صورت پیوسته در نظر گرفته می شود. در ادامه با برخی از مراحل این مدل VDMs بیشتر آشنا می شویم.

فرایند رو به جلو (Forward Process)

در این فرایند، قصد داریم به آرامی داده x را به یک نویز تصادفی تبدیل کنیم و این کار را نیز با استفاده از نویز گاوسی در هر مرحله انجام می دهیم. در اینصورت، نسخه نویزی شده x در گام t را متغیر نهان z_t می نامیم که مقادیر t در بازه $t \in [0, 1]$ قرار دارند. به این معنا که برای $t = 0$ ، کمترین تصویر نویزی و برای $t = 1$ ، بیشترین تصویر نویزی را خواهیم داشت. همچنین، فرض می کنیم که توزیع شرطی z_t برحسب x به صورت زیر می باشد:

$$q(z_t|x) = \mathcal{N}(a_t x, \sigma_t^2 I)$$

که در عبارت بالا a_t و σ_t^2 هر کدام توابعی مثبت برحسب t می باشند.

(آ) یکی از شرایطی که ما دوست داریم در فرایند فروارد رعایت شود این مورد می باشد که واریانس متغیرهای نهان ما در طول فرایند عوض نشود. نشان دهید برای اینکه چنین شرطی رعایت شود باید رابطه زیر برقرار باشد. (فرض شود داده ورودی استاندارد شده می باشد). (۲ نمره)

$$a_t = \sqrt{1 - \sigma_t^2}$$

Answer:

We have this relation:

$$z_t = a_t x + \sigma_t \epsilon_t; \quad \epsilon_t \sim \mathcal{N}(0, I)$$

So we could say,

$$V[z_t] = V[a_t x + \sigma_t \epsilon_t] = V[a_t x] + V[\sigma_t \epsilon_t] = a_t^2 V[x] + \sigma_t^2 V[\epsilon_t] = a_t^2 V[x] + \sigma_t^2$$

We want that $V[x] = V[z_t]$

$$V[x] = a_t^2 V[x] + \sigma_t^2 \rightarrow a_t^2 = \frac{V[x] - \sigma_t^2}{V[x]} \rightarrow a_t^2 = 1 - \frac{\sigma_t^2}{V[x]} \rightarrow a_t = \sqrt{1 - \frac{\sigma_t^2}{V[x]}}$$

Because our input data is standardized, we can rewrite the expression as below:

$$a_t = \sqrt{1 - \sigma_t^2}$$

مشابه DDPM فرض کنید زنجیره تغییرات ما از x تا z_1 یک زنجیره مارکوفی (Markov chain) باشد. زنجیره مارکوفی به این معنا می باشد که برای بدست آوردن تصویر نویزی شده در زمان t ما فقط نیاز به آخرین تصویر کمتر نویزی بدست آمده داریم. در حالت گسسته در زمان این زنجیره مارکوف به صورت زیر می باشد.

$$z_1 \leftarrow z_{\left(\frac{T-1}{T}\right)} \leftarrow \dots \leftarrow z_0 \leftarrow x$$

در حالت پیوسته $T \rightarrow \infty$ میل می کند، در نتیجه مقدار تغییر در هر گام بسیار ریز می باشد.

(ب) با استفاده از تعریف مارکوف بودن فرایند فروارد، نشان دهید که رابطه زیر برقرار می باشد. (۲ نمره)

$$q(z_s, z_t | x) = q(z_t | z_s) q(z_s | x)$$

Answer:

First, let's simplify it a little bit.

$$q(z_s, z_t) = q(z_s) q(z_t | z_s)$$

Then we can conditional the above relation:

$$q(z_s, z_t | x) = q(z_s | x) q(z_t | z_s, x)$$

Finally, we can use Bayes' theorem for our main relation:

$$q(z_s, z_t | x) = \frac{q(x) q(z_s, z_t | x)}{q(x)} = q(z_s | x) q(z_t | z_s, x) = q(z_s | x) q(z_t | z_s)$$

Since the forward process is a Markov chain $q(z_t | z_s, x) = q(z_t | z_s)$

(ج) نشان دهید توزیع $q(z_t | z_s)$ به ازای $t > s$ یک توزیع گاوسی به صورت زیر می باشد. (۴ نمره)

$$q(z_t | z_s) = \mathcal{N}(a_{t|s} z_s, \sigma_{t|s}^2 I)$$

$$a_{t|s} = \frac{a_t}{a_s}, \sigma_{t|s}^2 = \sigma_t^2 - a_{t|s}^2 \sigma_s^2$$

مجدداً مشابه DDPM ما علاقه مند هستیم تا فرم بسته ای برای توزیع $q(z_s | z_t, x); t > s$ بدست بیاوریم زیرا از آن در فرایند رو به عقب (Reverse process) و ساده سازی تابع هزینه مدل استفاده می شود.

Answer:

Let's focus on deriving $a_{t|s}$ first. By construction, we know that each z_t is given by:

$$z_t = a_t x + \sigma_t \epsilon_t = a_t \left(\frac{z_s - \sigma_s \epsilon_s}{a_s} \right) + \sigma_t \epsilon_t$$

Since $x = \frac{(z_s - \sigma_s \epsilon_s)}{a_s}$ for any $s < t$. The conditional mean of $q(z_t | z_s)$ is then given by:

$$E[z_t | z_s] = a_t \left(\frac{z_s - \sigma_s E[\epsilon_s]}{a_s} \right) + \sigma_t E[\epsilon_t] = \frac{a_t}{a_s} z_s = a_{t|s} z_s$$

To compute a closed-form expression for the variance of $q(z_t | z_s)$ we can start by rewriting the equation for z_t in terms of preceding latent z_s as follows:

$$z_t = a_{t|s} z_s + \sigma_{t|s} \epsilon_t = \frac{a_t}{a_s} (a_s x + \sigma_s \epsilon_s) + \sigma_{t|s} \epsilon_t = a_t x + \frac{a_t}{a_s} \sigma_s \epsilon_s + \sigma_{t|s} \epsilon_t \rightarrow z - a_t x = \frac{a_t}{a_s} \sigma_s \epsilon_s + \sigma_{t|s} \epsilon_t$$

$$\sigma_t \epsilon_t = \frac{a_t}{a_s} \sigma_s \epsilon_s + \sigma_{t|s} \epsilon_t$$

The above implication allows us to compute the variance $\sigma_{t|s}^2$ straightforwardly.

$$V[\sigma_t \epsilon_t] = V \left[\frac{a_t}{a_s} \sigma_s \epsilon_s + \sigma_{t|s} \epsilon_t \right]$$

$$\sigma_t^2 V[\epsilon_t] = \left(\frac{a_t}{a_s} \right)^2 \sigma_s^2 V[\epsilon_s] + \sigma_{t|s}^2 V[\epsilon_t]$$

$$\sigma_t^2 = \left(\frac{a_t}{a_s} \right)^2 \sigma_s^2 + \sigma_{t|s}^2 \rightarrow \sigma_{t|s}^2 = \sigma_t^2 - a_{t|s}^2 \sigma_s^2$$

(د) حالا با استفاده از نتایج بدست آمده از قسمت (ب) و قسمت (ج) نشان دهید روابط زیر برقرار می‌باشند. (۸ نمره)

$$q(z_s|z_t, x) = \mathcal{N}(\mu_Q(z_t, x; s, t), \sigma_Q^2(s, t)I)$$

$$\mu_Q(z_t, x; s, t) = \frac{(a_{t|s}\sigma_s^2)}{\sigma_t^2} z_t + \frac{(a_s\sigma_{t|s}^2)}{\sigma_t^2} x$$

$$\sigma_Q^2(s, t) = \frac{(\sigma_{t|s}^2\sigma_s^2)}{\sigma_t^2}$$

راهنمایی: ابتدا نشان دهید $q(z_s|z_t, x) \propto q(z_t|z_s)q(z_s|x)$ برقرار می‌باشد و سپس عبارت را باز کنید و همچنین زمانیکه دو گاوسی را در هم ضرب می‌کنید ضرایب را در نظر بگیرید.

Answer:

$$\begin{aligned} q(z_s|z_t, x) &= \frac{q(z_s, z_t|x)q(x)}{q(z_t|x)q(x)} = \frac{q(z_s, z_t|x)}{q(z_t|x)} = \frac{q(z_t|z_s)q(z_s|x)}{q(z_t|x)} \propto q(z_t|z_s)q(z_s|x) \\ &= \prod_{i=1}^D \frac{1}{\sigma_{t|s}\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma_{t|s}^2}(z_{t,i} - a_{t|s}z_{s,i})^2\right\} \cdot \prod_{i=1}^D \frac{1}{\sigma_s\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma_s^2}(z_{s,i} - a_s x_i)^2\right\} \\ &\propto \prod_{i=1}^D \exp\left\{-\frac{1}{2\sigma_{t|s}^2}(z_{t,i} - a_{t|s}z_{s,i})^2\right\} \prod_{i=1}^D \exp\left\{-\frac{1}{2\sigma_s^2}(z_{s,i} - a_s x_i)^2\right\} \\ &= \prod_{i=1}^D \exp\left\{-\frac{1}{2\sigma_{t|s}^2}(z_{t,i} - a_{t|s}z_{s,i})^2 - \frac{1}{2\sigma_s^2}(z_{s,i} - a_s x_i)^2\right\} \end{aligned}$$

With simplifying the above expression, we can reach this below one:

$$\begin{aligned} &= \prod_{i=1}^D \exp\left\{-\frac{1}{2\sigma_{t|s}^2}(z_{t,i} - a_{t|s}z_{s,i})^2 - \frac{1}{2\sigma_s^2}(z_{s,i} - a_s x_i)^2\right\} \\ &= \prod_{i=1}^D \exp\left\{-\frac{1}{2}\left[z_{s,i}^2\left(\frac{a_{t|s}^2}{\sigma_{t|s}^2} + \frac{1}{\sigma_s^2}\right) - 2z_{s,i}\left(\frac{a_{t|s}z_{t,i}}{\sigma_{t|s}^2} + \frac{a_s x_i}{\sigma_s^2}\right) + \frac{z_{t,i}^2}{\sigma_{t|s}^2} + \frac{a_s^2 x_i^2}{\sigma_s^2}\right]\right\} \end{aligned}$$

The next step is to match the above equation with what we expect to see in a Gaussian distribution.

$$N(x; \mu, \sigma^2) \propto \exp\left\{-\frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} - \frac{\mu^2}{2\sigma^2}\right\}$$

Without loss of generality, consider the D=1 dimensional case for brevity.

From the first term, we get:

$$\frac{1}{\sigma_Q^2} = \frac{a_{t|s}^2}{\sigma_{t|s}^2} + \frac{1}{\sigma_s^2} \rightarrow \frac{\sigma_s^2 a_{t|s}^2}{\sigma_s^2 \sigma_{t|s}^2} + \frac{\sigma_{t|s}^2}{\sigma_{t|s}^2 \sigma_s^2} \rightarrow \sigma_Q^2 = \frac{\sigma_{t|s}^2 \sigma_s^2}{a_{t|s}^2 \sigma_s^2 + \sigma_{t|s}^2}$$

From the second term, we get:

$$\frac{a_{t|s}z_t}{\sigma_{t|s}^2} + \frac{a_s x}{\sigma_s^2} = \frac{\mu_Q}{\sigma_Q^2} \rightarrow \mu_Q = \sigma_Q^2 \left(\frac{a_{t|s}z_t}{\sigma_{t|s}^2} + \frac{a_s x}{\sigma_s^2} \right)$$

With simplify the above equation, we can have:

$$\mu_Q = \frac{a_{t|s}\sigma_s^2}{a_{t|s}^2\sigma_s^2 + \sigma_{t|s}^2} z_t + \frac{a_s\sigma_{t|s}^2}{a_{t|s}^2\sigma_s^2 + \sigma_{t|s}^2} x$$

Using the fact $\sigma_{t|s}^2 = \sigma_t^2 - a_{t|s}^2\sigma_s^2$ from question (b) we can get:

$$\begin{aligned} \mu_Q(z_t, x; s, t) &= \frac{a_{t|s}\sigma_s^2}{\sigma_t^2} z_t + \frac{a_s\sigma_{t|s}^2}{\sigma_t^2} x \\ \sigma_Q^2(s, t) &= \frac{\sigma_{t|s}^2\sigma_s^2}{\sigma_t^2} \end{aligned}$$

فرایند رو به عقب و تابع هزینه مدل (Loss function & Reverse process)

فرض کنید ما دو مدل با نام‌های $\hat{x}_\theta(z_t; t)$ و $p_\theta(z_s|z_t)$ آموزش داده‌ایم که در ادامه عملکرد هر کدام را توضیح می‌دهیم. مدل $\hat{x}_\theta(z_t; t)$: ورودی آن تصویر نویزی و زمان مربوطه می‌باشد و خروجی آن تصویر بدون نویز x می‌باشد. مدل $p_\theta(z_s|z_t)$: شبکه‌ای هست که سعی می‌کند تا رفتار $q(z_s|z_t, x)$ را تقلید بکند به این معنی که اگر تصویر نویزی شده را به شبکه بدهیم، مدل به ما تصویر کمتر نویزی شده در لحظه قبلی را می‌دهد و برای انجام این کار نیز x را از مدل $\hat{x}_\theta(z_t; t)$ بدست می‌آورد و به صورت زیر تعریف می‌شود.

$$p_\theta(z_s|z_t) = \mathcal{N}(\mu_\theta(z_t; s, t), \sigma_Q^2(s, t)I)$$

$$\mu_\theta(z_t; s, t) = \frac{(a_{t|s}\sigma_s^2)}{\sigma_t^2}z_t + \frac{(a_s\sigma_{t|s}^2)}{\sigma_t^2}\hat{x}_\theta(z_t; t)$$

(ه) نشان دهید که KL divergence بین دو توزیع $p_\theta(z_s|z_t)$ و $q(z_s|z_t, x)$ به صورت زیر می‌باشد. (۲ نمره)

$$D_{KL}(q(z_s|z_t, x) \parallel p_\theta(z_s|z_t)) = \frac{1}{2\sigma_Q^2(s, t)} \|\mu_Q(z_t, x; s, t) - \mu_\theta(z_t; s, t)\|_2^2$$

راهنمایی: از رابطه زیر در ساده‌سازی استفاده کنید.

$$D_{KL}(N_0 \parallel N_1) = \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1}\Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1}(\mu_1 - \mu_0) - d + \log\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right) \right)$$

Answer:

$$\begin{aligned} D_{KL}(q(z_s|z_t, x) \parallel p_\theta(z_s|z_t)) &= \frac{1}{2} \left[\text{tr} \left(\frac{1}{\sigma_Q^2} I \sigma_Q^2 I \right) - D + (\mu_\theta - \mu_Q)^T \frac{1}{\sigma_Q^2} I (\mu_\theta - \mu_Q) + \log \frac{\det \sigma_Q^2 I}{\det \sigma_Q^2 I} \right] \\ &= \frac{1}{2} \left[D - D + \frac{1}{\sigma_Q^2} (\mu_\theta - \mu_Q)^T (\mu_\theta - \mu_Q) + 0 \right] = \frac{1}{2\sigma_Q^2} \sum_{i=1}^D (\mu_{Q,i} - \mu_{\theta,i})^2 = \frac{1}{2\sigma_Q^2(s, t)} \\ &\quad \parallel \mu_Q(z_t, x; s, t) - \mu_\theta(z_t; s, t) \parallel_2^2 \end{aligned}$$

همچنین SNR (نسبت سیگنال به نویز) نیز به صورت زیر تعریف می‌شود:

$$\text{SNR}(t) = \frac{a_t^2}{\sigma_t^2}$$

(و) نشان دهید که می‌توان رابطه بدست آمده در قسمت (ه) را نیز بیشتر ساده‌سازی کرد و به عبارت زیر رسید. (۳ نمره)

$$D_{KL}(q(z_s|z_t, x) \parallel p_\theta(z_s|z_t)) = \frac{1}{2} (\text{SNR}(s) - \text{SNR}(t)) \|x - \hat{x}_\theta(z_t; t)\|_2^2$$

راهنمایی: ابتدا مشابه کاری که در DDPM برای ساده‌سازی انجام میدادید را بر روی نتیجه قسمت (ه) اعمال کنید تا به عبارتی به فرم زیر برسید و سپس ضریب بدست آمده را براساس SNR بازنویسی بکنید تا به حکم مسئله برسید.

$$D_{KL}(q(z_s|z_t, x) \parallel p_\theta(z_s|z_t)) = \frac{1}{2} \gamma \|x - \hat{x}_\theta(z_t; t)\|_2^2$$

Answer:

$$\begin{aligned}
D_{KL}(q(z_s|z_t, x) \parallel p_\theta(z_s|z_t)) &= \frac{1}{2\sigma_Q^2(s, t)} \|\mu_Q(z_t, x; s, t) - \mu_\theta(z_t; s, t)\|_2^2 \\
&= \frac{1}{2\sigma_Q^2(s, t)} \left\| \frac{a_{t|s}^2 \sigma_s^2}{\sigma_t^2} z_t + \frac{a_s \sigma_{t|s}^2}{\sigma_t^2} x - \left(\frac{a_{t|s} \sigma_s^2}{\sigma_t^2} z_t + \frac{a_s \sigma_{t|s}^2}{\sigma_t^2} \widehat{x}_\theta(z_t; t) \right) \right\|_2^2 \\
&= \frac{1}{2\sigma_Q^2(s, t)} \left\| \frac{a_s \sigma_{t|s}^2}{\sigma_t^2} x - \frac{a_s \sigma_{t|s}^2}{\sigma_t^2} \widehat{x}_\theta(z_t; t) \right\|_2^2 \\
&= \frac{1}{2\sigma_Q^2(s, t)} \left(\frac{a_s \sigma_{t|s}^2}{\sigma_t^2} \right)^2 \|x - \widehat{x}_\theta(z_t; t)\|_2^2
\end{aligned}$$

Recall $\sigma_Q^2(s, t) = \frac{\sigma_{t|s}^2 \sigma_s^2}{\sigma_t^2}$

$$= \frac{\sigma_t^2}{2\sigma_{t|s}^2 \sigma_s^2} \left(\frac{a_s \sigma_{t|s}^2}{\sigma_t^2} \right)^2 \|x - \widehat{x}_\theta(z_t; t)\|_2^2$$

With some steps of simplification, we can have:

$$\begin{aligned}
&= \frac{1}{2} \left(\frac{a_s^2}{\sigma_s^2} - \frac{a_t^2}{\sigma_t^2} \right) \|x - \widehat{x}_\theta(z_t; t)\|_2^2 \\
&= \frac{1}{2} (\text{SNR}(s) - \text{SNR}(t)) \|x - \widehat{x}_\theta(z_t; t)\|_2^2
\end{aligned}$$

حالا ما برای اینکه تابع هزینه را بنویسیم نیاز داریم تا مقداری مسئله را سبک کنیم، فرض کنید ما بازه‌ی $[0, 1]$ را به T قسمت با اندازه برابر $l = 1/T$ تقسیم کرده باشیم. در اینصورت مسئله ما به همان حالت گسسته در زمان DDPM تبدیل می‌شود در اینصورت می‌توان تابع هزینه VDM را نوشت:

$$L_T(x) = \sum_{i=1}^T \mathbb{E}_{q(z_{t(i)}|x)} [D_{KL}(q(z_{s(i)}|z_{t(i)}, x) \parallel p_\theta(z_{s(i)}|z_{t(i)}))]$$

در عبارت بالا $t(i) = \frac{i}{T}$ و $s(i) = \frac{i-1}{T}$ می‌باشند. می‌توان با نتایج بدست آمده از قسمت‌های قبل و اندکی ساده‌سازی تابع هزینه بالا را به صورت زیر بازنویسی کرد.

$$L_T(x) = \frac{T}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I), i \sim \mathcal{U}\{1, T\}} [(\text{SNR}(s(i)) - \text{SNR}(t(i))) \|x - \widehat{x}_\theta(z_{t(i)}; t(i))\|^2]$$

اما ما دوست داریم که فضای نهان پیوسته داشته باشیم، یعنی می‌خواهیم $T \rightarrow \infty$ میل بکند تا بتوان $L_\infty(x)$ را بدست بیاوریم که براساس آن مدل را آموزش بدهیم.
(ر) نشان دهید $L_\infty(x)$ به صورت زیر می‌باشد. (۴ نمره)

$$L_\infty(x) = \lim_{T \rightarrow \infty} L_T(x) = -\frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I), t \sim \mathcal{U}(0, 1)} [\text{SNR}'(t) \|x - \widehat{x}_\theta(z_t; t)\|^2]$$

راهنمایی: برای اینکه بتوان به عبارت بالا رسید باید ابتدا s را براساس t بازنویسی بکنید و سپس از تعریف مشتق استفاده بکنید.

Answer:

For simplicity, I have removed the indices, but they exist.

$$\begin{aligned}
 L_T(x) &= \frac{1}{2} E_{\epsilon \sim N(0, I), i \sim U\{1, T\}} \left[T \left(SNR \left(t - \frac{1}{T} \right) - SNR(t) \right) \| x - \widehat{x}_\theta(z_t; t) \|^2 \right] \\
 &= \frac{1}{2} E_{\epsilon \sim N(0, I), i \sim U\{1, T\}} \left[\frac{SNR(t-l) - SNR(t)}{l} \| x - \widehat{x}_\theta(z_t; t) \|^2 \right] \\
 &= \lim_{T \rightarrow \infty} \frac{1}{2} E_{\epsilon \sim N(0, I), i \sim U\{1, T\}} \left[\frac{SNR(t-l) - SNR(t)}{l} \| x - \widehat{x}_\theta(z_t; t) \|^2 \right] \\
 &= \frac{1}{2} E_{\epsilon \sim N(0, I)} \left[\int_0^1 \frac{-dSNR(t)}{dt} \| x - \widehat{x}_\theta(z_t; t) \|^2 dt \right] \\
 &= -\frac{1}{2} E_{\epsilon \sim N(0, I), t \sim U(0, 1)} [SNR'(t) \| x - \widehat{x}_\theta(z_t; t) \|^2 dt]
 \end{aligned}$$

برای مطالعه بیشتر می توانید این دو مقاله را مطالعه فرمایید.

[مقاله اول](#)

[مقاله دوم](#)

سوال سوم: Score Matching (۲۰ نمره)

در این سوال به بررسی روش score matching در تولید نمونه جدید خواهیم پرداخت. همانطور که در کلاس درس نیز بیان شد، روش های مبتنی بر امتیاز یکی دیگر از روش هایی هستند که به منظور تولید نمونه جدید برای توزیع مورد نظر خود مورد استفاده قرار می گیرند. این روش ها از دو جزء اصلی تشکیل شده اند: Langevin dynamics و score-matching که در ادامه با هر کدام از این قسمت ها بیشتر آشنا خواهیم شد.

۱. Langevin Dynamics

فرض کنید که توزیع $p(x)$ در دسترس است و قصد داریم از این توزیع نمونه برداری کنیم. Langevin dynamics یک فرآیند تکرار شونده^۱ است که در انجام این کار به ما کمک می کند و به صورت عبارت زیر تعریف می شود:

$$x_{t+1} = x_t + \delta \nabla_x \log p(x_t) + \sqrt{2\delta} \epsilon; \epsilon \sim N(0, I) \quad (۴)$$

که در عبارت بالا x_1 نقطه شروع اولیه و δ اندازه گام است. همچنین t در بازه $[1 \dots T]$ قرار دارد.

(آ) فرض کنید در عبارت ۴ ترمی که مربوط به نویز است وجود نداشته باشد، در این صورت نشان دهید که عبارت حاصل معادل شروع از یک نقطه تصادفی در فضای R^d و رسیدن رسید به یکی از قله های توزیع $p(x)$ است. همچنین در مورد تفاوت عبارت حاصل با maximum likelihood نیز به صورت مختصر توضیح دهید. (۳ نمره)

^۱Iterative

Answer: The peak is a special place because it is where the probability is the highest. So, if we say that a sample x is drawn from a distribution $p(x)$, certainly the “optimal” location for x is where $p(x)$ is maximized. If $p(x)$ has multiple local minima, any one of them would be fine. So, naturally, the goal of sampling is equivalent to solving the optimization:

$$x^* = \operatorname{argmax}_x \log p(x).$$

To emphasize that this is not maximum likelihood estimation. In maximum likelihood, the data point x is fixed but the model parameters are changing. Here, the model parameters are fixed but the data point is changing. The table below summarizes the difference:

| Problem | Sampling | Maximum Likelihood |
|---------------------|---|---|
| Optimization target | A sample x | Model parameter θ |
| Formulation | $x^* = \operatorname{argmax}_x \log p(x; \theta)$ | $\theta^* = \operatorname{argmax}_\theta \log p(x; \theta)$ |

The easiest way to solve this optimization is, of course, gradient descent. For $\log p(x)$:

$$x_{t+1} = x_t + \delta \nabla_x \log p(x_t)$$

Where $\nabla_x \log p(x)$ denotes the gradient of $\log p(x)$ evaluated at x_t , and the δ is the step size. We use “+” instead of “-” because we are solving a maximization problem.

(ب) بر اساس نتیجه بدست آمده در قسمت آ می توان گفت که عبارت بدون ترم نویز، می تواند به ما در رسیدن به یکی از قله های توزیع $p(x)$ کمک کند، یا به عبارت دیگر داده بدست آمده احتمال بالایی خواهد داشت، در اینصورت دلیل اضافه کردن نویز در عبارت ۴ چیست؟ همچنین رفتار Langevin dynamic را در زمانی که قله یافت شده تیز یا صاف باشد را توصیف کنید. (۲ نمره)

Answer: The key is that we are not interested in solving the optimization problem. Instead, we are more interested in sampling from a distribution. By introducing the random noise. We randomly pick a sample that is closer to the peak and we will move left and right slightly. If the curvature around the peak is sharp, we will concentrate most of the steady state points x_T there. If the curvature around the peak is flat, we will spread around.

۲. Score Matching Techniques

یکی از چالش های تولید نمونه جدید، در درسترس نبودن توزیع نمونه ها یا همان $p(X)$ و عدم امکان استفاده از Langevin dynamic برای نمونه برداری از آن است. به منظور استفاده از این روش و رفع مشکل مذکور، سعی می کنند عبارت $\nabla_x \log p_\theta(x)$ را با استفاده از یک شبکه عصبی تخمین بزنند.

$$s_\theta(x) = \nabla_x \log p_\theta(x) \quad (۵)$$

به عبارت ۵ تابع امتیاز^۲ نیز گفته می شود. روش های مختلفی برای بدست آوردن تابع امتیاز وجود دارد که تفاوت آنها در تابع هزینه مدل است. فرض کنید دادگان آموزشی به صورت $X = \{x^{(1)}, \dots, x^{(m)}\}$ باشد و تابع هزینه نیز به صورت زیر تعریف شده باشد:

$$J_1(\theta) = E_{q(x)} \left[\frac{1}{2} \|s_\theta(x) - \nabla_x \log q(x)\|^2 \right] \quad (۶)$$

یک راحل این است که سعی کنیم از تخمین توزیع با استفاده از کرنل بهره ببریم:

$$q(x) = \frac{1}{M} \sum_{i=1}^M K(x | x^{(i)}) \quad (۷)$$

²Score function

در عبارت ۷ $K(\cdot)$ تابع کرنل است. یکی از کرنل‌هایی که می‌توان برای این کار مورد استفاده قرار داد، کرنل گاوسی^۳ است.

$$K(x | x^{(i)}) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^d \exp \left(-\frac{\|x - x^{(i)}\|^2}{2\sigma^2} \right) \quad (۸)$$

(آ) نشان دهید که با استفاده از کرنل عبارت ۸ و همچنین استفاده از توزیع تخمین زده شده $q(x)$ می‌توان به رابطه زیر رسید و با استفاده از این رابطه می‌توان مدل را آموزش داد. (۴ نمره)

$$\nabla_x \log q(x) = \frac{\sum_{i=1}^M \frac{1}{\sigma^2} (x^{(i)} - x) K(x | x^{(i)})}{\sum_{i=1}^M K(x | x^{(i)})}$$

Answer:

$$\begin{aligned} \nabla_x \log q(x) &= \nabla_x \log \frac{1}{M} \sum_{i=1}^M K(x | x^{(i)}) \\ &= \nabla_x \log \sum_{i=1}^M K(x | x^{(i)}) + \nabla_x \log \frac{1}{M} \\ &= \frac{\nabla_x \sum_{i=1}^M K(x | x^{(i)})}{\sum_{i=1}^M K(x | x^{(i)})} \\ &= \frac{\sum_{i=1}^M \nabla_x K(x | x^{(i)})}{\sum_{i=1}^M K(x | x^{(i)})} \\ &= \frac{\sum_{i=1}^M \nabla_x \left[\left(\frac{1}{\sqrt{2\pi}\sigma} \right)^d \exp \left(-\frac{\|x - x^{(i)}\|^2}{2\sigma^2} \right) \right]}{\sum_{i=1}^M K(x | x^{(i)})} \\ &= \frac{\sum_{i=1}^M \left[\frac{(x^{(i)} - x)}{\sigma^2} \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^d \exp \left(-\frac{\|x - x^{(i)}\|^2}{2\sigma^2} \right) \right]}{\sum_{i=1}^M K(x | x^{(i)})} \\ &= \frac{\sum_{i=1}^M \left[\frac{(x^{(i)} - x)}{\sigma^2} K(x | x^{(i)}) \right]}{\sum_{i=1}^M K(x | x^{(i)})} \end{aligned}$$

(ب) با وجود اینکه می‌توان با استفاده از کرنل گاوسی ۸ شبکه را آموزش داد، معایت این روش را بیان کنید. (۱ نمره)

Answer: The issue of explicit score matching is that the kernel density estimation is a fairly poor non-parametric estimation of the true distribution. Especially when we have limited number of samples and the samples live in a high dimensional space, the kernel density estimation performance can be poor. Moreover, this objective requires evaluating for each training step, which also fails to scale well to large datasets.

یک راه‌حل دیگر برای حل مسئله بالا استفاده از denoising score matching است. در این روش تابع هزینه به صورت زیر تعریف می‌شود:

$$J_2(\theta) = E_{q(x, x_0)} \left[\frac{1}{2} \|s_\theta(x) - \nabla_x \log q(x | x_0)\|^2 \right] \quad (۹)$$

(ج) نشان دهید بین دو تابع هزینه نوشته شده در عبارت‌های ۶ و ۹ رابطه زیر برقرار است. (۸ نمره)

$$J_2(\theta) = J_1(\theta) + C$$

³Gaussian kernel

Answer:

$$\begin{aligned} J_1(\theta) &= E_{q(x)} \left[\frac{1}{2} \| s_\theta(x) - \nabla_x \log q(x) \|^2 \right] \\ &= E_{q(x)} \left[\frac{1}{2} \| s_\theta(x) \|^2 - s_\theta(x)^T \nabla_x \log q(x) + \frac{1}{2} \| \nabla_x \log q(x) \|^2 \right] \end{aligned}$$

The last term is independent of θ , Let's focus on the second term:

$$\begin{aligned} &= E_{q(x)} [s_\theta(x)^T \nabla_x \log q(x)] = \int (s_\theta(x)^T \nabla_x \log q(x)) q(x) dx \\ &= \int \left(s_\theta(x)^T \frac{\nabla_x q(x)}{q(x)} \right) q(x) dx = \int s_\theta(x)^T \nabla_x q(x) dx \end{aligned}$$

Next, we consider conditioning by recalling $q(x) = \int q(x_0) q(x|x_0) dx_0$

$$\begin{aligned} \int s_\theta(x)^T \nabla_x q(x) dx &= \int s_\theta(x)^T \nabla_x \left(\int q(x_0) q(x|x_0) dx_0 \right) dx \\ &= \int s_\theta(x)^T \left(\int q(x_0) \nabla_x q(x|x_0) dx_0 \right) dx \\ &= \int s_\theta(x)^T \left(\int q(x_0) \nabla_x q(x|x_0) \times \frac{q(x|x_0)}{q(x|x_0)} dx_0 \right) dx \\ &= \int s_\theta(x)^T \left(\int q(x_0) \frac{\nabla_x q(x|x_0)}{q(x|x_0)} \times q(x|x_0) dx_0 \right) dx \\ &= \int s_\theta(x)^T \left(\int q(x_0) \nabla_x \log q(x|x_0) \times q(x|x_0) dx_0 \right) dx \\ &= \iint q(x, x_0) s_\theta(x)^T \nabla_x \log q(x|x_0) dx_0 dx \\ &= E_{q(x, x_0)} [s_\theta(x)^T \nabla_x \log q(x|x_0)] \end{aligned}$$

Now, if we put all together:

$$J_1(\theta) = E_{q(x)} \left[\frac{1}{2} \| s_\theta(x) \|^2 \right] - E_{q(x, x_0)} [s_\theta(x)^T \nabla_x \log q(x|x_0)] + C_1$$

If we compare that with $J_2(\theta)$ then we have:

$$\begin{aligned} J_2(\theta) &= E_{q(x, x_0)} \left[\frac{1}{2} \| s_\theta(x) - \nabla_x \log q(x|x_0) \|^2 \right] \\ &= E_{q(x, x_0)} \left[\frac{1}{2} \| s_\theta(x) \|^2 - s_\theta(x)^T \nabla_x \log q(x|x_0) + \frac{1}{2} \| \nabla_x \log q(x|x_0) \|^2 \right] \end{aligned}$$

Again, the last term is independent of θ :

$$J_2(\theta) = E_{q(x, x_0)} \left[\frac{1}{2} \| s_\theta(x) \|^2 \right] - E_{q(x, x_0)} [s_\theta(x)^T \nabla_x \log q(x|x_0)] + C_2$$

Therefore, we can write:

$$J_2(\theta) = J_1(\theta) - C_1 + C_2$$

(د) حال که نشان داده شد تابع هزینه denoising score matching معادل با تابع هزینه اصلی است، با استفاده از فرض زیر ابتدا تابع هزینه denoising را دوباره بازنویسی کنید و سپس فرآیند آموزش را توضیح دهید. در نهایت با استفاده از Langevin dynamics برای مدل بدست آمده نمونه جدید تولید کنید. (۲ نمره)

$$q(x | x_0) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^d \exp \left(-\frac{\|x - x_0\|^2}{2\sigma^2} \right)$$

Answer:

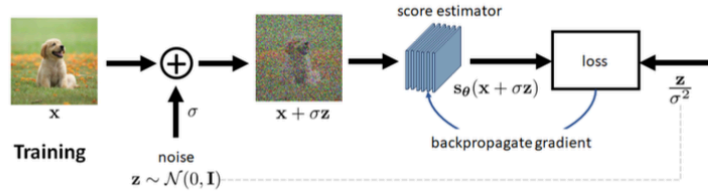
$$\nabla_x \log q(x|x_0) = \nabla_x \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^d \exp \left(\frac{-\|x - x_0\|^2}{2\sigma^2} \right) = \nabla_x \left\{ -\frac{\|x - x_0\|^2}{2\sigma^2} - C \right\} = -\frac{x - x_0}{\sigma^2} = -\frac{z}{\sigma^2}$$

As a result, the loss function of the denoising score matching becomes:

$$J_2(\theta) = E_{q(x, x_0)} \left[\frac{1}{2} \|s_\theta(x) + \frac{z}{\sigma^2}\|^2 \right] = E_{q(x_0)} \left[\frac{1}{2} \|s_\theta(x_0 + \sigma z) + \frac{z}{\sigma^2}\|^2 \right]$$

We have written the last expression based on reparameterization trick.

The quantity $x_0 + \sigma z$ is adding noise to a clean image. The score function s_θ is supposed to take this noisy image and predict the noise. Predicting the noise is equivalent to denoise, because any denoised image plus the predicted noise will give us the noisy observation. The architecture is like that:



The training step can simply described as follows:

$$\theta^* = \operatorname{argmin}_\theta \frac{1}{M} \sum_{i=1}^M \frac{1}{2} \|s_\theta(x^{(i)} + \sigma z^{(i)}) + \frac{z^{(i)}}{\sigma^2}\|^2 \text{ where } z^{(i)} \sim \mathcal{N}(0, I)$$

For inference, we assume that we have already trained the score estimator s_θ . To generate an image, we perform the following procedure for $t=1, \dots, T$:

$$x_{t+1} = x_t + \delta s_\theta(x_t) + \sqrt{2\delta} z_t; \text{ where } z_t \sim \mathcal{N}(0, I)$$

برای مطالعه بیشتر می توانید این دو مقاله را مطالعه فرمایید.

[مقاله اول](#)

[مقاله دوم](#)