### بسم الله الرحمن الرحيم



یادگیری عمیق نیمسال دوم ۰۳-۰۴ مدرس: مهدیه سلیمانی

دانشگاه صنعتی شریف دانشکدهی مهندسی کامپیوتر

تمرین چهارم ددلاین تمرین : ۹ خرداد

- برای ارسال هر تمرین تا ساعت ۲۳:۵۹ روز ددلاین فرصت دارید. مهلت تاخیر (مجاز و غیر مجاز) برای این تمرین، ۷ روز است (یعنی حداکثر تاریخ ارسال تمرین ۱۶ خرداد است).
- در هر کدام از سوالات، اگر از منابع خارجی استفاده کردهاید باید آن را ذکر کنید. در صورت همفکری با افراد دیگر هم باید نام ایشان را در سوال مورد نظر ذکر نمایید.
- پاسخ تمرین باید ماحصل دانستههای خود شما باشد. در صورت رعایت این موضوع، استفاده از ابزارهای هوش مصنوعی با ذکر نحوه و مصداق استفاده بلامانع است.
  - پاسخ ارسالی واضح و خوانا باشد. در غیر این صورت ممکن است منجر به از دست دادن نمره شود.
  - پاسخ ارسالی باید توسط خود شما نوشته شده باشد. به اسکرینشات از منابع یا پاسخ افراد دیگر نمرهای تعلق نمیگیرد.
- در صورتی که بخشی از سوالها را جای دیگری آپلود کرده و لینک آن را قرار داده باشید، حتما باید تاریخ آپلود مشخص و قابل اتکا باشد.
- محل بارگذاری سوالات نظری و عملی در هر تمرین مجزا خواهد بود. به منظور بارگذاری بایستی تمارین تئوری در یک فایل pdf با نام pdf با نام [Last-Name][Student-Id].pdf بارگذاری شوند. (یپ با نام Last-Name][Student-Id].zip بارگذاری شوند.
- در صورت وجود هرگونه ابهام یا مشکل، در کوئرای درس آن مشکل را بیان کنید و از پیغام دادن مستقیم به دستیاران آموزشی خودداری کنید.
  - طراحان این تمرین: پرهام رضایی، امیرحسین حاجیمحمدرضایی، علیرضا فرج تبریزی، محمدجواد احمدپور

# بخش نظری (۱۰۰ نمره)

# پرسش ۱. طراحی گام به گام VAE (۲۰ نمره)

در این تمرین قصد داریم به صورت گام به گام یک مدل VAE را توسعه دهیم. به دلیل وابستگی بین صورت سوال هر بخش به جواب بخش قبل، در صورت هرگونه ابهام یا سوال میتوانید در تلگرام ابهامات خود را بر طرف نمایید (ایدی طراح: Alireza\_1197@). در ادامه پیشفرضهای مسئله آمده است:

#### فرضيات:

ما مدل خود را بر پایهی یک معماری encoder-decoder با یک متغیر نهان پیوسته توسعه خواهیم داد. مؤلفههای کلیدی عبارتند از:

- داده ها: مجموعه ای از مشاهدات  $X^{(i)} \in \mathbb{R}^n$  که در آن هر داده ی  $X^{(i)} \in \mathbb{R}^n$  است.
- فضای نهان (Latent Space): یک متغیر نهان  $Z \in \mathbb{R}^m$  که اطلاعات فشرده شده ای از ورودی را نمایش می دهد.

• مدل مولد (Decoder): توزیع خروجی را یک توزیع گوسی به شرط توابعی از متغیر نهان تعریف میکند:

$$p_{\theta}(X \mid Z) = \mathcal{N}(X \mid \mu_{\theta}(Z), \Sigma_{\theta}(Z))$$

که در آن Decoder نگاشت Z به یارامترهای خروجی را انجام می دهد:

$$\mu_{\theta}(Z) = NN_{\alpha}(Z), \quad \Sigma_{\theta}(Z) = diag(\sigma(NN_{\beta}(Z)))$$

• توزیع پیشین (Prior): متغیر پنهان از یک توزیع گاوسی چندمتغیره نمونه گیری می شود:

$$p_{\theta}(Z) = \mathcal{N}(Z \mid \mu_z, \sigma_z^2 I)$$

• پارامترهای مدل: مجموعه پارامترهای قابل آموزش  $\{\mu_z,\sigma_z^2,\alpha,\beta\}$  که  $\alpha$  و  $\beta$  وزنهای شبکه Decoder

# گام ۱: برآورد درستنمایی (Likelihood Estimation)

در مدلهای مولد، هدف بیشینه کردن درستنمایی دادههای مشاهدهشده در توزیع خروجی مدل است.

(الف) یک عبارت برای **درستنمایی حاشیهای**  $p_{\theta}(X)$  برای یک نمونهی داده X استخراج کنید.

(ب) چالش اصلی در محاسبه یا بهینه سازی مستقیم  $p_{\theta}(X)$  در عمل چیست؟

(ج) از آنجا که محاسبه ی  $p_{\theta}(X)$  عملی نیست، یک تقریب برای آن پیشنهاد دهید. راهنمایی: سعی کنید درستنمایی حاشیه ای را به صورت امید ریاضی بیان کنید.

### گام ۲: کاهش واریانس با Importance Sampling

چالش اصلی در برآورد  $p_{\theta}(X)$  از طریق نمونه گیری، **واریانس بالا** است، زیرا نمونههای Z ممکن است از نواحیای بیایند که به داده ی مشاهده شده ی X ارتباطی ندارند.

استفاده می کنیم.  $p_{\theta}(X)$  را با استفاده از importance sampling استفاده می کنیم.  $p_{\theta}(X)$  را با استفاده از یک توزیع پیشنهادی q(Z) بازنویسی کنید.

برای توزیع پیشنهادی q(Z) معرفی کنید. q(Z) معرفی کنید.

# گام ۳: از Likelihood به ELBO

برای توجیه انتخاب توزیع پیشنهادی (Proposal Distribution)، دوباره به تخمین تابع هدف  $p_{ heta}(X)$  بازمی گردیم.

(الف) چرا در یادگیری ماشین معمولاً با  $\log p_{\theta}(X)$  کار میکنیم و نه با خود  $p_{\theta}(X)$ ?

ب کران پایین برای یک کوان پایین برای استفاده از Jensen's inequality یک کوان پایین برای  $\log p_{ heta}(X)$  با استخراج کنید (این کران به عنوان Evidence Lower Bound یا ELBO شناخته می شود).

(ج) آیا بیشینه کردن این کران پایین در حین آموزش بهتنهایی کافی است؟ چه مشکلاتی ممکن است ایجاد شود؟

(د) تفاضل بین  $\log p_{ heta}(X)$  و کران پایین استخراج شده چیست?

راهنمایی: این فاصله را به صورت KL divergence بنویسید.

# گام ۴: پارامتری کردن و بهینهسازی Posterior Approximation

 $p_{\theta}(Z \mid X)$  واقعی posterior تا اینجا انتخاب توزیع پیشنهادی را توجیه کردیم، اما با یک چالش جدید مواجه هستیم: posterior واقعی KL divergence غیرقابل محاسبه است. بنابراین نمی توانیم مستقیماً KL divergence را کمینه کنیم. برای رفع این مشکل باید فرم توزیع تقریبی q(Z) را مشخص کنیم. یک انتخاب رایج، توزیع گاوسی است:

$$q_{\lambda}(Z) = \mathcal{N}(Z \mid \mu, \sigma^2 I), \quad \lambda = \{\mu, \sigma^2\}$$

(الف) چگونه می توان KL divergence را به صورت غیر مستقیم کمینه کرد؟

راهنمایی: به گرادیان  $\nabla_{\lambda} \log p_{\theta}(X)$  فکر کنید.

(ب) گرادیان  $\nabla_{\theta}$  ELBO و تخمین Monte Carlo آن را محاسبه کنید.

Monte را نیز محاسبه کنید؟ چه چالشهایی در هنگام تخمین این گرادیان با نمونه گیری (ج) آیا می توانید  $\nabla_{\lambda}$  ELBO وجود دارد و چگونه می توان آنها را حل کرد؟

راهنمایی: بررسی کنید چگونه می توان نمونه گیری  $Z \sim q_{\lambda}(Z)$  را طوری نوشت که تصادفی بودن از پارامترهای  $\lambda$  جدا شود.

#### گام ۵: تعمیم به کل مجموعه داده

تمام محاسبات قبلی بر اساس یک نمونه داده انجام شدند. اما اگر بخواهیم  $\log p_{\theta}(X)$  را روی کل مجموعه داده بیشینه کنیم، چه تغییراتی ایجاد می شود؟

(الفِ) این کار چه تأِثیری روی **تابعِ هدف** و **گرادیانها** نسبت به  $\theta$  و  $\lambda$  دارد؟

 $q_{\lambda}(Z)$  به سادگی قابل تعمیم به کل مجموعه هستند. اما  $\lambda$ ، که پارامترهای توزیع تقریبی  $q_{\lambda}(Z)$  به سادگی قابل تعمیم به کل مجموعه هستند. اما  $\lambda$  به خورد که قابلیت تعمیم به همه داده ها را است، برای هر داده به طور جدا تعریف شده. چگونه میتوان  $\lambda$  را طوری مدل کرد که قابلیت تعمیم به همه داده الله داشته باشد؟

. راهنمایی: اینجا نقش شبکه encoder مطرح می شود (پارامترهای آن را با  $\phi$  نامگذاری کنید).

### گام ۶: ELBO نهایی و آموزش مدل

اكنون كه كل هدف VAE استخراج شده است، فرآيند آموزش را توصيف كنيد.

(الف) نشان دهید که ELBO را می توان به صورت زیر نوشت و مفهوم هر ترم از عبارات را توضیح دهید:

$$\mathbb{E}_{z \sim q_{\phi}(z|x)} \left[ \log p_{\theta}(x \mid z) \right] - \text{KL} \left( q_{\phi}(z \mid x) \parallel p_{\theta}(z) \right)$$

(ب) حال كه مدل VAE كامل شده است، الگوريتم آموزش آن را بنويسيد. در اين الگوريتم نشان دهيد:

- چگونه از encoder و decoder استفاده می شود؟
  - نمونه گیری چگونه انجام میشود؟
  - تابع هدف چگونه بهینه میشود؟

# پرسش ۲۰ Hierarchical VAE (۲۰ نمره)

یک مدل Hierarchical Variational Autoencoder (HVAE) با دو لایه را در نظر بگیرید که دارای متغیرهای نهان  $z_2$  و است، به طوری که  $z_1$  لایه یپایین و  $z_2$  لایه یبالا میباشد. مدل مولد به صورت زیر تعریف می شود:

$$p(x, z_1, z_2) = p(x \mid z_1) p(z_1 \mid z_2) p(z_2)$$

تقریب پسین (Encoder) به صورت زیر تعریف می شود:

$$q(z_1, z_2 \mid x) = q(z_1 \mid x) q(z_2 \mid z_1)$$

#### (الف)

کران پایین شواهد (ELBO) را به صورت L(x) برای این مدل HVAE استخراج کنید. (ELBO) را به شکل امید لال پایین شواهد (ELBO) را به صورت  $q(z_1 \mid x)q(z_2 \mid z_1)$  است که از ۳ بخش تشکیل شده است، لگاریتم درستنمایی و KL ریاضی روی توزیع های مربوطه.

**(ب)** 

معنای هر یک از اجزای موجود در ELBO بهدست آمده در قسمت (الف) را توضیح دهید و تأثیر آنها را بررسی کنید.

(ج)

در Hierarchical VAE ها، هدف این است که هر لایه ی نهان، اطلاعات مفیدی را کدگذاری کند که در بازسازی داده انقش داشته باشد. با این حال، در عمل اغلب مشاهده می شود که متغیرهای نهان لایه های بالایی مانند  $z_2$  معمولاً توسط مدل نادیده گرفته می شوند. به این معنا که خروجی encoder برای  $z_2$ ، یعنی  $q(z_2 \mid z_1)$ ، تقریباً مستقل از ورودی مده و با توزیع پیشین  $p(z_2)$  همراستا می شود.

- توضیح دهید که چرا این پدیده که به آن «فروریزش توزیع پسین» (Posterior Collapse) یا «متغیرهای نهان غیرفعال» (Inactice Latent Variables) گفته می شود، در زمان بهینه سازی ELBO رخ می دهد.
  - این رفتار چه چیزی را در مورد جریان اطلاعات در ساختار سلسلهمراتبی مدل نشان میدهد؟
- حداقل دو تکنیک (اعم از فرایند آموزش مدل یا تغییر در معماری مدل) پیشنهاد دهید که میتوانند به جلوگیری از این مسئله کمک کنند و توضیح دهید چرا این تکنیکها مؤثر هستند.

# پرسش ۳. GAN (۲۰ نمره)

در این سوال، قصد داریم بین مسئلهی "تشخیص داده واقعی و داده تولید شده" و مسئلهی "فاصله توزیع تصاویر واقعی و تصاویر تولید شده" ارتباطی برقرار کنیم.

توزیع تصاویر واقعی را با  $P_d$  و توزیع تصاویر تولید شده را با  $P_g$  نشان میدهیم. همچنین توزیع P که دادهای تصادفی را به ما میدهد، بهصورت زیر تعریف میشود:

$$P(x) = 0.5P_g(x) + 0.5P_d(x)$$

بنابراین، هر داده با احتمال مساوی از یکی از دو توزیع  $P_d$  یا  $P_d$  می آید. حال، مجموعهای از نمونهها (x,y) تشکیل می دهیم، به این ترتیب که:

- و را از توزیع P نمونه گیری می کنیم. x
- اگر x از توزیع واقعی  $P_d$  باشد، برچسب آن را y=1 قرار میدهیم، و در غیر این صورت (اگر از  $P_g$  آمده باشد) y=-1

اکنون،  $\mathcal D$  را مجموعه تمامی تمایزدهندهها در نظر میگیریم. هدف ما یافتن بهترین تمایزدهنده در این مجموعه است. فرض میکنیم این مجموعه هیچ محدودیتی ندارد و تمایزدهندهها دارای ظرفیت نامحدود هستند. منظور از ظرفیت نامحدود این است به ازای هر تابع y(x) دلخواه، تابعی در  $\mathcal D$  با چنین رفتاری وجود دارد. بنابراین اگر نقطه بهینه تابع برای هر x را بیابید میتوانید از وجود تمایزدهنده ای با این اتخاذها مطمئن باشید. بهترین تمایزدهنده  $\mathcal D$  در مجموعه  $\mathcal D$  تابع هزینه زیر را کمینه میکند:

$$R_l(D) = \mathbb{E}_{(x,y) \sim P}[l(yD(x))]$$

در اینجا، هزینه به صورت امیدریاضی روی داده های (x,y) تعریف شده است. اکنون مقدار بهینه این تابع را تعریف میکنیم:

$$R_l(\mathcal{D}) = \inf_{D \in \mathcal{D}} R_l(D)$$

#### بخش اول

ابتدا امیدریاضی موجود در رابطه بالا را به صورت انتگرال روی توزیع x بازنویسی کنید. عبارت شما باید مشابه فرم زیر ماشد:

$$\inf_{D \in \mathcal{D}} \left\{ \int \left( \left[ \mathsf{sth} \right] \, p_d(x) \, + \left[ \mathsf{sth} \right] \, p_g(x) \right) \, dx \right\}$$

سپس با توجه به ظرفیت نامحدود مجموعه  $\mathcal{D}$ ، بهینه عبارت بدست آمده را از روی بهینه نقطهای حساب کنید. توجه کنید که ظرفیت نامحدود چه تاثیری روی اینفیمم دارد.

نهایتا با جاگذاری بهینه و جابهجایی عبارتها نشان دهید که این مقدار برابر است با:

$$-\frac{1}{2}\int f\left(\frac{p_d(x)}{p_g(x)}\right)p_g(x)dx$$

توجه کنید که باید تابع f را به دست آورید. پس از رسیدن به این رابطه، بررسی کنید که تابع f به دست آمده بر حسب ورودی خود، محدب و نزولی باشد.

ازین به بعد عبارت بالاً را به استفاده از مفهوم واگرایی f به صورت

$$-\frac{1}{2}\mathbb{I}_f(\mathbb{P}_d,\mathbb{P}_{\eth})$$

نمایش میدهیم.

#### بخش دوم

در بخش قبل، نشان دادیم که این مسئله معادل کمینه کردن یک واگرایی بین دو توزیع است. اکنون میخواهیم بررسی کنیم که با انتخاب توابع هزینه مختلف l(x)، چه نوع واگراییهایی حاصل میشوند. برای هر تابع هزینه l(x):

- ۱. ابتدا محدب بودن و یکنوایی تابع f متناظر را بررسی کنید.
- ۲. سپس، تمایزدهنده بهینه، تابع f، و واگرایی متناظر را استخراج کنید.

## حالتهاي مورد بررسي:

- $l(x) = \mathbb{I}(x \leq 0)$  (الف
  - $l(x) = (1-x)^2$  •
- $l(x) = \log(1 + e^{-x})$  (z •

راهنمایی: واگرایی های متناظر توابع هزینه به صورت نامرتب در ادامه آمدهاند. میتوانید جهت بررسی درستی حل خود از این بخش استفاده کنید.

- Total Variation که منظور از TV تابع  $R(\mathcal{D}) = \frac{1}{2}(1 \mathbb{I}_{\mathsf{TV}}(\mathbb{P}_d, \mathbb{P}g))$  •
- است. Jensen Shannon که منظور از  $R(\mathcal{D}) = \log 2 \mathbb{I}_{\mathsf{JS}}(\mathbb{P}_d, \mathbb{P}_g)$ 
  - است.  $\frac{-4t}{t+1}$  است. که منظور از g تابع  $R(\mathcal{D})=1-\mathbb{I}_g(\mathbb{P}_d,\mathbb{P}_g)$  •

# پرسش ۴. یادگیری تابع امتیاز در Diffusion (۲۰ نمره)

در اسلایدهای درس دیدیم که هدف اصلی ما یادگیری تابع امتیاز است. یک راه ساده برای رسیدن به این هدف، تعریف تابع هزینهای بین مقدار حقیقی تابع امتیاز و خروجی شبکه عصبی است:

$$l_1(\theta) = \mathbb{E}_{q(x)} \left[ \frac{1}{2} ||s_{\theta}(x) - \nabla_x \log q(x)||_2^2 \right]$$

اما مشکل اینجاست که به  $\nabla_x \log q(x)$  دسترسی نداریم.

در ادامه یاد گرفتیم که اگر مدل بتواند نویز را تخمین بزند، میتوان از آن برای تخمین تابع امتیاز استفاده کرد. قبل از پرداختن به آن، رابطهای سادهتر را بررسی میکنیم.

بخش (الف): ساده سازی تابع هزینه در این بخش میخواهیم مقدار ناشناخته تابع امتیاز را از تابع هزینه حذف کنیم. با استفاده از انتگرال جز به جز و حذف بخشهای مستقل از  $\theta$  نشان دهید که:

$$l_1(\theta) = \mathbb{E}_{q(x)} \left[ \frac{1}{2} ||s_{\theta}(x)||_2^2 + \text{Tr}(\nabla_x s_{\theta}(x)) \right] + C_1$$

که  $C_1$  مستقل از  $\theta$  است. در اثبات خود فرض کنید که وقتی  $x \to \infty$  داریم که  $q(x)s_{\theta}(x) \to 0$ . در پایان، توضیح دهید چرا این تابع هزینه در عمل ممکن است برای آموزش مناسب نباشد؟

بخش (ب): ارتباط با تخمین نویز اکنون هدف ما این است که نشان دهیم تابع هزینه بالا معادل با تابع هزینه ای است که بر اساس تخمین نویز تعریف می شود:

$$l_3(\theta) = \mathbb{E}_{x \sim q(x), \ \epsilon \sim \mathcal{N}(0, \mathcal{I})} \left[ \frac{1}{2} \left\| s_{\theta}(\underbrace{x + \sigma \epsilon}_{\tilde{x}}) + \frac{\epsilon}{\sigma} \right\|^2 \right]$$

جفت داده سالم و نویزی را  $(x, \tilde{x})$  مینامیم. طبق رابطه زیر، برای کرنل گاوسی داریم:

$$\frac{\partial \log q(\tilde{x}|x)}{\partial \tilde{x}} = \frac{1}{\sigma^2} (x - \tilde{x})$$

ر نتىجە:

$$l_3(\theta) = \mathbb{E}_{q(x,\tilde{x})} \left[ \frac{1}{2} \left\| s_{\theta}(\tilde{x}) - \frac{\partial \log q(\tilde{x}|x)}{\partial \tilde{x}} \right\|^2 \right]$$

ثابت كنيد كه:

$$l_3 = l + C$$

که C مستقل از  $\theta$  است.

اگر نیاز به راهنمایی دارید، میتوانید به پیوست مقاله زیر مراجعه کنید:

A Connection Between Score Matching and Denoising Autoencoders

اما سعی کنید مراحل اثبات را خودتان کامل نوشته و توضیح دهید.

پرسش ۵. آیا فرض مارکوف در Diffusion الزامی است؟ (۲۰ نمره)

در ساختار فروارد دیفیوژن، فرض کردیم که:

$$q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) = \prod_{t=1}^{T} q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$$

این فرض مارکوف باعث می شود که برای تولید، نیاز به انجام T مرحله به صورت ترتیبی داشته باشیم. برای درستی فرض گاوسی بودن می دانیم، T باید بزرگ باشد که باعث هزینه زمانی زیاد می شود. آیا می توانیم بدون فرض مارکوف بودن هم به همان توزیع حاشیه  $q(x_t \mid x_0)$  برسیم؟ دقت کنید که در فرم تابع هدف (Objective) نهایی، یعنی:

$$L(\epsilon_{\theta}) := \sum_{t=1}^{T} \gamma_{t} \mathbb{E}_{\mathbf{x}_{0} \sim q(\mathbf{x}_{0}), \, \boldsymbol{\epsilon}_{t} \sim \mathcal{N}(0, \mathbf{I})} \left[ \left\| \boldsymbol{\epsilon}_{\theta}^{(t)} \left( \sqrt{\bar{\alpha}_{t}} \mathbf{x}_{0} + \sqrt{1 - \bar{\alpha}_{t}} \boldsymbol{\epsilon}_{t} \right) - \boldsymbol{\epsilon}_{t} \right\|_{2}^{2} \right]$$

تنها این حاشیه تاثیر دارد و توزیع مشترک  $x_i$  ها تاثیری ندارد. در مدل استاندارد، این مارجین به صورت زیر است:

$$q(x_t \mid x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

بنابراین، میتوان ساختاری غیرمارکوفی تعریف کرد که همین مارجین را تولید کند. برای یک بردار  $\sigma \in \mathbb{R}^T_{\geq 0}$  توزیع اینفرنس را به صورت زیر تعریف میکنیم:

$$q_{\sigma}(x_{1:T} \mid x_0) := q_{\sigma}(x_T \mid x_0) \prod_{t=2}^{T} q_{\sigma}(x_{t-1} \mid x_t, x_0)$$

که در آن:

$$q_{\sigma}(x_T \mid x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_T}x_0, (1 - \bar{\alpha}_T)\mathbf{I})$$

: t > 1 و برای

$$q_{\sigma}(x_{t-1} \mid x_t, x_0) = \mathcal{N}\left(\sqrt{\bar{\alpha}_{t-1}}x_0 + \sqrt{1 - \bar{\alpha}_{t-1}} - \sigma_t^2 \cdot \frac{x_t - \sqrt{\bar{\alpha}_t}x_0}{\sqrt{1 - \bar{\alpha}_t}}, \sigma_t^2 \mathbf{I}\right)$$

الف) اثبات برابری مارجینها با استفاده از خواص توزیع گاوسی و به کمک استقرا ثابت کنید:

$$q_{\sigma}(x_t \mid x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

توجه کنید که فرایند فروارد ما دیگر مارکوف نیست، زیرا  $x_0$  به  $x_0$  نیز وابسته است:

$$q_{\sigma}(x_t \mid x_{t-1}, x_0) = \frac{q_{\sigma}(x_{t-1} \mid x_t, x_0) \, q_{\sigma}(x_t \mid x_0)}{q_{\sigma}(x_{t-1} \mid x_0)}$$

مشاهده جالب: اگر  $0 \to 0$ ، واریانس گاوسی به صفر میل میکند و  $x_{t-1}$  به صورت قطعی از  $x_t$  به دست میآید. بنابراین میتوان فرآیند جنریشن را به یک فرآیند قطعی تبدیل کرد که تنها منبع تصادفی آن نویز اولیه  $x_t$  است. این دیدگاه امکان نوعی تناظر بین فضای پنهان و خروجی را فراهم میکند، و دقت کنید که توزیع حاشیه ای که حاصل می شود همان توزیع مدل دیفیوژن تصادفی خواهد بود.

فرآیند جنریشن ابتدا با استفاده از عبارت زیر که مشابه مدل استاندارد است، نمونه ای نویزی ساخته و بدون دسترسی به  $x_0$  نویز آن را تخمین میزنیم:

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

با مدل تخمین نویز می توانیم نمونهی بدون نویز را به صورت زیر بازسازی کنیم:

$$f_{\theta}^{(t)}(x_t) := \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon_{\theta}^{(t)}(x_t)}{\sqrt{\bar{\alpha}_t}}$$

بنابراین ابتدا داریم که

$$p_{\theta}(x_T) = \mathcal{N}(0, \mathbf{I})$$

و پس از آن

$$p_{\theta}^{(t)}(x_{t-1} \mid x_t) = \begin{cases} \mathcal{N}(f_{\theta}^{(1)}(x_1), \sigma_1^2 \mathbf{I}) & t = 1 \\ q_{\sigma}(x_{t-1} \mid x_t, f_{\theta}^{(t)}(x_t)) & \text{غير اين صورت} \end{cases}$$
غير اين صورت

حال با این مدل جنریشن، آموزش به چه صورتی میشود؟ ابتدا دقت میکنیم که هدف آموزش مدل ما چه بود؟ مشابه قبل، ما به دنبال پارامترهایی هستیم که عبارت زیر اتخاذ شود.

$$\arg\max_{\theta} \mathbb{E}_{x_0 \sim q}[\log p_{\theta}(x_0)]$$

بنابراین میتوان همان استدلالهای مدل استاندارد را برای یافت تابع هزینه در این مدل نیز کرد.

ب تطابق تابع هدف (Objective) با دیفیوژن کلاسیک با تکرار مراحل ذکر شده در اسلایدهای درس و تحلیل ELBO نشان دهید که تابع هدف این مدل مولد با مدل بررسی شده در اسلایدها یکسان است (به جز یک ثابت مستقل از  $\theta$ ). در برابری فرض کنید که در تابع هدف مدل استاندارد تابع هزینه t های مختلف، ضریب برابری دارند.

پرسش: به نظر شما چگونه می توان از این نگاه برای کاهش زمان طولانی جنریشن استفاده کرد؟