



به نام خدا



دانشگاه تهران
دانشکده مهندسی برق و کامپیوتر
مدل‌های مولد عمیق

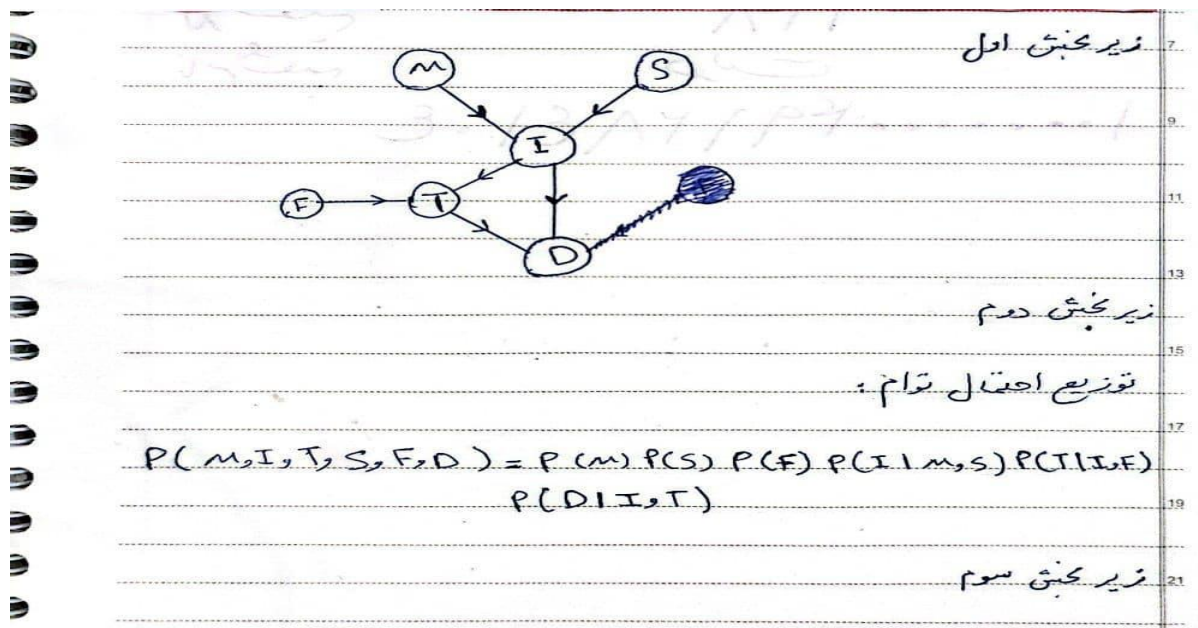
تمرین اول

| | |
|--------------------|-----------|
| نام و نام خانوادگی | محمد متقی |
| شماره دانشجویی | 159403018 |
| تاریخ ارسال گزارش | |

فهرست

سؤال اول

بخش اول



(A) نادرست، میان F و D مسیر فعالی وجود دارد. $F \rightarrow T \rightarrow D$: چون شرطی بر T نیست، مسیر باز است و این دو متغیر وابسته‌اند.

(B) درست، مسیرهای S به I به D و S به I به T به D همگی به‌خاطر شرط بر I بسته می‌شوند چون I گره‌ای غیر برخوردارکننده در میانه مسیر است؛ پس S و D (فصل و مرگ با دانستن شدت بیماری) به شرط I مستقل‌اند.

(C) درست، میان M و F فقط مسیرهایی وجود دارد که شامل برخوردارکننده‌ها هستند، مانند M به I به T از یک طرف و F از طرف دیگر که در T به هم می‌رسند. از آنجا که برخوردارکننده‌ها تا وقتی رویشان شرط نگذاریم مسیر را می‌بندند، این دو متغیر مستقل‌اند.

(D) نادرست، وقتی روی برخوردارکننده T شرط می‌گذاریم، مسیر M به I به T از یک سو و F از سوی دیگر باز می‌شود و وابستگی القا می‌گردد؛ پس استقلال از بین می‌رود..

(E) درست، شرط‌گذاری هم‌زمان بر D و I باعث بسته شدن همه مسیرهای میان M و T می‌شود؛ مسیر مستقیم M به I به T با شرط بر I بسته است و مسیر غیرمستقیم M به I به D به T نیز به‌دلیل وجود همان شرط‌ها بسته می‌ماند. به این ترتیب در این حالت M و T مستقل‌اند.

بخش دوم

زیربخش اول

در شبکه‌های بیزی، هر متغیر فقط به والد‌های خودش وابسته است پس توزیع توأم همه متغیرها حاصل ضرب احتمال هر متغیر مشروط به والد‌های خودش است.

$$p(O, S, A, T, B, M, C) = p(O)p(A)p(C)p(S|O)p(T|O, A)p(B|S)p(M|A, T, B)$$

زیربخش دوم

Markov blanket یک گره مجموعه‌ای از گره‌هاست که اگر مقدارشان را بدانی، آن گره از تمام گره‌های دیگر شبکه مستقل می‌شود. یعنی فرزندان، والدین و سایر والدین فرزندان آن گره.

$$MB(T) = \{O, A, M, B\}$$

زیربخش سوم

خیر، گراف مارکوف حاصل از این شبکه‌ی بیزی Perfect I-Map نیست چون در فرایند moralization، والد‌های مشترک هر گره به هم وصل می‌شوند و ساختارهای برخوردکننده مثل $A \rightarrow T \leftarrow O$ یا $A \rightarrow M \leftarrow T$ از بین می‌روند. این باعث می‌شود بعضی استقلال‌هایی که در شبکه‌ی بیزی وجود داشت (مثل استقلال بین O و A یا بین A و B هنگام شرط‌گذاری) دیگر در گراف مارکوف برقرار نباشد. در نتیجه گراف فقط یک I-Map معمولی است، نه Perfect I-Map، چون بخشی از استقلال‌های اصلی مدل جهت‌دار را از دست می‌دهد.

زیربخش چهارم

خیر، گراف مارکوف حاصل chordal نیست. در تعریف، گرافی chordal است که هیچ چرخه‌ای با چهار گره یا بیشتر بدون یال میان‌بر نداشته باشد. در گراف ما، پس از مارال‌سازی، چرخه‌هایی مانند $O - S - B - T - O$ یا $A - T - M - B - A$ وجود دارند که بین رأس‌های غیرمتوالی مثل O و B یا S و T یالی وجود ندارد. نبود این یال‌ها یعنی چرخه‌های بدون chord داریم، بنابراین گراف غیر chordal است و نمی‌توان آن را دقیقاً با یک شبکه‌ی بیزی معادل بازسازی کرد.

زیربخش پنجم

کلیک‌های ماکسیمال دقیقاً $\{A, O, T\}$ ، $\{A, T, M\}$ ، $\{T, B, M\}$ ، $\{O, S\}$ و $\{S, B\}$ هستند. مجموعه $\{O, S, B\}$ کلیک نیست چون یال‌های OB و ST وجود ندارند و بنابراین $\{O, S\}$ و $\{S, B\}$ هر کدام به‌تنهایی ماکسیمال باقی می‌مانند. در نتیجه توزیع توأم به‌صورت گیبس روی کلیک‌های ماکسیمال فاکتور می‌شود.

$$p(A, O, T, S, B, M, C) \propto \phi_{AOT}(A, O, T) \phi_{ATM}(A, T, M) \phi_{TBM}(T, B, M) \phi_{OS}(O, S) \phi_{SB}(S, B) \phi_C(C)$$

زیربخش ششم

الف) نگره C از بقیه جداست پس در فاکتور بندی توأم به‌صورت $p(C)$ درمی‌آید.

$$P(A, O, T, S, B, M, C) p(C) = P(A, O, T, S, B, M, C)$$

$$\sum_C p(C) P(A, O, T, S, B, M, C) = P(A, O, T, S, B, M)$$

* زیرا $p(C)$ خود یک توزیع نرمال است که $\sum_C p(C) = 1$ و در نتیجه با حذف $p(C)$ هیچ عامل دیگری تغییری نمی‌کند پس در BN ، حاشیه‌گیری دقیقاً منجر به حذف فاکتور $p(C)$ می‌شود.

$$P(A, O, T, S, B, M, C) = \frac{1}{Z} \phi_C(C) \phi(A, O, T, S, B, M) \quad \text{ب)}$$

$$\phi_C(C) \leq \frac{\phi}{Z} \quad \phi \text{ ضرب همی پتانسیل های مستقل اند } C \text{ می باشد. حال}$$

$$\text{ما ثابت } k \text{ در نظر می گیریم. } Z' = \frac{Z}{k} \text{ پس داریم:}$$

$$\phi(A, O, T, S, B, M) \propto \frac{1}{Z'} = P(A, O, T, S, B, M)$$

در نتیجه در MRF چون ϕ_C فقط ثابت مقیاس تولید می‌کند، آن ثابت در Z جذب می‌شود و ساختار $P(A, O, T, S, B, M)$ تغییری نمی‌کند.

بخش سوم

زیربخش اول

$\{A, B, C\}$ ، $\{B, C, D\}$ ، $\{C, D, F\}$ ، $\{B, E\}$ ، $\{E, F\}$ ، $\{E, G\}$.
 $p(A, B, C, D, E, F, G) \propto$

$$\phi_{ABC}(A, B, C) \phi_{BCD}(B, C, D) \phi_{CDF}(C, D, F) \phi_{BE}(B, E) \phi_{EF}(E, F) \phi_{EG}(E, G)$$

زیربخش دوم

الف) نادرست، A و G مستقل نیستند زیرا مسیر $A B E G$ باز است و وابستگی بین آن‌ها برقرار می‌شود.

ب) نادرست، F و A با داشتن مجموعه D و C مستقل نیستند. مسیرهای $A C F$ و $A B D F$ با شرط‌های C و D بسته می‌شوند اما مسیر $A B E F$ همچنان باز است چون روی B و E شرط نکرده‌ایم و این مسیر وابستگی را حفظ می‌کند.

ج) درست، G و C با داشتن E مستقل هستند زیرا تمام مسیرهای میان آن‌ها مانند $C F E G$ و $C B E G$ از E عبور می‌کنند و شرط بر E این مسیرها را می‌بندد.

د) درست، شرط دانستن B و C برای A کافی است و افزودن E اطلاعات جدیدی نمی‌دهد زیرا در گراف بدون جهت، A با شرط همسایه‌هایش B و C از سایر گره‌ها مستقل می‌شود.

زیربخش سوم

همانطور که در بخش قبل ثابت کردیم، ضرب یک ثابت مثبت در پتانسیل ثابت نرمال‌سازی Z را مقیاس می‌کند و شکل توزیع پس از نرمال‌سازی تغییر نمی‌کند؛ توزیع نهایی روی (A, \dots, G) همان قبلی است.

$$\begin{array}{l} \textcircled{Z} \quad P(z) = e^{-z} \quad z > 0 \\ \downarrow \\ \textcircled{X} \quad P(x|z) = ze^{-x} \quad x > 0 \\ P(z|x) = ? \end{array}$$

تقریب توزیع پسین $Eq(z) = \frac{\gamma}{\theta} z e^{-\theta z} \quad z > 0$ و $Eg(z) = \frac{\gamma}{\theta}$

$$P(z|x) \propto P(x|z) P(z) = ze^{-(\alpha+1)z}$$

ما به ثابت K پیدا کنیم به گونه ای که

$$\int_0^{\infty} K z e^{-(\alpha+1)z} dz = 1$$

$$\xrightarrow{z(\alpha+1)=u} \frac{1}{(\alpha+1)^2} \int_0^{\infty} \frac{u^2}{(\alpha+1)^2} e^{-u} \frac{du}{\alpha+1} = \frac{1}{\alpha+1} \int_0^{\infty} u^2 e^{-u} du = \frac{\gamma}{\alpha+1}$$

طبی فرض داریم: $q(z) = \theta^2 z e^{-\theta z}$

$$\int_0^{\infty} \theta^2 z e^{-\theta z} dz = \theta^2 \times \frac{1}{\theta^2} = 1$$

حالا این توزیع نیز $Eq(z) = \frac{\gamma}{\theta}$ است.

ما ضایع KL را نسبت به $q(z)$ میگیریم: $KL(q||P) = Eq \left[\log \frac{q(z)}{P(z|x)} \right]$

ما به این ضایع KL minimal میگیریم: $\log q(z) = \gamma \log \theta + \log z - \theta z$

$\log P(z|x) = \gamma \log (\alpha+1) + \log z - (\alpha+1)z$

$$\log \frac{q(z)}{P(z|x)} = \gamma (\log \theta - \log (\alpha+1)) + ((\alpha+1) - \theta) z$$

ما امید ریاضی تحت q میگیریم: $Eg[z] = \frac{\gamma}{\theta}$

$$KL(\theta) = \gamma (\log \theta - \log (\alpha+1)) + \gamma \left(\frac{\alpha+1}{\theta} - 1 \right)$$

ما مشتق میگیریم نسبت به θ :

$$\frac{dKL(\theta)}{d(\theta)} = \frac{\gamma}{\theta} - \frac{\gamma(\alpha+1)}{\theta^2} = 0 \rightarrow \gamma\theta - \gamma(\alpha+1) = 0$$

پس $\theta^* = \alpha+1$

برای $\theta = \alpha+1$ ما ضایع KL کمینه می شود و در این مثال خاص دقیقاً صفر

می شود به ازای این مقدار.

$$P(z|x) = (\alpha+1)^2 z e^{-(\alpha+1)z}$$

$$q(z) = \theta^2 z e^{-\theta z} \rightarrow \theta^* = \alpha+1$$

توزیع صافی
و q دقیقاً یکسان می شوند

سؤال دوم

بخش اول

زیربخش اول

به طور کلی می توان گفت در مدل های مولد نهفته، محاسبه مستقیم احتمال داده ها $\log p_{\theta}(x)$ ممکن نیست و برای رفع این مشکل از یک توزیع تقریبی $q_{\phi}(z|x)$ استفاده می کنیم و با نامساوی ینسن (Jensen) کران پایینی به دست می آوریم که به آن ELBO می گویند.

این کران جمع دو مؤلفه است. مؤلفه اول کیفیت بازسازی را می سنجد. مؤلفه دوم میزان انحراف توزیع کمکی از پرایر را جریمه می کند. هر چه کران بالاتر باشد مدل بهتر داده را توضیح می دهد.

$$E_{q(z|x)}[\log p(x|z)] - D_{KL}(q(z|x)||p(z))$$

The first term = reconstruction, the second term = regularization

زیربخش دوم

با توجه به توضیحات گیت هاب، dSprites یک دیتاست دوبعدی است که برای ارزیابی میزان Disentanglement در مدل های Unsupervised Representation Learning طراحی شده است. هر تصویر شامل یک شکل ساده در پس زمینه سیاه با ابعاد 64×64 پیکسل است و همه تصاویر به صورت رویه ای ساخته شده اند.

این تصاویر بر اساس 6 عامل نهفته مستقل (Independent Latent Factors) تولید شده اند که در دیتاست dSprites تمام ترکیب های ممکن از این 6 عامل دقیقاً یک بار ظاهر می شوند، بنابراین تعداد کل تصاویر برابر با 737,280 نمونه می شود.

فاکتور ها: color, shape, scale, orientation, position X, and position Y

هدف از ایجاد dSprites این بوده که مانند یک تست استاندارد برای ارزیابی Disentanglement عمل کند و بسنجد که آیا مدل های یادگیری بدون نظارت قادرند عوامل نهفته واقعی را به درستی از داده استخراج کنند یا نه و در مقاله معروف β -VAE: Learning basic visual concepts with a constrained variational framework منتشر شده در کنفرانس ICLR 2017، از این دیتاست برای محاسبه ی متریک های سنجش Disentanglement استفاده شده است.

Three random samples from dSprites



شکل 1: سه نمونه تصادفی از دیتاست dSprites

زیربخش سوم

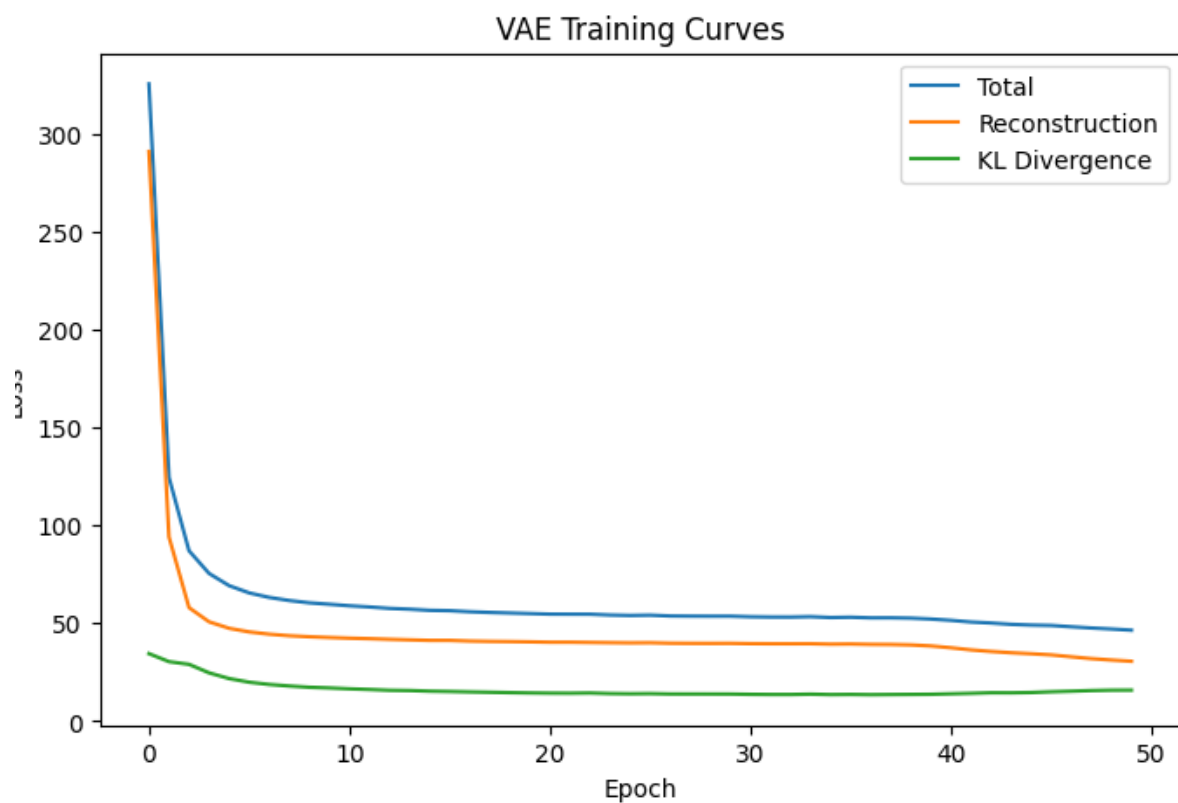
در مدل VAE برخلاف یک Autoencoder معمولی، خروجی بخش Encoder یک بردار ثابت نیست، بلکه یک توزیع احتمالی روی فضای نهفته است که معمولاً نرمال در نظر گرفته می‌شود. در واقع، به جای این که Encoder مستقیماً مقدار z را برگرداند، دو بردار میانگین و انحراف معیار تولید می‌کند و سپس برای هر نمونه، یک z تصادفی از این توزیع انتخاب می‌کنیم.

اما این عمل تصادفی (Sampling) درون شبکه باعث می‌شود مسیر گرادیان‌ها از Decoder به Encoder قطع شود، زیرا عمل نمونه‌گیری قابل مشتق‌گیری نیست و به همین دلیل آموزش مدل با Backpropagation دیگر ممکن نیست.

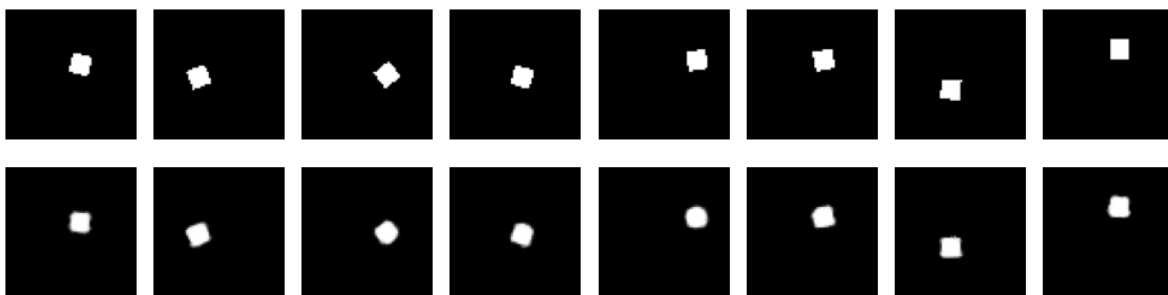
راه حل این مشکل استفاده از Reparameterization Trick است. در این روش، به جای آن که مستقیماً از توزیع نرمال نمونه بگیریم، یک نویز استاندارد را از توزیع نرمال استاندارد $N(0, I)$ می‌کشیم و سپس z را به صورت تابعی از سه پارامتر μ ، σ و ε بازنویسی می‌کنیم:

$$z = \mu + \sigma \odot \varepsilon, \varepsilon \sim \mathcal{N}(0, I)$$

به این ترتیب، تمام بخش‌های μ و σ که از شبکه می‌آیند، در مسیر گرادیان قرار می‌گیرند و تنها ε یک نویز ثابت است. این بازنویسی باعث می‌شود عملیات نمونه‌گیری به شکل یک عمل مشتق‌پذیر دربیاید و بتوان از Backpropagation استفاده کرد. در نتیجه، مشکل قطع شدن گرادیان‌ها حل می‌شود و Encoder و Decoder هر دو می‌توانند به طور همزمان با روش گرادیان نزولی آموزش ببینند.



شکل 2 آموزش VAE با استفاده از 10000 داده در Epoch 50



شکل 3 نمونه های تصادفی و بازسازی شده

نتایج آموزش مدل VAE نشان می‌دهد که فرآیند یادگیری به‌صورت پایدار و مؤثر انجام شده است. در نمودار تابع هزینه، مشاهده می‌شود که Total Loss در چند epoch نخست به‌شدت کاهش یافته و سپس به مقدار تقریباً ثابت حدود ۴۶ رسیده است، که نشان‌دهنده همگرایی مدل است. کاهش سریع در مراحل ابتدایی به این معناست که مدل در همان دوره‌های اولیه ساختار کلی داده‌ها را یاد گرفته و در ادامه فقط جزئیات بازسازی را بهبود داده است. همچنین، Reconstruction Loss به حدود ۳۰ کاهش یافته که حاکی از بازسازی مناسب شکل‌ها و موقعیت اجسام در تصاویر dSprites است. در مقابل، مقدار KL Divergence در حدود ۱۵ مانده که نشان می‌دهد فضای نهان منظم و نزدیک به توزیع نرمال استاندارد باقی مانده است.

از نظر کیفی، بازسازی‌های نهایی بسیار مشابه ورودی‌ها هستند و مدل موفق شده ویژگی‌های اصلی نظیر شکل و موقعیت را به‌درستی حفظ کند. در تصاویر بازسازی‌شده تفاوت چشمگیری با نمونه‌های اصلی دیده نمی‌شود و تنها مقدار کمی افت جزئیات در لبه‌ها وجود دارد، که در VAE‌های پایه طبیعی است. در مجموع، مدل تعادل مناسبی میان دقت بازسازی و منظم‌بودن فضای نهان برقرار کرده است.

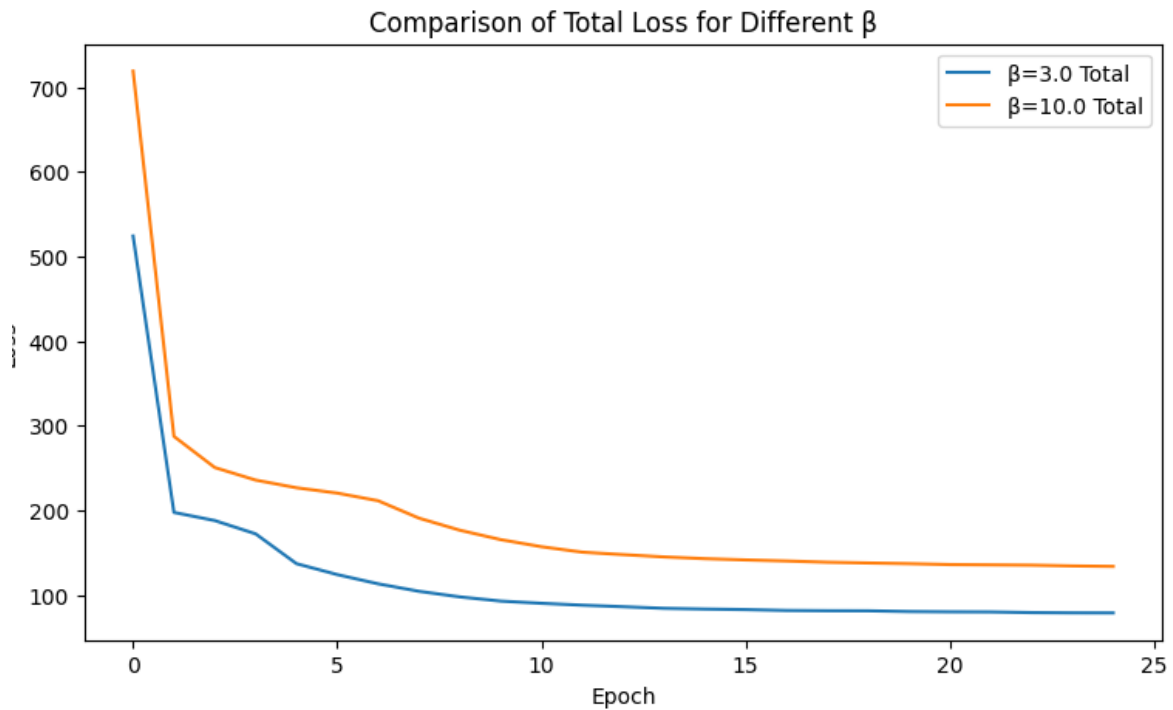
زیربخش پنجم

ترم اول یعنی بازسازی (Reconstruction Term) تلاش می‌کند خروجی مدل تا حد ممکن شبیه ورودی باشد. ترم دوم یعنی KL Divergence باعث می‌شود توزیع نهفته $q(z|x)$ خیلی از توزیع پرایر $p(z)$ معمولاً $\text{Normal}(0,1)$ دور نشود.

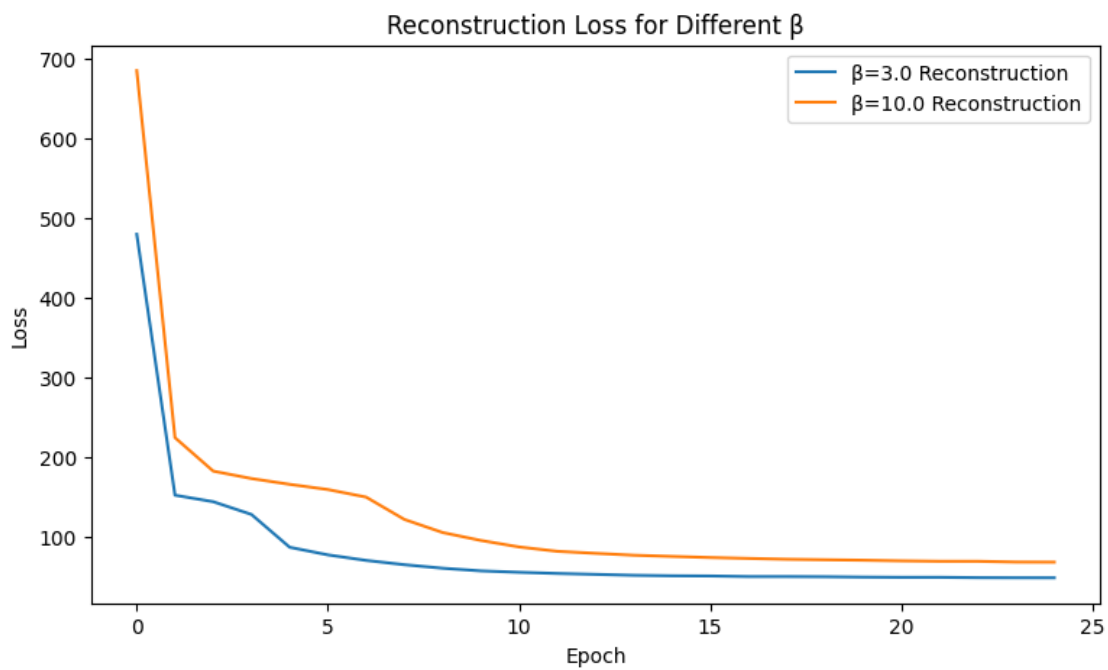
اگر $\beta = 1$ باشد، مدل دقیقاً همان VAE کلاسیک است ولی وقتی $\beta > 1$ باشد، مقدار ترم KL در تابع هدف بزرگتر می‌شود و مدل مجازات بیشتری برای تفاوت بین $q(z|x)$ و $p(z)$ در نظر می‌گیرد. به زبان ساده مدل را مجبور می‌کنیم که متغیرهای نهفته‌اش (z) تا حد ممکن شبیه توزیع نرمال استاندارد و از هم مستقل باشند.

β -VAE همان VAE است که در آن با بزرگتر کردن β ، به مدل فشار بیشتری وارد می‌کنیم تا فضای نهفته‌اش منظم‌تر و مستقل‌تر باشد. این باعث می‌شود عوامل پنهان داده‌ها به صورت تفکیک‌شده در متغیرهای z نمایش داده شوند.

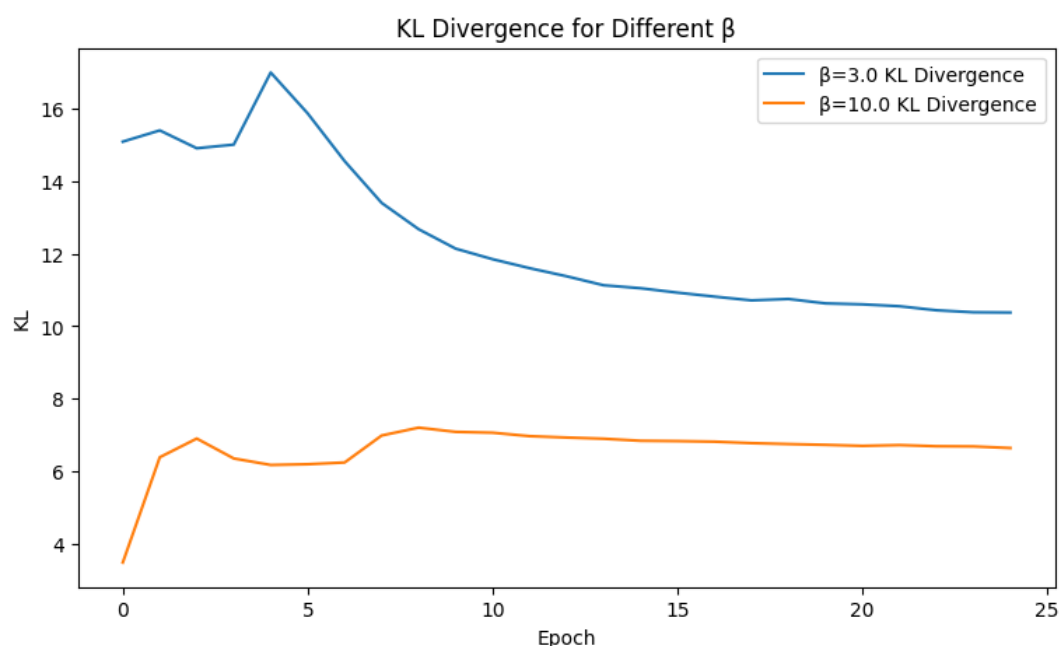
وقتی $\beta > 1$ باشد، مدل بیشتر وادار می‌شود که $q(z | x)$ به $p(z)$ نزدیکتر شود.



شکل 4 مقایسه $loss$ بین دو مقدار β

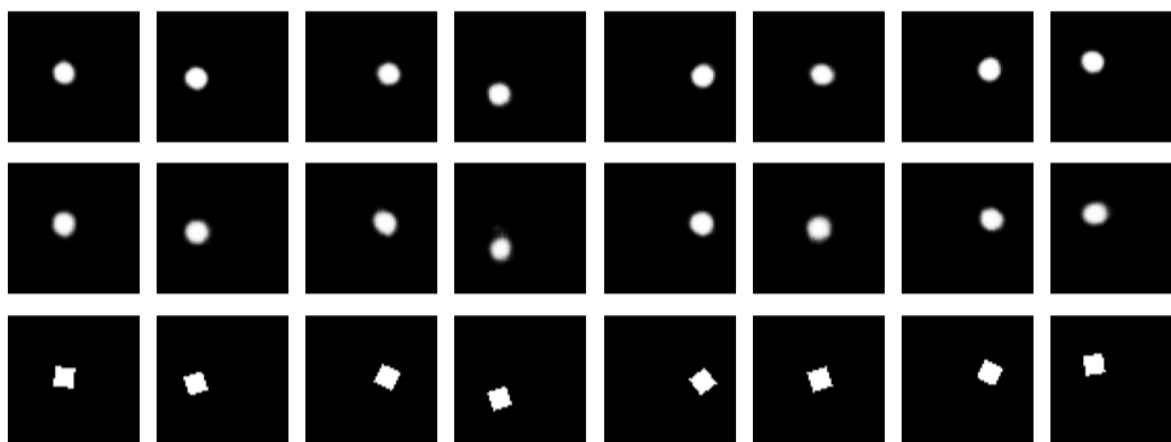


شکل 5 مقایسه $loss$ بازسازی



شکل 6 مقایسه KL فاصله به ازای دو مقدار β

می توان از نمودار های بالا دریافت که بر اساس منحنی های به دست آمده، $\beta=3$ بازسازی بهتری ارائه می دهد ولی فضای نهان کمتر منظم است؛ $\beta=10$ بازسازی ساده تر اما فضای نهان منظم تر و مناسب تر برای disentanglement می سازد. انتخاب مقدار β باید بر اساس هدف نهایی صورت گیرد: برای dSprites اگر نمایش جدایش پذیر معیار اصلی است، $\beta=10$ انتخاب ترجیحی است؛ اگر کیفیت بازسازی اولویت دارد، $\beta=3$ بهتر عمل می کند.



شکل 7 مقایسه اشکال بازی سازی شده با دو مقدار β 3 و 10 و عکس اصلی در پایین

همان‌طور که در شکل مشاهده می‌شود، مدل با $\beta=3$ بازسازی‌های واضح‌تری ارائه می‌دهد، در حالی‌که $\beta=10$ تصاویر نرم‌تر و کم‌جزئیات‌تر تولید کرده است. این رفتار مطابق انتظار VAE است با افزایش β ، مدل فشار بیشتری برای منظم‌سازی فضای نهان اعمال می‌کند، در نتیجه دقت بازسازی کاهش یافته ولی نمایش نهان جدایی‌پذیرتر می‌شود.

به طوری کلی می‌توان اینگونه جمع‌بندی کرد که نتایج به‌دست‌آمده از اجرای مدل‌های β -VAE با مقادیر مختلف β نشان می‌دهد که این پارامتر نقش کلیدی در تنظیم توازن میان دقت بازسازی و منظم‌بودن فضای نهان دارد. با مقدار $\beta=3$ مدل توانست شکل‌ها و موقعیت اجسام را با وضوح بالا و خطای بازسازی کمتر بازتولید کند، در حالی‌که افزایش β به 10 باعث شد تصاویر خروجی نرم‌تر و ساده‌تر شوند اما فضای نهان ساختاریافته‌تر و منظم‌تر گردد. همان‌طور که در نمودارهای آموزشی دیده می‌شود، در β بزرگتر، مقدار Reconstruction Loss افزایش و KL Divergence کاهش یافته است، که بیانگر محدودتر شدن اطلاعات انتقال‌یافته از ورودی به فضای نهان است. این رفتار مطابق با تئوری β -VAE است که با اعمال ضریب بزرگتر بر ترم KL، مدل را وادار می‌کند بازنمایی‌های فشرده‌تر و مستقل‌تری بیاموزد. در نتیجه، مدل با β کوچک‌تر در بازسازی دقیق‌تر عمل می‌کند، ولی مدل با β بزرگتر نمایش نهانی فراهم می‌سازد که برای disentanglement و شناسایی عوامل پنهان در داده‌ها مناسب‌تر است.

زیربخش هفتم

معیار MIG (Mutual Information Gap) یکی از پرکاربردترین سنج‌ها برای ارزیابی میزان جدایش‌پذیری نمایش نهان (disentanglement) در مدل‌های مولد مانند β -VAE است. ایده‌ی اصلی این معیار بر پایه‌ی اندازه‌گیری ارتباط میان هر فاکتور پنهان واقعی داده مثل موقعیت، مقیاس یا چرخش در dSprites و هر بعد از فضای نهان z است.

برای هر فاکتور، ابتدا اطلاعات متقابل بین آن فاکتور و تمام ابعاد z محاسبه می‌شود، سپس دو مقدار بیشینه انتخاب می‌شوند: بیشترین اطلاعات متقابل و دومین بیشترین مقدار. اختلاف بین این دو مقدار بیان می‌کند که آیا یک فاکتور عمدتاً توسط یک بُعد خاص در z توضیح داده می‌شود یا بین چند بُعد پخش شده است. MIG به‌صورت میانگین نرمال‌شده‌ی این اختلاف‌ها بر کل فاکتورها تعریف می‌شود.

اگر MIG بالا باشد، یعنی هر ویژگی پنهان داده با یک بُعد از فضای نهان متناظر است (disentangled)، در حالی‌که MIG پایین یا نزدیک صفر نشان می‌دهد ابعاد نهان درهم و غیرمستقل‌اند. بنابراین، این معیار سنج‌ای کمی از میزان تفکیک‌پذیری عوامل مولد در فضای نمایش مدل فراهم می‌کند.

```

β=1  MIG=0.0108 {'shape': 0.0, 'scale': 0.0, 'orientation': np.float64(0.0013), 'posX': np.float64(0.0195), 'posY': np.float64(0.0331)}
β=3  MIG=0.0198 {'shape': 0.0, 'scale': 0.0, 'orientation': np.float64(0.001), 'posX': np.float64(0.0502), 'posY': np.float64(0.048)}
β=10 MIG=0.0138 {'shape': 0.0, 'scale': 0.0, 'orientation': np.float64(0.0), 'posX': np.float64(0.0492), 'posY': np.float64(0.0196)}

```

شکل 8 محاسبه MIG به ازای β های مختلف

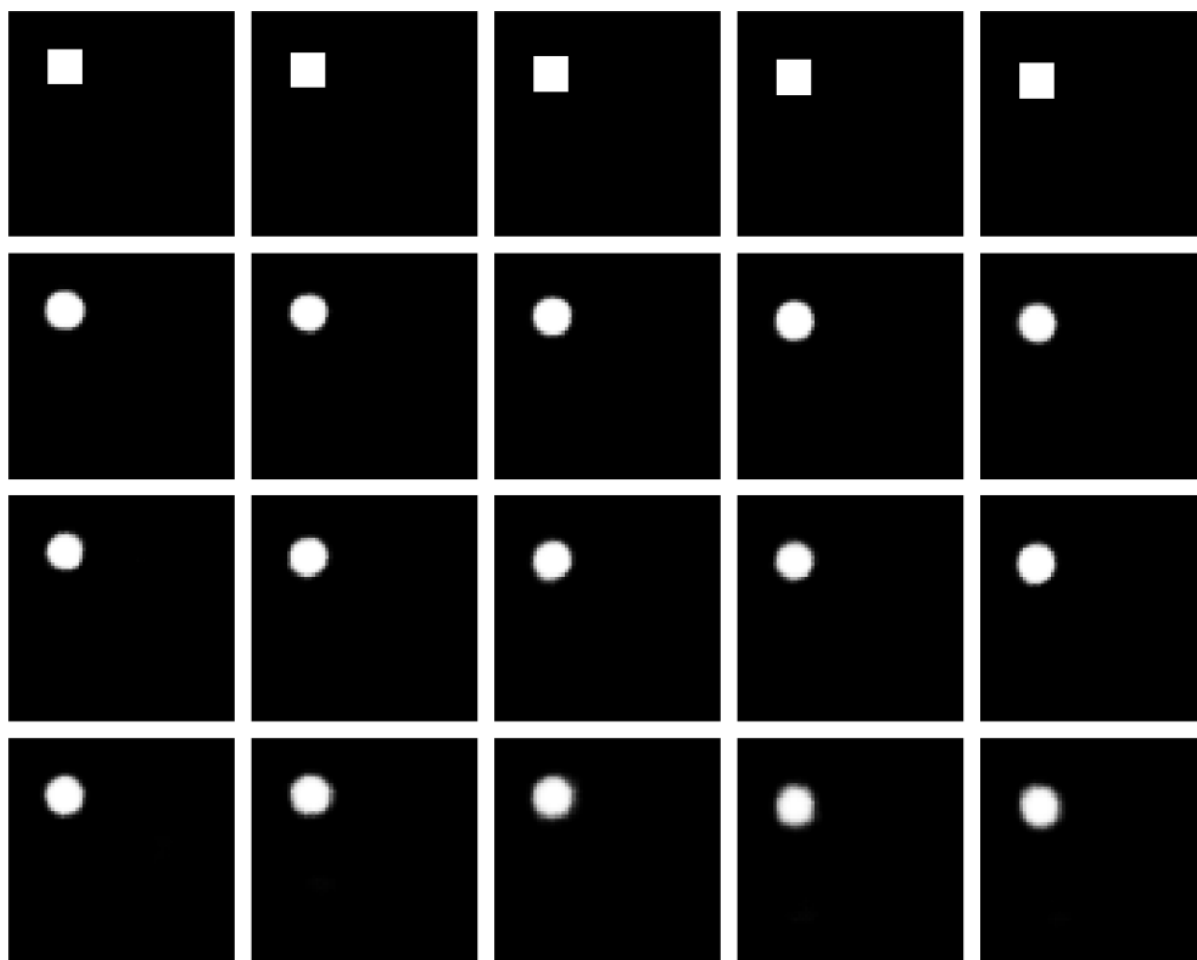
در این بخش، عملکرد مدل های β -VAE با مقادیر مختلف β مورد بررسی قرار گرفت تا میزان جدایش متغیرهای نهان (disentanglement) با استفاده از شاخص Mutual Information Gap (MIG) ارزیابی شود. سه مدل با β های 1، 3 و 10 آموزش داده شدند و نتایج عددی به صورت زیر به دست آمدند:

$\beta=10 \rightarrow \text{MIG}=0.0138$ و $\beta=3 \rightarrow \text{MIG}=0.0198$ ، $\beta=1 \rightarrow \text{MIG}=0.0108$

این نتایج نشان می دهد که با افزایش β از 1 به 3، مقدار MIG افزایش یافته است، که بیانگر بهبود در جدایش متغیرهای پنهان و افزایش منظم سازی فضای نهان است. با این حال، افزایش بیش از حد β تا مقدار 10 باعث کاهش MIG می شود، زیرا مدل بیش از حد محدود شده و بخشی از توانایی بازسازی دقیق تصویر را از دست می دهد.

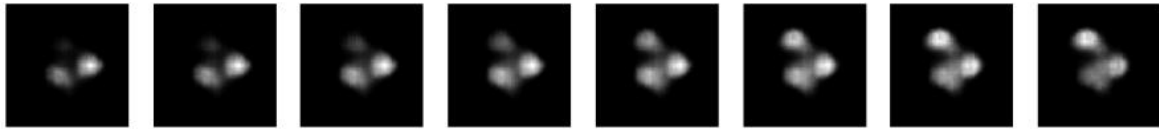
بررسی تفکیکی عوامل (per-factor) نیز نشان داد که بیشترین اطلاعات در ابعاد مربوط به موقعیت مکانی (posX و posY) متمرکز است، در حالی که متغیرهای مربوط به شکل و اندازه سهم ناچیزی دارند. این نتیجه منطقی است زیرا در داده های dSprites، تنوع شکل ها نسبت به تغییرات مکانی کمتر است.

در مجموع، می توان نتیجه گرفت که مدل با $\beta=3$ بهترین تعادل میان بازسازی مناسب و جدایش مؤثر ویژگی های نهان را برقرار کرده است و مقدار MIG در این حالت بیشترین مقدار مثبت را دارد، که با پیش بینی های نظری مقاله β -VAE سازگار است..



شکل 9 مقایسه‌ی بازسازی تصاویر ورودی در مدل‌های β -VAE با مقادیر مختلف β .

Latent dim 0 traversal



شکل 10 تغییر تدریجی در بُعد صفر فضای نهان (Latent traversal) برای مدل β -VAE با $\beta=10$.

در این شکل، مقدار یکی از ابعاد فضای نهان (z_0) از مقدار منفی به مثبت تغییر کرده است. مشاهده می‌شود که با تغییر این بُعد، تنها یک ویژگی از تصویر (مانند موقعیت یا شکل جسم) تغییر می‌کند در حالی که سایر ویژگی‌ها ثابت می‌مانند. این رفتار نشان‌دهنده جدایش‌پذیری (disentanglement) در نمایش نهانی مدل β -VAE است.

برای بررسی کیفی عملکرد مدل‌ها، از دو روش استفاده شد: بازسازی تصاویر ورودی و مشاهده traversal فضای نهان. در بازسازی‌ها مشاهده شد که برای $\beta=1$ ، خروجی‌ها از لحاظ بصری واضح‌تر و نزدیک‌تر به داده‌های واقعی هستند اما مدل درک تفکیک‌پذیری از فاکتورهای مختلف را ندارد. در $\beta=3$ ، بازسازی اندکی ساده‌تر شده ولی مدل ساختار منظم‌تری از ویژگی‌ها در فضای پنهان یافته است. در $\beta=10$ ، بازسازی‌ها نرم‌تر و انتزاعی‌تر هستند و برخی جزئیات از بین رفته‌اند، اما هر بُعد از فضای نهان مسئول یک ویژگی خاص تصویر (مثل موقعیت یا چرخش) شده است.

نتایج traversal نیز این الگو را تأیید می‌کنند: با تغییر تدریجی مقدار یک بُعد از z در مدل $\beta=10$ ، تنها یک ویژگی مشخص در خروجی تغییر می‌کند در حالی که سایر ویژگی‌ها ثابت می‌مانند. این رفتار نشانه‌ی واضحی از یادگیری بازنمایی disentangled است. در مجموع، هرچند افزایش β باعث افت جزئی در کیفیت بازسازی می‌شود، ولی منجر به نمایش‌های پنهان منظم‌تر و قابل‌تجزیه‌تری می‌شود که ویژگی‌های مستقل داده را به‌طور جداگانه در فضای نهان مدل می‌نمایانند.

بخش دوم

زیربخش اول

مدل (VQ-VAE (Vector Quantized VAE یکی از نسخه‌های پیشرفته VAE است که فضای نهان را از حالت پیوسته به Discrete latent space تبدیل می‌کند. در این روش، خروجی Encoder به جای یک بردار پیوسته، به نزدیک‌ترین بردار از میان مجموعه‌ای از codebook vectors نگاشت می‌شود. به این ترتیب، مدل یاد می‌گیرد که نمایش داده‌ها را در قالب کدهای گسسته ذخیره کند. در فرآیند آموزش، چون عمل quantization مشتق‌پذیر نیست، از تکنیک Straight-Through Estimator برای عبور گرادیان‌ها استفاده می‌شود. مزیت اصلی VQ-VAE این است که از posterior collapse جلوگیری می‌کند و باعث می‌شود نمایش‌های نهفته معنادارتر و فشرده‌تر یاد گرفته شوند. این ویژگی باعث شده VQ-VAE در داده‌هایی مثل گفتار و ویدئو عملکرد بهتری داشته باشد، زیرا ساختارهای سطح بالا را بهتر حفظ می‌کند. در عوض، به دلیل گسسته بودن فضای نهان، فرآیند آموزش پیچیده‌تر و حساس‌تر به انتخاب اندازه‌ی codebook است.

زیربخش دوم

مدل VampPrior توسط Tomczak و Welling برای بهبود توزیع پیشین در VAE معرفی شد. در VAE معمولی، پرایر معمولاً یک توزیع نرمال استاندارد است $p(z) = \mathcal{N}(0, I)$ که ممکن است

بیش از حد محدودکننده باشد و منجر به غیرفعال شدن برخی ابعاد نهان (Inactive latent units) شود. در VampPrior، پرایر به صورت ترکیب چند توزیع پسین تعریف می شود که در آن u_k ها ورودی های مصنوعی (pseudo-inputs) قابل یادگیری اند.

$$p_{\lambda}(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^K q_{\phi}(\mathbf{z} | \mathbf{u}_k),$$

شکل 11 فرمول prior در مدل vampPrior

فرمول بالا تعریف اصلی پرایر در مدل VampPrior است و نشان می دهد که برخلاف VAE معمولی که از پرایر ثابت و ساده ای مثل نرمال استاندارد $\mathcal{N}(0, I)$ استفاده می کند، در VampPrior پرایر به صورت میانگینی از چند توزیع پسین تعریف می شود. در این مدل، مجموعه ای از pseudo-inputs با نماد u_k در نظر گرفته می شود که مدل در طول آموزش آن ها را یاد می گیرد. سپس برای هر u_k ، انکو در توزیع $q_{\phi}(\mathbf{z} | u_k)$ را تولید می کند و پرایر نهایی برابر با میانگین این توزیع ها است.

این طراحی باعث می شود پرایر بتواند چندحالتی (multi-modal) و منطبق تر با داده های واقعی باشد. در نتیجه، مدل ظرفیت بیشتری برای بازنمایی ساختارهای پیچیده پیدا می کند، ابعاد نهان بیشتری فعال می شوند و بازسازی ها کیفیت بالاتری دارند. البته در مقابل، مدل VampPrior از نظر محاسباتی سنگین تر است و نیاز به حافظه بیشتری برای ذخیره ی pseudo-inputs دارد.

زیربخش سوم

مدل SC-VAE ترکیبی است از اصول sparse coding و VAE؛ یعنی فرض می کند که کد نهانی می تواند به صورت ترکیبی خطی از تعداد کمی atoms زده شود، و به کمک نسخه یادگرفتنی از ISTA (Iterative Shrinkage-Thresholding Algorithm) این ترکیب را حل می کند. نقش ISTA در این مدل آن است که فرآیندی تکراری برای یافتن ترکیبی از اتم ها با مقدار آستانه گذاری دارد و با یادگیری پارامترهای آن، کد نهانی منظم تر، فشرده تر و قابل تفسیرتر می شود. مزیت اصلی SC-VAE نسبت به مدل پایه VAE این است که بازسازی دقیق تر و با وضوح بالاتر ارائه می دهد و نمایشی تنک

از داده ایجاد می‌کند که مفاهیم زیرین را بهتر جدا می‌کند. از سوی دیگر، معایب ممکن است پیچیدگی محاسباتی بیشتر، نیاز به تنظیم دقیق آستانه‌ها و اتصالات الگوریتمی بیشتر باشد.

1. Chen, R. T. Q., Li, X., Grosse, R. B., & Duvenaud, D. (2018). Isolating sources of disentanglement in variational autoencoders. In Advances in Neural Information Processing Systems (NeurIPS 2018).
2. van den Oord, A., Vinyals, O., & Kavukcuoglu, K. (2017). Neural discrete representation learning. In Advances in Neural Information Processing Systems (NeurIPS 2017).
3. Tomczak, J. M., & Welling, M. (2018). VAE with a VampPrior. In Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics (AISTATS 2018) (Vol. 84, pp. 1214-1223). PMLR.
4. Xiao, P., Qiu, P., Ha, S., Bani, A., Zhou, S., & Sotiras, A. (2024). Sparse Coding based Variational Autoencoder with Learned ISTA (SC-VAE). Pattern Recognition, 161, Article 111187.