



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)
دانشکده مهندسی کامپیوتر

سمینار کارشناسی ارشد
گرایش هوش مصنوعی و رباتیکز

آشنایی با مدل‌های پخشی

نگارش

عطیه غفارلوی مقدم

استاد درس

دکتر رضا صفابخش

۱۴۰۲/۶/۱۰

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

چکیده

مدل‌های پخشی به عنوان یک خانواده‌ی قدرتمند از مدل‌های مولد با نتایج قابل توجه در حوزه‌های مختلف تصویر، صوت، متن و گراف، در سال‌های اخیر مورد توجه بسیاری از محققین یادگیری ماشین قرار گرفته‌اند. همه‌ی مدل‌های مطرح شده در این خانواده از یک رویکرد مشابه پیروی می‌کنند و آن هم تخریب داده با نویز در یک فرآیند گام به گام و سپس بازگرداندن این فرآیند برای تولید نمونه‌های جدید از داده می‌باشد. در این گزارش ابتدا به بررسی دقیق سه دسته‌ی مطرح مدل‌های پخشی و شباهت‌ها و تفاوت‌های آنها می‌پردازیم. پس از معرفی این مدل‌ها، به بررسی رابطه‌ی میان آنها و انواعی دیگر از مدل‌های مولد می‌پردازیم و در نهایت به تعدادی از چالش‌های مطرح در خصوص این مدل‌ها اشاره می‌کنیم.

واژه‌های کلیدی:

مدل‌های مولد، مدل‌های پخشی، مدل‌های مولد مبتنی بر امتیاز، مدل‌های پخشی نویززدای احتمالاتی

فهرست مطالب

صفحه

عنوان

۱	مقدمه	۱
۴	مدل‌های مولد مبتنی بر امتیاز	۲
۵	۱-۲ مفاهیم اولیه	
۵	۱-۱-۲ تابع امتیاز و فرآیند تطبیق امتیاز	
۶	۲-۱-۲ تطبیق امتیاز نويززدا	
۷	۳-۱-۲ دینامیک لانژوین	
۷	۲-۲ تولید نمونه با روش تطبیق امتیاز و دینامیک لانژوین	
۷	۱-۲-۲ فرضیه‌ی روبه‌ها	
۷	۲-۲-۲ نواحی کمتر چگال	
۸	۳-۲-۲ راه‌حل پیشنهادی	
۹	۴-۲-۲ شبکه‌ی امتیاز شرطی شده با نويز	
۱۰	۵-۲-۲ به کارگیری تطبیق امتیاز برای آموزش شبکه‌ی امتیاز شرطی شده با نويز	
۱۰	۶-۲-۲ نمونه‌گیری از شبکه‌ی امتیاز شرطی شده با نويز	
۱۱	۳-۲ بهبود مدل‌های مبتنی بر امتیاز	
۱۲	۱-۳-۲ انتخاب دنباله‌ی نويز	
۱۳	۳ مدل‌های پخشی نويززدای احتمالاتی	
۱۴	۱-۳ تعریف مدل	
۱۴	۱-۱-۳ مسیر رو به جلو	
۱۵	۲-۱-۳ مسیر رو به عقب	
۱۶	۳-۱-۳ تابع مدل مولد	
۱۶	۲-۳ تعریف تابع خطا	
۱۷	۳-۳ تغییر پارامتر تابع خطا	
۱۸	۴-۳ فرآیند یادگیری	
۱۹	۵-۳ بهبود مدل‌های پخشی نويززدای احتمالاتی	
۲۰	۱-۵-۳ یادگیری واریانس	
۲۰	۲-۵-۳ بهبود گام‌های نويز	
۲۲	۴ مدل‌های مولد مبتنی بر معادلات دیفرانسیل تصادفی	
۲۳	۱-۴ معرفی فرآیند پخشی با معادلات دیفرانسیل تصادفی	
۲۴	۱-۱-۴ مسیر روبه جلو	
۲۴	۲-۱-۴ مسیر معکوس	

۲۵ تخمین تابع امتیاز ۳-۱-۴
۲۵ مزیت‌های مدل‌های مولد مبتنی بر معادله دیفرانسیل تصادفی ۲-۴
۲۶ ارتباط روش مبتنی بر معادلات دیفرانسیل تصادفی با روش‌های پیش‌تر معرفی شده ۳-۴
۲۷ انواع روش‌های نمونه برداری ۴-۴
۲۷ حل کننده‌های عام منظوره ی معادلات دیفرانسیل تصادفی ۱-۴-۴
۲۷ تبدیل به معادلات دیفرانسیل ساده ۲-۴-۴
۲۸ مقایسه‌ی روش‌های مولد ۵
۲۹ بررسی ارتباط مدل‌های پخشی با سایر مدل‌های مولد ۱-۵
۲۹ ارتباط خودکدگذارهای تغییراتی با مدل‌های پخشی ۱-۱-۵
۳۰ ارتباط شبکه‌های مولد تقابلی با مدل‌های پخشی ۲-۱-۵
۳۱ مشکل سه‌گانه‌ی مدل‌های مولد ۳-۱-۵
۳۲ مقایسه‌ی انواع روش‌های پخشی ۲-۵
۳۳ نتایج و عملکرد ۱-۲-۵
۳۳ معماری ۲-۲-۵
۳۴ ویژگی‌های کلی ۳-۲-۵
۳۵ جمع‌بندی و کارهای آتی ۶
۳۶ جمع‌بندی ۱-۶
۳۷ کارهای آتی ۲-۶
۳۸ کتاب‌نامه
۴۱ واژه‌نامه‌ی فارسی به انگلیسی
۴۴ واژه‌نامه‌ی انگلیسی به فارسی

شکل	فهرست تصاویر	صفحه
۱-۱	روند کلی یک مدل پخشی	۳
۱-۲	میدان برداری تابع امتیاز برای یک توزیع یک بعدی	۶
۲-۲	مشکل نواحی کمتر چگال در تخمین تابع امتیاز	۸
۳-۲	مشکل نواحی کمتر چگال در دینامیک لانژوین	۹
۴-۲	حل مشکل نواحی کمتر چگال در دینامیک لانژوین تابکاری شده	۱۱
۱-۳	فرآیند پخش در مپناها	۱۴
۲-۳	معماری $pixelCNN++$	۱۹
۳-۳	تفاوت برنامه‌ی نویز خطی و کسینوسی در فرآیند پخشی	۲۰
۴-۳	تفاوت برنامه‌ی نویز خطی و کسینوسی در حذف بخشی از فرآیند نویززدایی	۲۱
۱-۴	مدل‌های پخشی مبتنی بر معادلات دیفرانسیل	۲۳
۱-۵	مدل‌های پخشی مبتنی بر معادلات دیفرانسیل	۳۰
۲-۵	مشکل سه‌گانه‌ی مدل‌های مولد	۳۲

صفحه	فهرست جداول	جدول
۳۲	مقایسه‌ی مدل‌های مولد با یکدیگر	۱-۵
۳۳	مقایسه‌ی کیفیت تصاویر تولیدی با سه مدل پخشی روی مجموعه داده‌ی $CIFAR-10$	۲-۵

فهرست اختصارات

عنوان اختصاری عنوان کامل

مپنا مدل پخشی نويززدای احتمالاتی

فصل اول

مقدمه

مدل‌های مولد کاربردهای متنوعی در یادگیری ماشین دارند، از جمله‌ی این کاربردها می‌توان به تولید عکس‌های باکیفیت، تولید قطعه‌های موسیقی و یا صوت، افزایش کارایی یادگیری نیمه نظارتی^۱، تشخیص ناهنجاری^۲، یادگیری تقلیدی و کمک به یادگیری تقویتی اشاره کرد. پیشرفت‌های اخیر این حوزه دو رویکرد کلی را دنبال می‌کنند: رویکردهای مبتنی بر درست‌نمایی^۳ و شبکه‌های مولد تقابلی^۴. در رویکردهای مبتنی بر درست‌نمایی از لگاریتم درست‌نمایی یا خطاهای جایگزین به عنوان تابع هدف یادگیری استفاده می‌شود؛ در حالیکه در شبکه‌های مولد تقابلی تلاش می‌شود تا یادگیری با یک رویکرد تقابلی بین دو شبکه‌ی تمایزگر^۵ و تولیدکننده^۶ صورت گیرد.

با وجود عملکرد قابل توجه مدل‌های مبتنی بر درست‌نمایی و شبکه‌های مولد تقابلی، این مدل‌ها به واسطه‌ی نوع تعریفشان درگیر برخی محدودیت‌های ذاتی هستند. به عنوان مثال مدل‌های مبتنی بر درست‌نمایی برای ساخت توزیع‌های نرمال شده نیازمند به کارگیری معماری‌های خاص یا استفاده از توابع خطای جایگزین (مانند حد پایین مشاهده^۷ در خودکدگذارهای تغییراتی^۸) در فرآیند یادگیری هستند. شبکه‌های مولد تقابلی مشکلات موجود در مدل‌های مبتنی بر درست‌نمایی را ندارند اما فرآیند یادگیری آنها به خاطر ماهیت تقابلی آن می‌تواند دچار ناپایداری شود. همچنین تابع هدف مدل‌های مولد تقابلی برای مقایسه و ارزیابی انواع این مدل‌ها چندان مناسب نیست [۱۴].

مدل‌های پخشی^۹ یک دسته‌ی نوظهور از مدل‌های مولد هستند که با توجه به عملکرد مناسب آنها در حوزه‌های تصویر، صوت، متن و گراف به سرعت مورد استقبال قرار گرفتند. وجه مشترک همه‌ی مدل‌های پخشی این است که در یک فرآیند رو به جلو^{۱۰} داده را با نویز تخریب می‌کنند و در نتیجه از توزیع داده به یک توزیع دلخواه قابل محاسبه می‌رسند و سپس در یک مسیر رو به عقب که توسط مدل یاد گرفته می‌شود، تلاش می‌کنند تا داده را بازسازی و از یک نمونه از توزیع دلخواه به یک نمونه از توزیع داده‌ها برسند. این موضوع در شکل ۱-۱ نشان داده شده است.

با وجود اینکه چهارچوب گفته شده در بند قبل در همه‌ی مدل‌های پخشی وجود دارد؛ رویکردهای مختلفی برای فرموله سازی و پیاده سازی این روش وجود دارد. در مدل‌های مولد مبتنی بر تابع امتیاز^{۱۱} تلاش می‌شود با استفاده از گرادیان لگاریتم توزیع داده، که به آن تابع امتیاز^{۱۲} می‌گوییم، و یک فرآیند

¹Semi-supervised learning

²Anomaly

³Likelihood

⁴Generative adversarial networks

⁵Discriminator

⁶Generator

⁷Evidence lower bound

⁸Variational autoencoders

⁹Diffusion models

¹⁰Forward process

¹¹Score-based generative models

¹²Score function



شکل ۱-۱: روند کلی یک مدل پخششی

در هر مدل پخششی یک فرآیند ثابت گام به گام رو به جلو برای تخریب داده و یک فرآیند معکوس برای تولید یک نمونه داده از نویز وجود دارد [۹].

گام به گام حرکت در جهت این گرادیان نمونه‌ی جدید تولید کنیم. از طرفی در مدل‌های پخششی نویززدای احتمالاتی فرآیند تخریب داده و تولید نمونه با یک زنجیره‌ی مارکوفی احتمالاتی مدل می‌شود و در نهایت در مدل‌های پخششی مبتنی بر معادلات دیفرانسیل تصادفی، فرآیند تخریب به صورت یک معادله دیفرانسیل تصادفی پیوسته مدل می‌شود که با پیدا کردن معکوس این معادله و گسسته سازی آن امکان تولید نمونه‌ی جدید از توزیع داده‌ها فراهم می‌گردد.

می‌توان نشان داد که مدل‌های پخششی بسیاری از محدودیت‌های مدل‌های مولد دیگر مانند نیاز به تمایزگر در شبکه‌های مولد تقابلی، نیاز به مقیدسازی شبکه در جریان‌های نرمال‌ساز^{۱۳} و نیاز به تنظیم توزیع‌های پسین در خودکدگذارهای تغییراتی را ندارند و قابلیت تولید تصاویر باکیفیت و متنوع را دارند. مشکل اصلی در مدل‌های پخششی فرآیند نمونه‌برداری طولانی آنهاست که باعث عدم کارایی این مدل‌ها در کاربردهای تعاملی می‌شود. با رفع این موضوع می‌توان اطمینان حاصل کرد که این دسته از مدل‌های مولد از توانایی همه‌جانبه‌ای برای به کارگیری در حوزه‌های مختلف برخوردار هستند.

در این گزارش ابتدا هر یک از سه مدل پخششی و فرموله سازی هر کدام را بیان می‌کنیم و نشان می‌دهیم که تحت شرایطی خاص این روش‌ها هم ارز و قابل تبدیل به یکدیگر هستند. سپس در فصل پنجم به مقایسه‌ی کلی روش‌های پخششی با سایر مدل‌های مولد و همچنین مقایسه‌ی انواع مدل‌های پخششی مطالعه شده در این گزارش با یکدیگر می‌پردازیم و در نهایت به برخی جنبه‌هایی از مدل‌های پخششی می‌پردازیم که قابلیت مطالعه‌ی بیشتر و اعمال نوآوری را دارند.

¹³Normalizing flows

فصل دوم

مدل‌های مولد مبتنی بر امتیاز

در این فصل به معرفی اولین دسته از مدل‌های پخش می‌پردازیم. این مدل‌ها برای تولید نمونه از گرادیان توزیع نسبت به داده، که تابع امتیاز نام دارد، کمک می‌گیرند. اما از آنجایی که تابع امتیاز در اکثر داده‌های دنیای واقعی خوش تعریف نیست، برای تولید نمونه‌ها نیازمند سازوکاری هستیم که بتواند ما را به تابع امتیاز مناسب‌تری برساند. در این فصل ابتدا مفاهیم اولیه‌ی به کارگیری تابع امتیاز برای تولید نمونه‌ی جدید را بررسی می‌کنیم و سپس مشکلات این روش و راه‌حل‌های پیشنهادی ارائه شده را مطرح می‌کنیم و در نهایت تکنیک‌هایی برای افزایش مقیاس پذیری این مدل معرفی می‌کنیم.

۱-۲ مفاهیم اولیه

۱-۱-۲ تابع امتیاز و فرآیند تطبیق امتیاز

فرض کنید مجموعه داده‌ی ما از تعدادی نمونه‌ی مستقل و با توزیع یکسان از توزیع ناشناخته‌ی $p_{data}(x)$ تشکیل شده است. امتیاز یک توزیع احتمالاتی $p(x)$ را گرادیان لگاریتم آن توزیع نسبت به داده‌ها تعریف می‌کنیم که به فرم $\nabla_x \log p(x)$ نشان داده می‌شود. امتیاز یک توزیع احتمالاتی در واقع یک میدان برداری^۱ است که به سمتی اشاره می‌کند که درست‌نمایی داده‌ها بیشترین رشد را دارد. شکل ۱-۲ این میدان برداری را برای یک توزیع یک بعدی نشان می‌دهد. با داشتن تابع امتیاز یک توزیع می‌توانیم از آن توزیع نمونه‌گیری کنیم و خود توزیع را تخمین بزنیم.

شبکه‌ی امتیاز $s_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^D$ یک شبکه‌ی عصبی با پارامترهای θ می‌باشد که برای تخمین تابع امتیاز یک توزیع مانند $p_{data}(x)$ آموزش داده می‌شود. در صورتی که بتوانیم تابع امتیاز توزیع را بدست آوریم می‌توانیم از این تابع برای مدل مولد و همچنین تولید نمونه‌های جدید استفاده کنیم. فرآیند تطبیق امتیاز^۲ [۷] فرآیندی است که برای آموزش شبکه‌ی امتیاز به کار گرفته می‌شود و هدف اصلی آن تخمین $\nabla_x \log p(x)$ بدون نیاز به تخمین مستقیم $p(x)$ است. تابع هدف تخمین امتیاز به صورت زیر قابل تعریف است:

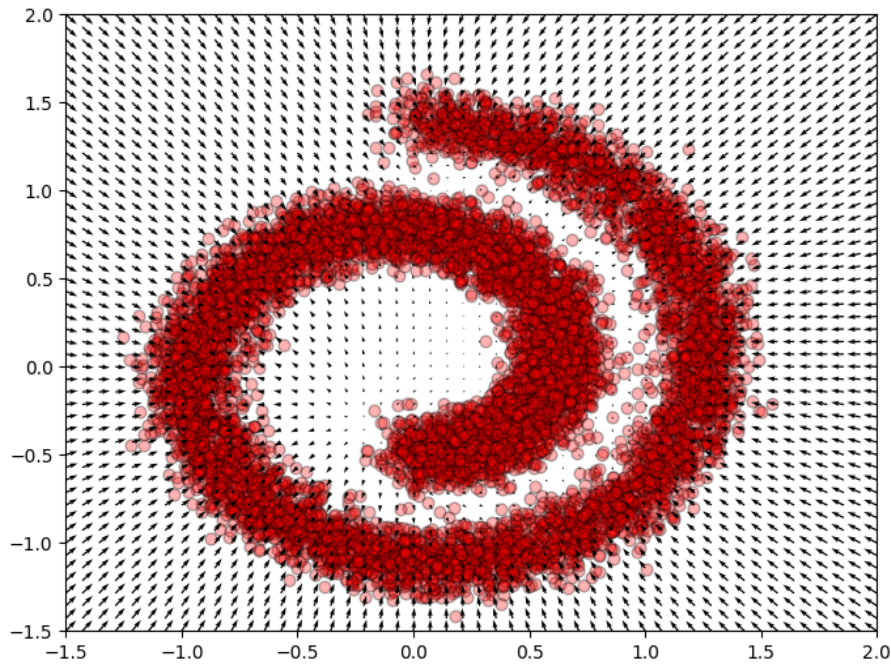
$$\frac{1}{2} \mathbb{E}_{p_{data}} [\|s_\theta(x) - \nabla_x \log p_{data}(x)\|_2^2] \quad (۱-۲)$$

که می‌توان به صورت ریاضی نشان داد که با رابطه‌ی زیر هم‌ارز است:

$$\mathbb{E}_{p_{data}(x)} [tr(\nabla_x s_\theta(x)) + \frac{1}{2} \|s_\theta(x)\|_2^2] \quad (۲-۲)$$

^۱ Vector field

^۲ Score matching



شکل ۲-۱: میدان برداری تابع امتیاز برای یک توزیع یک بعدی

معادله‌ی بالا به دلیل محاسبه‌ی ژاکوبین تابع امتیاز قابلیت مقیاس‌پذیری و پیاده‌سازی در داده‌ی با بعد بالا را ندارد بنابراین باید به دنبال جایگزینی برای تابع هدف باشیم به گونه‌ای که بتوانیم در مقیاس بزرگتر و ابعاد بالاتر نیز از رویکرد تطبیق امتیاز برای تخمین تابع امتیاز استفاده کنیم.

۲-۱-۲ تطبیق امتیاز نویززا

تطبیق امتیاز نویززا^۳ [۱۷] یک رویکرد تطبیق امتیاز است که به هدف رفع مشکل محاسبه‌ی ژاکوبین توسعه داده شد. در این فرآیند ابتدا نقطه‌ی x با توزیع نویز از پیش تعریف شده‌ی $q_\sigma(\tilde{x}|x)$ تخریب می‌شود و سپس تخمین امتیاز برای توزیع داده‌ی تخریب شده‌ی $q_\sigma(\tilde{x}) \triangleq \int_{q_\sigma}(\tilde{x}|x)p_{data}(x)dx$ صورت می‌گیرد. در این حالت تابع هدف به صورت زیر می‌باشد:

$$\frac{1}{2} \mathbb{E}_{q_\sigma(\tilde{x}|x)p_{data}(x)} [\|s_\theta(\tilde{x}) - \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x)\|_2^2] \quad (۳-۲)$$

که با کمینه کردن آن تابع بهینه‌ی $s_{\theta^*}(x)$ بدست می‌آید. حال می‌توان نشان داد که اگر نویز به اندازه‌ی کافی کوچک باشد، آنگاه $s_{\theta^*}(x)$ تخمینی از تابع امتیاز توزیع داده است.

³Denoising score matching

۳-۱-۲ دینامیک لانژوین

دینامیک لانژوین^۴ برای تولید نمونه از توزیع ناشناخته ی $p(x)$ با استفاده از تابع امتیاز این توزیع به کار می‌رود. با داشتن یک اندازه ی گام $\epsilon > 0$ و یک مقدار اولیه ی $\tilde{x}_0 \sim \pi(x)$ که π یک توزیع اولیه است. دینامیک لانژوین به صورت تکراری فرآیند زیر را تکرار می‌کند:

$$\tilde{x}_t = \tilde{x}_{t-1} + \frac{\epsilon}{2} \nabla_x \log p(\tilde{x}_{t-1}) + \sqrt{\epsilon} z_t \quad (۴-۲)$$

که $z_t \sim \mathcal{N}(0, I)$ است. توزیع \tilde{x}_T با توزیع $p(x)$ برابر می‌شود اگر $\epsilon \rightarrow 0$ و $T \rightarrow \infty$ باشد. با توجه به معادله ی (۴-۲) برای بدست آوردن نمونه ی جدید تنها به تابع امتیاز توزیع نیاز داریم از این رو می‌توانیم ابتدا یک شبکه ی امتیاز برای تخمین امتیاز در هر نقطه به دست آوریم و سپس با دینامیک لانژوین و با استفاده از تخمین‌های امتیازات به یک نمونه ی جدید از توزیع برسیم.

۲-۲ تولید نمونه با روش تطبیق امتیاز و دینامیک لانژوین

همان طور که پیش‌تر هم اشاره شد، با داشتن تخمین تابع امتیاز به کمک تطبیق امتیاز و دینامیک لانژوین می‌توان به یک نمونه ی جدید رسید. اما در ادامه نشان خواهیم داد که دوماً برای پیاده‌سازی این فرم از تولید نمونه‌ها وجود دارد که ما را ملزم می‌سازد تا روش پیشنهادی را بهبود دهیم.

۱-۲-۲ فرضیه ی رویه‌ها

فرضیه ی رویه‌ها بیان می‌کند که اکثر داده‌های دنیای واقعی در یک فضای با بعد بالا روی یک رویه ی با بعد کمتر قرار دارند. به صورت تجربی نشان داده شده است که این موضوع برای بسیاری از مجموعه داده‌ها برقرار است. در این صورت روش‌های مولد مبتنی بر امتیاز با دو مشکل اساسی روبرو هستند. مشکل اول این است که تابع امتیاز گرادین گرفته شده در فضای با بعد بالاتر است و زمانی که x محدود به یک رویه ی با بعد پایین‌تر است این گرادین تعریف نشده است. و دوم، تابع هدف فرآیند تطبیق امتیاز تنها در صورتی تخمین سازگاری از امتیاز می‌دهد که فضای پشتیبان^۵ توزیع داده‌ها با کل فضای تعریف شده یکی باشد که در حالتی که داده‌ها روی رویه ی با بعد کمتر هستند چنین نیست.

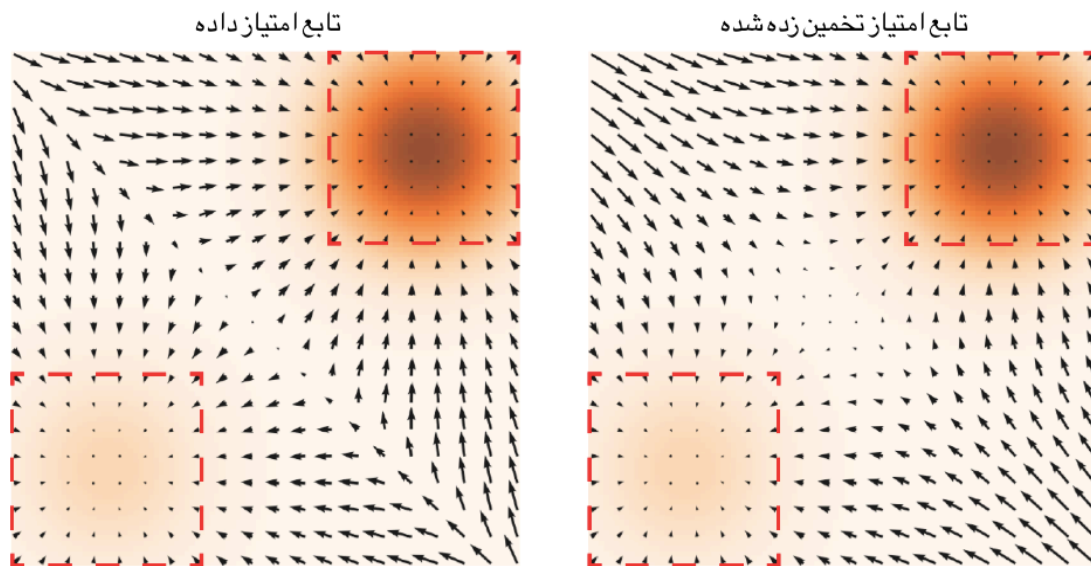
۲-۲-۲ نواحی کمتر چگال

پراکندگی داده در نواحی کمتر چگال هم برای تخمین تابع امتیاز و هم برای دینامیک لانژوین مشکل ایجاد می‌کند. در تخمین تابع امتیاز در نواحی که نمونه‌های کافی از داده‌ها وجود ندارد، به دلیل کم

^۴Langevin dynamics

^۵Support

بودن شواهد توزیع داده، خطای زیادی وجود دارد. برای توضیح این موضوع از یک مثال توزیع گاوسی مخلوط استفاده می‌کنیم. همانطور که در شکل ۲-۲ مشاهده می‌شود، در نواحی بین دو قله ی گاوسی تخمین تابع امتیاز تخمین قابل اتکایی نیست.



شکل ۲-۲: مشکل نواحی کمتر چگال در تخمین تابع امتیاز

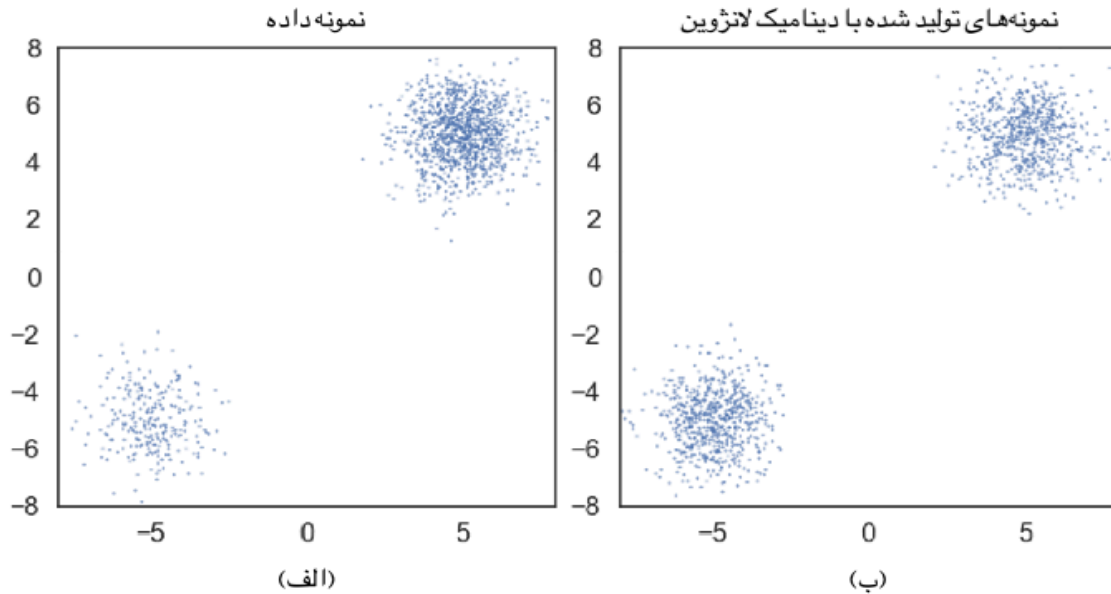
شکل سمت چپ جهت‌های میدان برداری امتیاز را در یک توزیع ترکیب گاوسی نشان می‌دهد و شکل سمت راست تخمین به دست آمده از تابع امتیاز را نشان می‌دهد. همانطور که مشاهده می‌کنید تخمین امتیاز در بین دو قله ی گاوسی که ناحیه ی کمتر چگال است با خطا همراه است [۱۴].

همچنین زمانی که دو قله ی توزیع یا دو ناحیه ی چگال داده با یک ناحیه ی کمتر چگال از یکدیگر جدا شده باشند، دینامیک لانژوین قابلیت درک درست وزن نسبی قله‌ها نسبت به یکدیگر را از دست می‌دهد و در واقع قدرت تمایز میان نواحی چگال تر را از دست می‌دهد و در نتیجه به توزیع درست همگرا نمی‌شود. همچنین در صورتی که دو مد داده فضای پشتیبان مجزایی داشته باشند نیز دینامیک لانژوین برای تولید نمونه‌های درست نیازمند گام‌های کوچکتر و در نتیجه زمان ترکیب^۶ زیادی می‌باشد. شکل ۳-۲ این موضوع را نشان می‌دهد.

۳-۲-۲ راه‌حل پیشنهادی

می‌توان مشاهده کرد که تخریب داده با نویز گاوسی تصادفی می‌تواند مجموعه داده ی مناسب تری برای اعمال روش‌های مبتنی بر امتیاز بسازد. زیرا اولاً اضافه کردن نویز گاوسی محدودیت داده به یک رویه ی

^۶Mixing



شکل ۲-۳: مشکل نواحی کمتر چگال در دینامیک لانژوین

شکل (الف) نمونه‌های یک توزیع مخلوط گاوسی را نشان می‌دهد که از دو قله با وزن‌های متفاوت ایجاد شده است. شکل (ب) نمونه‌های تولید شده با دینامیک لانژوین از همین توزیع را نشان می‌دهد. همانطور که مشاهده می‌شود، دینامیک لانژوین از هر دو قله به یک نسبت نمونه تولید کرده است و قادر به درک وزن نسبی قله‌ها نسبت به یکدیگر نیست [۱۴].

با بعد پایین‌تر را از بین می‌برد چراکه فضای پشتیبان یک توزیع گاوسی کل فضای با بعد بالاست و همین امر سبب می‌شود که مشکلات ایجاد شده به واسطه‌ی فرضیه‌ی رویه‌ها برطرف شود. و دوم اینکه نویز گاوسی بزرگ امکان پرکردن ناحیه‌های کمتر چگال در توزیع داده‌ی تخریب نشده را دارد و همین امر می‌تواند مشکل تخمین امتیاز در نواحی کمتر چگال و مشکلات دینامیک لانژوین را برطرف کند. همچنین اضافه کردن نویز در چندین مرحله تا همگرا شدن به توزیع اصلی داده و بعد اعمال دینامیک لانژوین روی هر یک از این توزیع‌های میانی می‌تواند سرعت همگرایی دینامیک لانژوین را افزایش دهد. با توجه به فرضیه‌ی مطرح شده در بند قبل می‌توانیم فرآیند تولید نمونه‌ی مبتنی بر امتیاز را به دو صورت بهبود دهیم: اول، داده را با سطوح مختلف نویز تخریب کنیم و دوم، از یک شبکه‌ی شرطی شده بر روی نویز برای تخمین امتیاز در همه‌ی سطوح نویز استفاده کنیم. سپس از دینامیک لانژوین برای تولید نمونه‌ها از سطوح نویز بزرگ به سطوح کوچک استفاده کنیم. در ادامه جزئیات این روش را بیشتر توضیح خواهیم داد.

۲-۲-۴ شبکه‌ی امتیاز شرطی شده با نویز

مجموعه‌ی $\{\sigma_i\}_{i=1}^L$ را یک دنباله‌ی هندسی مثبت در نظر بگیرید که $\frac{\sigma_1}{\sigma_2} = \dots = \frac{\sigma_{L-1}}{\sigma_L} > 1$ و توزیع داده‌ی تخریب شده را نیز به صورت $q_\sigma(x) \triangleq \int p_{data}(t) \mathcal{N}(x|t, \sigma^2 I) dt$ تعریف می‌کنیم. مجموعه‌ی $\{\sigma_i\}$ را به گونه‌ای تعریف می‌کنیم که σ_1 به اندازه‌ی کافی بزرگ باشد تا بتوانیم مشکلات مطرح

شده در بخش قبل را برطرف کنیم و σ_L را به اندازه‌ای کوچک تعریف می‌کنیم که توزیع تخریب شده تقریباً با توزیع داده‌ها برابر باشد. هدف آموزش شبکه‌ی امتیاز شرطی شده با نویز این است که تابع امتیاز تمامی توزیع‌های تخریب شده را به گونه‌ای به دست آورد که داشته باشیم: $\forall \sigma \in \{\sigma_i\}_{i=1}^L$: $s_\theta(x, \sigma) \approx \nabla_x \log q_\sigma(x)$. به عبارت $s_\theta(x, \sigma)$ شبکه‌ی امتیاز شرطی شده با نویز می‌گوییم. در [۱۴] شبکه‌ی به کارگرفته شده برای معماری شبکه‌ی امتیاز شرطی شده با نویز یک شبکه‌ی یو^۷ [۱۲] می‌باشد که یک شبکه‌ی بسیار موفق در امر قطعه بندی معنایی^۸ تصویر است.

۵-۲-۲ به کارگیری تطبیق امتیاز برای آموزش شبکه‌ی امتیاز شرطی شده با نویز

برای آموزش شبکه‌ی امتیاز، از تابع هدف مشابه تطبیق امتیاز نویز زدا استفاده شده است و با توجه به اینکه $q_\sigma(\tilde{x}|x) = \mathcal{N}(\tilde{x}|x, \sigma^2 I)$ است و در نتیجه داریم $\nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x) = -(\tilde{x} - x)/\sigma^2$ و بنابراین برای یک σ مشخص داریم:

$$l(\theta; \sigma) \triangleq \frac{1}{2} \mathbb{E}_{p_{data}(x)} \mathbb{E}_{\tilde{x} \sim \mathcal{N}(x, \sigma^2 I)} \left[\left\| s_\theta(\tilde{x}, \sigma) + \frac{\tilde{x} - x}{\sigma^2} \right\|_2^2 \right] \quad (۵-۲)$$

و از ترکیب معادله‌ی بالا برای همه‌ی سطوح نویز داریم:

$$\mathcal{L}(\theta; \{\sigma_i\}_{i=1}^L) \triangleq \frac{1}{L} \sum_{i=1}^L \lambda(\sigma_i) l(\theta; \sigma_i) \quad (۶-۲)$$

که $\lambda(\sigma_i)$ یک تابع ضریب از سطوح نویز است. تابع هدف تعریف شده برای این مساله به هیچ گونه یادگیری تقابلی^۹ و هیچ گونه نمونه برداری از تابع در فرآیند یادگیری نیاز ندارد و همین امر سبب برتری آن نسبت به سایر شبکه‌های مولد می‌شود.

۶-۲-۲ نمونه گیری از شبکه‌ی امتیاز شرطی شده با نویز

پس از آموزش شبکه‌ی امتیاز شرطی شده با نویز از یک دینامیک لانژوین تابکاری شده^{۱۰} برای تولید نمونه‌ها استفاده می‌کنیم. که ایده‌ی اصلی آن از تابکاری شبیه سازی شده^{۱۱} گرفته شده است. فرآیند این روش به این صورت است که با شروع از یک توزیع اولیه، دینامیک لانژوین را روی نمونه‌های این توزیع

^۷U-net

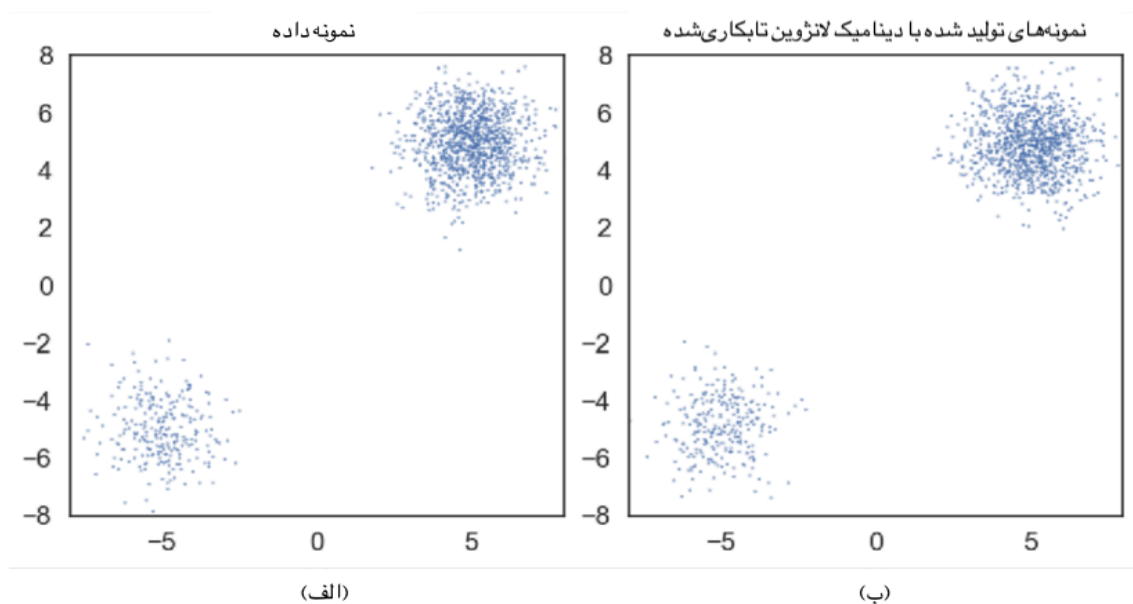
^۸Semantic segmentation

^۹Adversarial

^{۱۰}Annealed langevin dynamics

^{۱۱}Simulated annealing

شروع می‌کنیم و سپس آن را اجرا می‌کنیم تا از توزیع $q_{\sigma_1}(x)$ نمونه‌گیری کنیم. سپس نمونه‌های تولید شده را به عنوان نقطه‌ی شروع نمونه برداری با دینامیک لانژوین در توزیع $q_{\sigma_2}(x)$ در نظر می‌گیریم و در هر گام اندازه‌ی حرکت را کاهش می‌دهیم تا زمانی که دینامیک لانژوین را برای نمونه‌گیری از $q_{\sigma_L}(x)$ انجام دهیم که بسیار نزدیک به توزیع اصلی داده است و می‌تواند نمونه‌هایی از داده را تولید کند. نتایج رویکرد ارائه شده برای یک توزیع مخلوط گاوسی دو بعدی در شکل ۲-۴ آمده است که نشان می‌دهد که برخلاف دینامیک لانژوین این روش به خوبی امکان تشخیص نسبت وزنی قله‌ها به یکدیگر را دارد.



شکل ۲-۴: حل مشکل نواحی کمتر چگال در دینامیک لانژوین تابکاری شده

شکل (الف) نمونه‌های یک توزیع مخلوط گاوسی را نشان می‌دهد که از دو قله با وزن‌های متفاوت ایجاد شده است. شکل (ب) نمونه‌های تولید شده با دینامیک لانژوین تابکاری شده از همین توزیع را نشان می‌دهد. همانطور که مشاهده می‌شود، دینامیک لانژوین تابکاری شده برخلاف دینامیک لانژوین قابلیت درک وزن نسبی قله‌ها را دارد [۱۴].

۳-۲ بهبود مدل‌های مبتنی بر امتیاز

در بخش‌های قبل به معرفی و شرح جزئیات مدل‌های مبتنی بر امتیاز پرداختیم. این مدل‌ها دو مزیت اصلی نسبت به سایرین دارند که عبارتست از:

- در فرآیند آموزش نیازی به نمونه‌گیری نداریم و همین موضوع سبب می‌شود که برای آموزش شبکه‌های عمیق، کاراتر و مقیاس پذیرتر باشند.

- فرآیند نمونه برداری و آموزش کاملاً جدا از یکدیگر هستند و این امر امکان خلاقیت و استفاده از سایر روش‌های نمونه برداری مبتنی بر امتیاز را می‌دهد.

اما با وجود این مزایا و علیرغم عملکرد مناسب این روش‌ها، مدل‌های مبتنی بر امتیاز معرفی شده تنها روی تصاویر با رزولوشن پایین عملکرد مناسبی دارند و در رزولوشن‌های بالاتر این عملکرد با عوامل محدودکننده‌ای روبرو می‌شود. از جمله‌ی این عوامل می‌توان به این موضوع اشاره کرد که برای یادگیری شبکه‌ی امتیاز شرطی شده با نویز از روش تطبیق امتیاز نویزدا استفاده می‌شود و در نتیجه در این فرآیند نیازمند تخریب داده با سطوح نویز مختلف هستیم و این سطوح نویز باید به گونه‌ای باشند که شبکه‌ی امتیاز هم کلیات و هم جزئیات تصویر را بیاموزد. از این رو انتخاب سطوح نویز بسیار اهمیت دارد. در پیاده سازی توصیف شده در قسمت‌های قبل، این عملکرد برای تصاویر با رزولوشن پایین مناسب بود اما هیچ تکنیک یا ترفندی برای رزولوشن‌های بالاتر ارائه نشد. در ادامه دو تکنیک موثر و با پشتوانه‌ی تئوری را مختصراً توضیح می‌دهیم که مدل‌های مولد مبتنی بر امتیاز را توانمند می‌سازد تا تصاویری با رزولوشن بالاتر تولید کنند.

۱-۳-۲ انتخاب دنباله‌ی نویز

انتخاب دنباله‌ی نویز برای موفقیت شبکه‌های امتیاز شرطی شده با نویز بسیار مهم است. اگر شبکه تنها با یک نویز آموزش داده شود و یا نویزها بسیار کوچک باشند، تصاویر تولید شده کیفیت مناسبی ندارند چراکه تخمین تابع امتیاز به درستی صورت نمی‌گیرد و اگر نویزها خیلی بزرگ باشند نمونه‌های تولیدی خرابی دارند. برای تعیین دنباله‌ی نویز $\{\sigma_i\}_{i=1}^L$ باید دو پارامتر نویز اولیه یا σ_1 و L یا طول دنباله‌ی نویز را تعیین کنیم. پارامتر نویز اولیه تنوع داده‌ها را کنترل می‌کند و هر چقدر که بزرگتر باشد تنوع بیشتری را تضمین می‌کند و از طرفی بزرگتر شدن σ_1 به مقادیر نویز کاهش یافته‌ی بیشتری نیاز دارد و این امر دینامیک لانژوین را پرهزینه‌تر می‌کند. به طور تئوری می‌توان نشان داد که تعیین σ_1 باید به گونه‌ای باشد که از نظر عددی قابل مقایسه با حداکثر فاصله‌ی جفت داده‌ی موجود در مجموعه داده باشد [۱۵]. انتخاب دو نویز متوالی σ_i و σ_{i-1} باید به گونه‌ای باشد که $p_{\sigma_i}(x)$ نمونه‌های کافی در ناحیه‌های چگال $p_{\sigma_{i-1}}(x)$ داشته باشد. این امر اصل اساسی توسعه‌ی دینامیک لانژوین تابکاری شده است. می‌توان به صورت تئوری نشان داد که برای اینکه اصل بالا برقرار باشد باید یک دنباله‌ی هندسی برای ساخت نویز در نظر بگیریم به صورتی که نسبت $\lambda = \frac{\sigma_{i-1}}{\sigma_i}$ آن طوری تعیین شود که رابطه‌ی زیر را ارضا کند:

$$\Phi(\sqrt{2D}(\gamma - 1) + 3\gamma) - \Phi(\sqrt{2D}(\gamma - 1) - 3\gamma) \approx 0.5 \quad (۷-۲)$$

که D بعد داده‌ها و $\phi(\cdot)$ تابع چگالی احتمال گاوسی است [۱۵].

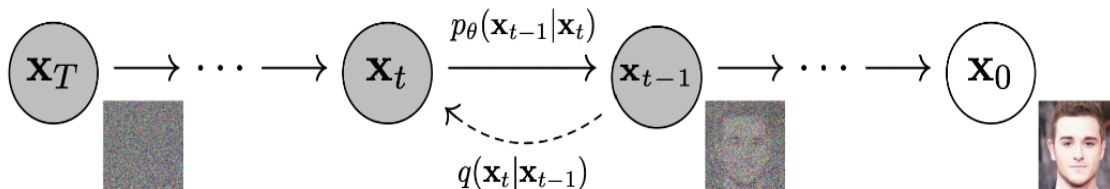
فصل سوم

مدل‌های پخش‌ی نويززدای احتمالاتی

در اين فصل به معرفي و بررسي مدل‌هاي پخش‌ي نويززداي احتمالاتي^۱ يا به اختصار مپنا مي‌پردازيم. اين مدل‌ها نيز از همان قاعده‌ي كلي مدل‌هاي پخش‌ي كه اضافه كردن نويز در يك فرآيند رو به جلو و حذف آن به صورت معكوس است پيروي مي‌كنند با اين تفاوت كه براي اين فرآيند از يك زنجيره‌ي احتمالاتي ماركفي استفاده مي‌كنند. در اين فصل ابتدا به تعريف خود مدل و تابع هزينه‌ي آن مي‌پردازيم، سپس با پارامتر ي سازي مجدد تابع هزينه ، هم ارزي اين روش با روش معرفي شده در فصل دوم را نشان مي‌دهيم و در نهايت در خصوص فرآيند يادگيري و بهبودهاي اين مدل مطلبي را مطرح مي‌كنيم.

۱-۳ تعريف مدل

مپناها مدل‌هايي مبتني بر متغير پنهان^۲ هستند كه از يك زنجيره‌ي ماركفي براي فرآيند پخش‌ي استفاده مي‌كنند. رويكرد كلي اين مدل‌ها به اين صورت است كه در يك مسير رو به جلو هر توزيع پيچيده‌اي از داده‌ها را مي‌گيرد و آن را به يك توزيع اوليه كه ساده و قابل محاسبه است تبديل مي‌كند و سپس يك مسير روبه عقب محدود به زمان را ياد مي‌گيرد كه نمونه‌هايي از توزيع اوليه را به نمونه‌هاي توزيع داده تبديل كند و اين همان مدل مولد مي‌باشد. تصوير ۱-۳ اين فرآيند را نشان مي‌دهد. بنابر اين اين مدل سه مشخصه‌ي اصلي دارد: مسير رو به جلو، مسير رو به عقب و تابع مدل مولد كه در ادامه هر يك را توضيح مي‌دهيم.



شكل ۱-۳: فرآيند پخش در مپناها

مدل گرافي جهت‌دار يك مپنا كه به صورت يك زنجيره‌ي ماركف روبه جلو و معكوس آن مي‌باشد [۶].

۱-۱-۳ مسير رو به جلو

اگر توزيع داده را با $q(x^{(0)})$ نشان دهيم، اين توزيع به صورت مرحله به مرحله به يك توزيع $\pi(y)$ كه از نظر آناليزي قابل بيان است تبديل مي‌شود و اين تبديل توسط يك هسته‌ي $T_\pi(y|y; \beta)$ صورت مي‌گيرد

^۱Denoising diffusion probabilistic models

^۲Latent variable models

که β را نرخ پخش^۲ می‌گوییم و داریم:

$$\pi(y) = \int dy' T_{\pi}(y|y'; \beta) \pi(y') \quad (۱-۳)$$

$$q(x^{(t)}|x^{(t-1)}) = T_{\pi}(x^{(t)}|x^{(t-1)}; \beta_t) \quad (۲-۳)$$

و بنابراین مسیر رو به جلو^۴ با شروع از توزیع داده‌ها و طی T گام پخششی به صورت زیر است:

$$q(x^{(0...T)}) = q(x^{(0)}) \prod_{t=1}^T q(x^{(t)}|x^{(t-1)}) \quad (۳-۳)$$

در ادامه‌ی این فصل هسته‌ی $q(x^{(t)}|x^{(t-1)})$ هسته‌ی گاوسی در نظر گرفته می‌شود اما هسته‌های دیگری چون هسته‌ی دوجمله‌ای^۵ نیز وجود دارند. یک ویژگی حائز اهمیت مسیر رو به جلو این است که می‌توان به جای استفاده از فرم زنجیره‌ی مارکوفی آن از یک فرم بسته استفاده کرد به این معنا که برای یافتن $x^{(t)}$ در زمان دلخواه t می‌توان از رابطه‌ی زیر کمک گرفت:

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I) \quad (۴-۳)$$

که $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$ و $\alpha_t := 1 - \beta_t$ تعریف می‌شود.

۲-۱-۳ مسیر رو به عقب

برای یادگیری توزیع مولد نیاز است که بتوانیم برعکس مسیر رو به جلو را یاد بگیریم و این به این معناست که باید بتوانیم از یک توزیع دلخواه قابل محاسبه به توزیع داده برسیم. بنابراین با تعریف توزیع اولیه به

^۳Diffusion rate

^۴Forward trajectory

^۵Binomial kernel

صورت $p(x^{(T)}) = \pi(x^{(T)})$ داريم:

$$p(x^{(0...T)}) = p(x^{(T)}) \prod_{t=1}^T p(x^{(t-1)} | x^{(t)}) \quad (5-3)$$

مي‌توان نشان داد كه اگر فرآيند پخش گاوسي باشد و گام‌هاي پخش β خيلي كوچك باشند آنگاه معكوس فرآيند پخش نيز گاوسي مي‌باشد [2]. و بنابر اين هر چقدر كه طول اين زنجيره بيشتر باشد β كوچكتر مي‌شود و مي‌توان از تابع گاوسي براي مسير معكوس نيز كمك گرفت. در فرآيند يادگيري با توجه به اينكه هسته‌ي گذر⁶ زنجيره‌ي ماركفي يك تابع گاوسي است، براي ساخت زنجيره تنها نياز به ميانگين و كوواريانس اين تابع گاوسي داريم و در ادامه سعي مي‌كنيم آنها را تخمين بزنيم.

3-1-3 تابع مدل مولد

تابع احتمال توليد شده توسط مپنا به صورت زير مي‌باشد:

$$p(x^{(0)}) = \int dx^{(1...T)} p(x^{(0...T)}) \quad (6-3)$$

از آنجا يكيه انتگرال در حالت كلي غير قابل محاسبه⁷ است، مي‌توان با ايده گرفتن از نمونه برداري اهميت⁸ به جاي محاسبه‌ي انتگرال از ميانگين احتمال نسبي مسير رو به جلو و مسير رو به عقب استفاده كرد كه به صورت زير قابل بيان است:

$$p(x^{(0)}) = \int dx^{(1...T)} q(x^{(1...T)} | x^{(0)}) \prod_{t=1}^T \frac{p(x^{(t-1)} | x^{(t)})}{q(x^{(t)} | x^{(t-1)})} \quad (7-3)$$

3-2 تعريف تابع خطا

با توجه به مدل تعريف شده در قسمت قبل، آموزش شبكه معادل بيشينه كردن لگاريتم درست‌نمايي⁹ تابع مدل مولد و يا كمينه كردن منفي اين تابع است. بنابر اين تابع هزينه را مي‌توان به صورت زير نوشت:

⁶Transition kernel

⁷Intractable

⁸Importance sampling

⁹Log-likelihood

$$L = \int dx^{(0)} q(x^{(0)}) \log p(x^{(0)}) \quad (۸-۳)$$

در اینجا برای آموزش به جای بهینه سازی مستقیم لگاریتم درست‌نمایی از بهینه سازی حد تغییراتی^{۱۰} آن کمک می‌گیریم که به صورت زیر قابل تعریف است:

$$\mathbb{E}[-\log p_{\theta}(x^0)] \leq \mathbb{E}_q[-\log \frac{p_{\theta}(x^{0:T})}{q(x^{1:T}|x^0)}] = \mathbb{E}_q[-\log p(x^T) - \sum_{t \geq 1} \log \frac{p_{\theta}(x^{t-1}|x^t)}{q(x^t|x^{t-1})}] =: L \quad (۹-۳)$$

۳-۳ تغییر پارامتر تابع خطا

از آنجایی که ترکیب مسیر رو به جلو و مسیر معکوس شبیه به یک کدگذار تغییراتی است، می‌توان تابع هزینه را به صورت زیر بازنویسی کرد:

$$\mathcal{L}_{vib} = \mathcal{L}_0 + \mathcal{L}_1 + \dots + \mathcal{L}_{T-1} + \mathcal{L}_T \quad (۱۰-۳)$$

$$\mathcal{L}_0 = -\log p_{\theta}(x_0|x_1) \quad (۱۱-۳)$$

$$\mathcal{L}_{t-1} = \mathcal{D}_{KL}(q(x_{t-1}|x_t, x_0) || p_{\theta}(x_{t-1}|x_t)) \quad (۱۲-۳)$$

$$\mathcal{L}_T = \mathcal{D}_{KL}(q(x_T|x_0) || p(x_T)) \quad (۱۳-۳)$$

و از آنجایی که در صورت کوچک بودن β_t هر دوی $q(x_{t-1}|x_t, x_0)$ و $p_{\theta}(x_{t-1}|x_t)$ توزیع‌هایی گاوسی هستند در نتیجه دیورجنس KL ^{۱۱} آنها را می‌توان به فرم تفاضل میانگین آنها در نظر گرفت و در نتیجه برای معادله ی (۱۲-۳) داریم:

^{۱۰}Variational bound

^{۱۱}KL-divergence

$$\mathcal{L}_{t-1} = \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t)\|^2 \right] + C \quad (14-3)$$

که در واقع نشان می‌دهد که در فرآیند آموزش یک مدل پخش‌ی بهینه سازی خطا شامل تخمین میانگین پسین^{۱۲} فرآیند رو به جلو می‌باشد. از طرفی با تغییر پارامتر $x_t(x_0, \epsilon) = \sqrt{\alpha_t}x_0 + \sqrt{(1 - \alpha_t)}\epsilon$ برای $\epsilon \sim \mathcal{N}(0, I)$ و قراردادی آن در فرمول معادله‌ی (۱۴-۳) می‌توان نشان داد که تخمین $\mu_\theta(x_t, t)$ به صورت زیر قابل پارامتری سازی است:

$$\mu_\theta(x_t, t) = \tilde{\mu}_t \left(x_t, \frac{1}{\sqrt{\alpha_t}} (x_t - \sqrt{1 - \alpha_t} \epsilon_\theta(x_t)) \right) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(x_t, t) \right) \quad (15-3)$$

که ϵ_θ یک تخمین از نویز موجود در x_t می‌باشد و در نتیجه برای محاسبه ی یک نمونه ی x_{t-1} تنها کافیست که عبارت زیر را محاسبه کنیم:

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(x_t, t) \right) + \sigma_t z \quad (16-3)$$

که z از توزیع گاوسی با میانگین صفر و کوواریانس همانی می‌آید. این تغییر پارامتر دو موضوع مهم را نشان می‌دهد: اول اینکه برای آموزش یک فرآیند رو به عقب می‌توان از تخمین $\tilde{\mu}_t$ با تابع میانگین وابسته به زمان μ_θ استفاده کرد و یا مستقیماً نویز اضافه شده به داده را با تخمین ϵ محاسبه کرد. دوم اینکه در حالتی از پارامتری سازی که از تخمین نویز در هر گام زمانی استفاده می‌کنیم در واقع فرآیندی مشابه تطبیق امتیاز نويززدا با دینامیک لانژوین را طی می‌کنیم و در این حالت تابع هدف شبکه را تبدیل به یک مدل تطبیق امتیاز نويززدا می‌کنیم.

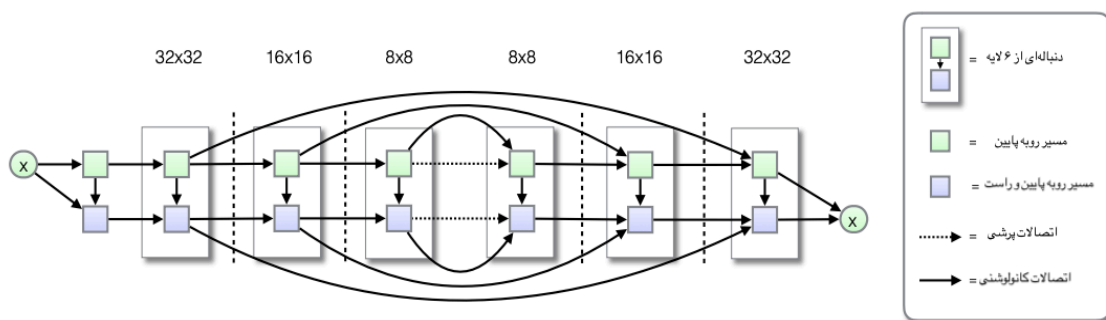
۴-۳ فرآیند یادگیری

با توجه به تغییر پارامتر انجام شده در قسمت قبل و همچنین به کارگیری آن در معادله‌ی (۱۰-۳) می‌توان به تابع هزینه ای رسید که کاملاً نسبت به θ مشتق پذیر است و برای آموزش مناسب است اما برای سادگی پیاده سازی و افزایش کیفیت نمونه‌های تولیدی می‌توان تابع هزینه را به شکل زیر نیز بازنویسی کرد:

¹²Posterior mean

$$L_{simple}(\theta) = \mathbb{E}_{t, x_0, \epsilon} \left[\left\| \epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t) \right\|^2 \right] \quad (17-3)$$

که t از یک توزیع یونیفرم بین ۱ و T می‌آید و واریانس β_t در این فرآیند یک پارامتر ثابت بدون یادگیری می‌باشد. در ساختار مپناها از هسته‌ی اصلی $++ pixelCNN$ [۱۳] استفاده شده است که در واقع نوعی شبکه‌ی یو با اتصالات گسترده‌ی باقی‌مانده‌ای می‌باشد. شکل ۲-۳ ساختار کلی شبکه‌ی به کار گرفته شده را نشان می‌دهد.



شکل ۲-۳: معماری $++ pixelCNN$

شبکه‌ی $++ pixelCNN$ که نوعی شبکه‌ی یو با اتصالات باقی‌مانده‌ای گسترده است به عنوان شبکه‌ی اصلی در مپناها استفاده می‌شود [۱۳].

۳-۵ بهبود مدل‌های پخش‌ی نويززداي احتمالاتي

با وجود عملکرد فوق‌العاده‌ی مپناها در تولید نمونه‌های با کیفیت بالا، جنبه‌های مختلفی از این مدل‌ها نیاز به بازنگری و ارزیابی دارند. از جمله اینکه این مدل‌ها تا چه حد توانایی پیدا کردن همه‌ی مد^{۱۳}های یک توزیع احتمالاتی را دارند و یا اینکه در مجموعه داده‌های متنوع مانند *ImageNET*، آیا قابلیت تولید نمونه‌های متنوع را دارند؟ لگاریتم درست‌نمایی، یک معیار پرکاربرد در بحث مدل‌های مولد است و بهینه‌سازی آن سبب می‌شود تا مدل مولد بتواند همه‌ی مد‌های یک توزیع را تولید کند [۱۱]. همچنین اخیراً نشان داده شده است که بهبود ناچیزی در بهینه‌سازی لگاریتم درست‌نمایی موجب تغییر چشم‌گیری در کیفیت نمونه و بازنمایی ویژگی‌های آموخته شده می‌شود [۴]. در مپنا با وجود عملکرد بسیار عالی در تولید تصاویر با کیفیت بالا مقدار لگاریتم درست‌نمایی قابل مقایسه با سایر روش‌های مولد نیست. اولین رویکردی که مسلماً باعث بهبود لگاریتم درست‌نمایی در مپناها می‌شود، افزایش T یا تعداد گام است.

¹³Mode

هر چقدر T بیشتر باشد لگاریتم درست‌نمایی مپناها بهینه‌تر خواهد بود. از طرفی این افزایش موجب ناکارآمدی زمانی و هزینه‌ای مپناها می‌شود. در ادامه‌ی این فصل دو رویکرد برای بهینه‌تر کردن لگاریتم درست‌نمایی در مپناها را به اختصار معرفی می‌کنیم.

۳-۵-۱ یادگیری واریانس

در قسمت‌های قبل و در هنگام محاسبه‌ی تابع هزینه گفتیم که واریانس را ثابت و بدون یادگیری در نظر می‌گیریم. ثابت در نظر گرفتن σ_t از دیدگاه افزایش کیفیت نمونه‌های تولیدی انتخاب مناسب و کارایی است اما از آنجاییکه میزان تاثیر جملات مختلف حد تغییراتی در گام‌های ابتدایی بیشتر است پس برای بهبود لگاریتم درست‌نمایی یادگیری واریانس می‌تواند باعث بهبود این مقدار شود.

۳-۵-۲ بهبود گام‌های نويز

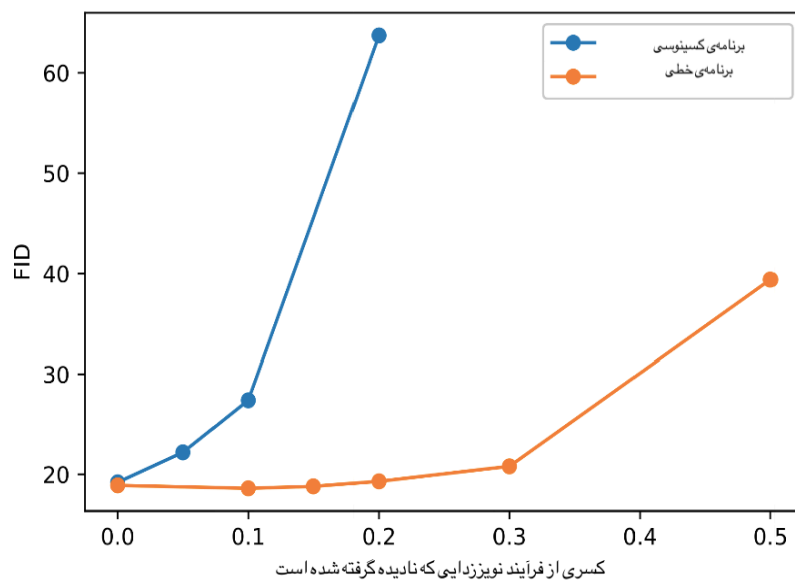
در مپنا مقدار نويز β_t با یک تابع خطی تغییر می‌کند. این تابع خطی برای تصاویر با رزولوشن بالا خوب عمل می‌کند اما برای تصاویر با رزولوشن پایین مناسب نیست. این امر به این علت است که انتهای فرآیند پخش شدت نويز تا حدی افزایش می‌یابد که عملاً در کیفیت نمونه بی‌تاثیر می‌شود. این موضوع در شکل ۳-۳ به خوبی نشان داده شده است. همچنین همانطور که در شکل ۳-۴ مشاهده می‌کنید در مدلی که با نويز خطی پخش شده است، نادیده گرفتن ۲۰ درصد از فرآیند معکوس تاثیر مخربی بر کیفیت نمونه‌ها نمی‌گذارد.



شکل ۳-۳: تفاوت برنامه‌ی نويز خطی و کسینوسی در فرآیند پخش‌ی

نمونه‌های فرآیند پخش‌ی با برنامه‌ی نويز خطی (بالا) و نويز کسینوسی (پایین). نمونه‌های یک چهارم آخر برنامه‌ی نويز خطی عملاً همگی نويز هستند و فرآیند نويزافزایی خیلی آهسته صورت گرفته است [۱۰].

برای حل این مشکل می‌توانیم از انواع دیگری از برنامه‌های نويز مانند برنامه‌ی نويز کسینوسی استفاده کنیم. همان طور که در شکل ۳-۴ می‌بینید، این برنامه‌ی نويز با سرعت آهسته‌تری از برنامه‌ی نويز خطی به تصویر نويز اضافه می‌کند. همچنین همانطور که در شکل ۳-۴ قابل مشاهده است حذف بخشی از فرآیند رو به عقب در این برنامه‌ی نويز باعث افت کیفیت تصویر می‌شود.



شکل ۳-۴: تفاوت برنامه‌ی نویز خطی و کسینوسی در حذف بخشی از فرآیند نویززدایی در مدلی که با نویز خطی پخش شده است، نادیده گرفتن ۲۰ درصد از فرآیند معکوس تاثیر مخربی بر کیفیت نمونه‌ها نمی‌گذارد در حالیکه در مدل کسینوسی حذف هر کسری از فرآیند معکوس کیفیت را بدتر می‌کند [۱۰].

فصل چهارم

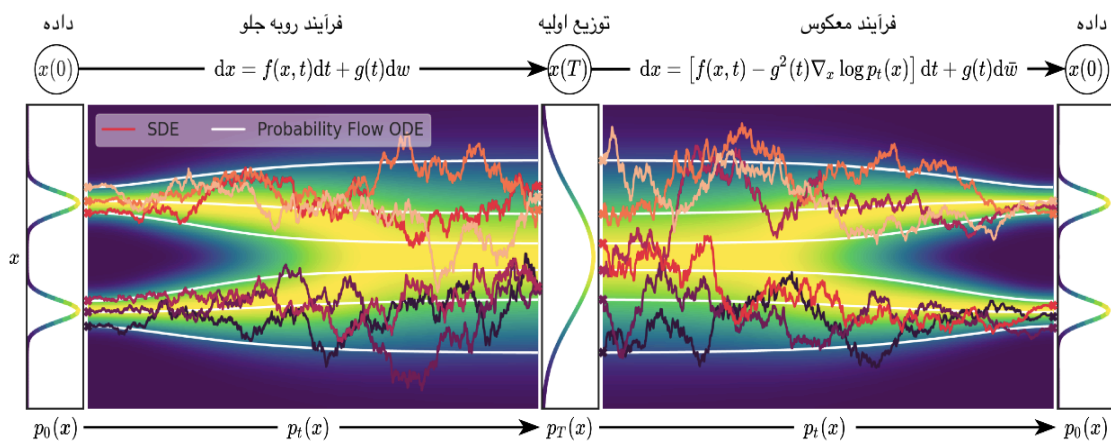
مدل‌های مولد مبتنی بر معادلات دیفرانسیل

تصادفی

در دو فصل گذشته دو گونه از مدل‌های پخش را معرفی کردیم که نقطه‌ی اشتراک هر دو این بود که در یک فرآیند پخش داده را با نویز فزاینده تخریب می‌کردند و سپس در یک فرآیند معکوس از این داده‌ی تخریب شده به مدلی مولد برای تولید داده‌ی جدید می‌رسیدند. در این فصل یک چهارچوب کلی برای فرآیند پخش تعریف می‌کنیم که به جای تخریب داده با یک سری گام گسسته‌ی متناهی، داده را با بی نهایت گام پیوسته تخریب می‌کنیم. این فرآیند تخریب از طریق یک معادله دیفرانسیل تصادفی تعریف شده صورت می‌گیرد و سپس نشان داده می‌شود که برای فرآیند تصادفی یک معادله دیفرانسیل تصادفی معکوس نیز وجود دارد که تنها وابسته به تابع امتیاز می‌باشد. در این بخش ابتدا به معرفی این چهارچوب کلی و مزایای آن می‌پردازیم، سپس ارتباط آن با دو روش پیش‌تر معرفی شده را بیان می‌کنیم و در نهایت به انواع روش‌های نمونه برداری از این توزیع می‌پردازیم.

۴-۱ معرفی فرآیند پخش با معادلات دیفرانسیل تصادفی

ایده‌ی اصلی در مدل‌هایی که مبتنی بر معادلات دیفرانسیل تصادفی هستند این است که تعداد گام‌های حرکت از توزیع داده به سمت توزیع نویز بی نهایت و به صورت حدی کوچک باشد به طوریکه توزیع داده‌های تخریب شده با افزایش نویز به صورت یک معادله دیفرانسیل تصادفی تغییر کند. مانند مدل‌های قبلی این مدل نیز از سه بخش مسیر رو به جلو، مسیر معکوس و تابع هدف یادگیری تشکیل شده است. تصویر ۴-۱ شمای کلی این روش را نشان می‌دهد.



شکل ۴-۱: مدل‌های پخش مبتنی بر معادلات دیفرانسیل

مدل‌های پخش مبتنی بر معادلات دیفرانسیل از فرآیندهای تصادفی و این معادلات برای مدل کردن فرآیند رو به جلو و فرآیند معکوس استفاده می‌کنند [۱۶].

۴-۱-۱ مسیر روبه جلو

همانطور که پیش‌تر هم اشاره شد در این مدل هدف توسعه‌ی فرآیند پخش $\{x(t)\}_{t=0}^T$ است به طوریکه $x(0) \sim p_0$ یا همان توزیع داده‌ها و $x(T) \sim p_T$ یا همان توزیع اولیه^۱ باشد و متغیر t یک متغیر پیوسته در بازه‌ی $[0, T]$ باشد. در چنین حالتی فرآیند پخش را یک راه حل معادله دیفرانسیل تصادفی به فرم زیر در نظر می‌گیریم:

$$dx = f(x, t)dt + g(t)dw \quad (۴-۱)$$

که در عبارت بالا $f(t, \cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ ضریب یا تابع رانش^۲ و $g(\cdot)$ یک مقدار عددی وابسته به زمان است که ضریب پخش^۳ گفته می‌شود. w نیز فرآیند وینر استاندارد^۴ که معادل همان نویز سفید تصادفی در فضای پیوسته می‌باشد. به طور کلی توزیع اولیه p_T هر توزیع دلخواه قابل محاسبه‌ای است که هیچ اطلاعاتی از توزیع داده در خود ندارد. مانند یک توزیع گاوسی با میانگین و کوواریانس ثابت. حال برای هر توزیع اولیه‌ی مشخص یا هر شکلی از فرآیند پخش یک معادله دیفرانسیل تصادفی طراحی می‌شود که امکان تبدیل توزیع داده به توزیع مشخص شده را دارد.

۴-۱-۲ مسیر معکوس

می‌دانیم که معکوس یک فرآیند پخشی که با معادلات دیفرانسیل مدل می‌شود نیز یک فرآیند پخشی است که در جهت عکس زمانی حرکت می‌کند و با معادله دیفرانسیل تصادفی معکوس زیر مدل می‌شود:

$$dx = [f(x, t) - g(t)^2 \nabla_x \log p_t(x)]dt + g(t)d\bar{w} \quad (۴-۲)$$

که \bar{w} همان فرآیند وینر استاندارد در جهت معکوس از T به 0 است و dt نیز در جهت منفی حرکت می‌کند^[۱]. همانطور که در معادله‌ی (۴-۲) مشاهده می‌شود این عبارت تنها نیازمند تابع امتیاز توزیع $p_t(x)$ می‌باشد و با داشتن این تابع می‌توانیم معکوس فرآیند پخشی را بسازیم و از توزیع داده‌ها نمونه تولید کنیم.

^۱Prior distribution

^۲Drift coefficient

^۳Diffusion coefficient

^۴Standard wiener process

۳-۱-۴ تخمین تابع امتیاز

تابع امتیاز یک توزیع همانطور که در فصل دوم مشاهده شد با فرآیند تطبیق امتیاز قابل تخمین و یادگیری است. تنها نکته‌ی حائز اهمیت این است که در این حالت نیاز به یک تعمیم پیوسته از تابع خطای این شبکه داریم. این تعمیم به صورت زیر قابل بیان است:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{x(0)} \mathbb{E}_{(x(t)|x(0))} \left[\left\| s_{\theta}(x(t), t) - \nabla_{x(t)} \log p_{0t}(x(t)|x(0)) \right\|_2^2 \right] \right\} \quad (3-4)$$

که $\lambda(t)$ یک تابع وزن‌دهی مثبت است و t به صورت یکنواخت از بازه‌ی $[0, T]$ نمونه برداری شده است و $x(0) \sim p_0(x)$ و $x(t) \sim p(x(t)|x(0))$ می‌باشد. اگر ظرفیت مدل مناسب و تعداد داده‌ها کافی باشد معادله‌ی (۳-۴) ما را به پارامترهای بهینه‌ی تابع امتیاز می‌رساند [۱۶].

۲-۴ مزیت‌های مدل‌های مولد مبتنی بر معادله دیفرانسیل تصادفی

سه مزیت کلی برای این دسته از روش‌ها می‌توان نام برد که عبارتند از:

- نمونه برداری و محاسبه‌ی درست‌نمایی منعطف: در این روش از هر حل‌کننده‌ی کلی معادلات دیفرانسیل تصادفی برای نمونه‌گیری از معادله دیفرانسیل تصادفی معکوس می‌توان کمک گرفت. علاوه بر این می‌توان از روش‌های دیگری حل‌کننده‌های قطعی نیز برای نمونه‌گیری کمک گرفت.
- تولید قابل کنترل: از آنجاییکه معادله دیفرانسیل تصادفی معکوس شرطی شده از توابع امتیاز غیر شرطی قابل محاسبه است، می‌توانیم فرآیند تولید نمونه را روی اطلاعاتی شرطی کنیم که در حین یادگیری در اختیار مدل قرار نگرفته است. این امر باعث می‌شود تا کاربردهایی مثل ترمیم تصویر^۵، تولید شرطی شده روی کلاس و رنگ آمیزی^۶ با یک مدل امتیاز غیرشرطی و بدون آنکه نیاز به آموزش مجدد باشد، قابل انجام باشند.
- ایجاد یک چهارچوب یکپارچه: این روش قابلیت جستجو و به کارگیری انواع مختلفی از معادلات دیفرانسیل تصادفی را برای بهبود مدل‌های پخش‌ی دارد. همچنین هر دو روش پیش‌تر معرفی شده نیز در قالب این روش قابل بیان و بهینه‌سازی هستند.

⁵Image inpainting

⁶Colorization

۳-۴ ارتباط روش مبتنی بر معادلات دیفرانسیل تصادفی با روش‌های

پیش‌تر معرفی شده

به طور کلی می‌توان گفت فرآیند پخش به کار گرفته شده در مپنا و تطبیق امتیاز با دینامیک لانژوین، گسسته‌سازی‌هایی از روش ارائه شده در این فصل هستند. هر کدام از این روش‌ها از یک معادله دیفرانسیل تصادفی برخوردار هستند که با گسسته‌سازی آن می‌توان هم ارزی روش‌های پیش‌تر معرفی شده را با روش مبتنی بر معادله دیفرانسیل تصادفی نشان داد. در روش تطبیق امتیاز با دینامیک لانژوین با استفاده از N سطح نویز می‌توان تابع تخریب $p_{\sigma_i}(x|x_0)$ را با توزیع x_i حاصل از یک زنجیره‌ی مارکف به شکل زیر به دست آورد:

$$x_i = x_{i-1} + \sqrt{\sigma_i^2 - \sigma_{i-1}^2} z_{i-1}, \quad i = 1, \dots, N \quad (۴-۴)$$

که $z_{i-1} \sim \mathcal{N}(0, I)$ و در صورتی که $N \rightarrow \infty$ آنگاه زنجیره‌ی مارکفی بالا به یک فرآیند تصادفی پیوسته به صورت $\{x(t)\}_{t=0}^1$ تبدیل می‌شود که خود t هم یک متغیر پیوسته در بازه‌ی $[0, T]$ است. در چنین صورتی می‌توانیم این فرآیند تصادفی را با یک معادله دیفرانسیل تصادفی به فرم زیر مدل کنیم:

$$dx = \sqrt{\frac{d[\sigma^2(t)]}{dt}} dw \quad (۵-۴)$$

این معادله، که به آن معادله‌ی انفجار واریانس^۷ می‌گوییم، در واقع حالت خاصی از چهارچوب مطرح شده در این فصل می‌باشد.

به طور مشابه برای مپنا نیز یک زنجیره‌ی مارکفی به شکل زیر داریم:

$$x_i = \sqrt{1 - \beta_i} x_{i-1} + \sqrt{\beta_i} z_{i-1}, \quad i = 1, \dots, N \quad (۶-۴)$$

که در صورتی که $N \rightarrow \infty$ ، به معادله دیفرانسیل تصادفی زیر همگرا می‌شود:

$$dx = -\frac{1}{2}\beta(t)xdt + \sqrt{\beta(t)}dw \quad (۷-۴)$$

و بنابراین تخریب‌های نویز به کار برده شده در مپنا و در تطبیق امتیاز با دینامیک لانژوین هر دو گسسته سازی‌هایی از یک معادله دیفرانسیل تصادفی هستند.

^۷Variance exploding

۴-۴ انواع روش‌های نمونه برداری

بعد از آموزش یک شبکه‌ی امتیاز وابسته به زمان می‌توان از آن برای ساخت یک معادله دیفرانسیل تصادفی معکوس استفاده کرد و سپس با استفاده از روش‌های عددی حل این معادله دیفرانسیل آن را شبیه‌سازی و نمونه‌هایی از توزیع p_0 تولید کرد. در این قسمت به دو دسته روش حل این معادلات اشاره خواهیم کرد:

۴-۴-۱ حل کننده‌های عام منظوره ی معادلات دیفرانسیل تصادفی

حل کننده‌های عددی مسیر تقریبی معادلات دیفرانسیل تصادفی را تخمین می‌زنند. بسیاری از این حل کننده‌های عام منظوره مانند روش اویلر-ماریاما^۸ برای تولید نمونه به کار گرفته می‌شوند.

۴-۴-۲ تبدیل به معادلات دیفرانسیل ساده

با دانستن تابع امتیاز یک راه دیگر برای حل معادلات دیفرانسیل تصادفی و نمونه برداری، تبدیل این معادلات به معادلات دیفرانسیل ساده^۹ و تبدیل فرآیند تصادفی به فرآیند قطعی^{۱۰} می‌باشد. به ازای هر فرآیند پخشی تصادفی یک فرآیند قطعی متناظر آن وجود دارد که هر دو مسیر توزیع‌های حاشیه‌ای یکسانی دارند. این فرآیند قطعی به صورت زیر قابل بیان است:

$$dx = [f(x, t) - \frac{1}{2}g^2(t) \nabla_x \log p_t(x)]dt \quad (۸-۴)$$

در صورتی که تابع امتیاز را بدانیم این مسیر به راحتی از معادله دیفرانسیل تصادفی قابل محاسبه است. استفاده از معادلات دیفرانسیل ساده به جای معادلات دیفرانسیل تصادفی مزیت‌های متعددی از جمله توانایی محاسبه‌ی دقیق درست‌نمایی، کدگذاری یکتا و نمونه برداری کارا دارد.

^۸Euler-maruyama

^۹Ordinary differential equations

^{۱۰}Deterministic process

فصل پنجم

مقایسه‌ی روش‌های مولد

در این فصل ابتدا به بررسی ارتباط بین مدل‌های پخشی با سایر مدل‌های مولد و مقایسه‌ی آنها می‌پردازیم؛ سپس مدل‌های پخشی پیش‌تر معرفی شده را از دو جنبه‌ی عملکرد و نتایج به دست آمده و معماری با یکدیگر مقایسه می‌کنیم.

۱-۵ بررسی ارتباط مدل‌های پخشی با سایر مدل‌های مولد

۱-۱-۵ ارتباط خودکدگذارهای تغییراتی با مدل‌های پخشی

خودکدگذارهای تغییراتی [۸] شامل یک کدگذار^۱ و یک کدگشا^۲ هستند که هر دو برای پیدا کردن نگاشتی بین داده‌های ورودی و یک فضای پنهان پیوسته آموزش داده می‌شوند. در این مدل‌ها تعبیه^۳ به دست آمده به عنوان یک متغیر پنهان در یک مدل مولد احتمالاتی در نظر گرفته می‌شود و یک کدگشای احتمالاتی آموزش داده می‌شود تا از این متغیر پنهان به داده برسد. بنابراین فرض می‌شود که داده‌ی x از یک متغیر پنهان مشاهده نشده‌ی z با یک توزیع شرطی $p_\theta(x|z)$ تولید می‌شود و $q_\phi(z|x)$ برای تخمین z به کار گرفته می‌شود. در این مدل برای تضمین دستیابی به تخمین مناسبی از z از یک رویکرد تغییراتی بیزی برای بیشینه کردن حد پایین مشاهده به فرم $\mathcal{L}(\phi, \theta, x) \leq \log p_\theta(x)$ که داریم:

$$\mathcal{L}(\phi, \theta; x) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x, z) - \log q_\phi(z|x)] \quad (۱-۵)$$

استفاده می‌شود. برای این نوع فرموله سازی این مدل انتخاب‌های متعددی برای مدل‌های کدگذار و کدگشا وجود دارد، معمولاً این مدل‌ها با خانواده‌ی توزیع‌های نمایی مدل می‌شوند که پارامترهای این توزیع از شبکه‌های عصبی چند لایه تولید می‌شود. مپناها را می‌توان به عنوان خودکدگذارهای تغییراتی سلسله مراتبی^۴ صورت بندی کرد که در آنها کدگذار ثابت است. در چنین صورت بندی مسیر رو به جلو نقش کدگذار را ایفا می‌کند و به عنوان یک مدل گاوسی خطی در نظر گرفته می‌شود و مسیر روبه عقب نیز نقش کدگشا را ایفا می‌کند که در گام‌های مختلف نویززدایی به کار گرفته می‌شود و همچنین اندازه‌ی متغیر پنهان به کار گرفته شده در کدگشا نیز با اندازه‌ی داده‌ی ورودی برابر است. شکل ۱-۵ این موضوع را نشان می‌دهد.

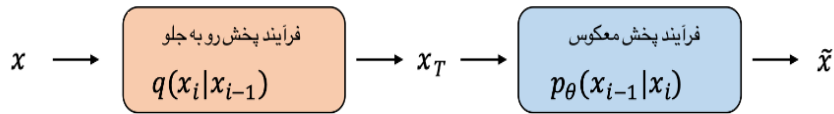
^۱Encoder

^۲Decoder

^۳Embedding

^۴Hierarchical variational autoencoders

خودکدگذار تغییراتی



شکل ۵-۱: مدل‌های پخشی مبتنی بر معادلات دیفرانسیل

مدل‌های پخشی مشابه یک خودکدگذار تغییراتی سلسله مراتبی هستند که کدگذار در آنها ثابت است [۲۰].

۵-۱-۲ ارتباط شبکه‌های مولد تقابلی با مدل‌های پخشی

شبکه‌های مولد تقابلی [۳] از دو زیر شبکه‌ی تولید کننده‌ی G و تمایزگر D تشکیل شده‌اند که معمولاً به فرم شبکه‌های عصبی پیاده سازی می‌شوند اما در هر فرمی از سیستم‌های مشتق پذیر که داده را از یک فضا به فضای دیگر می‌برند نیز قابل پیاده سازی هستند. تابع هزینه‌ی شبکه‌های مولد تقابلی به صورت زیر قابل بیان است:

$$\min_G \max_D \mathbb{E}_{x \sim p_{data}} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (۲-۵)$$

شبکه‌ی تولید کننده تلاش می‌کند با تولید نمونه‌های جدید به صورت ضمنی توزیع داده را مدل کند و در مقابل شبکه‌ی تمایزگر تلاش می‌کند که با حداکثر دقت نمونه‌های تولید شده را از نمونه‌های واقعی جدا کند. فرآیند بهینه سازی این شبکه تا جایی ادامه پیدا می‌کند که به یک نقطه‌ی زینی^۵ که کمینه برای شبکه‌ی تولید کننده و بیشینه برای شبکه‌ی تمایزگر است برسد. در چنین نقطه‌ای شبکه‌ی مولد توزیع صحیح نمونه‌های داده را یاد گرفته است.

در ساختار پایه‌ای شبکه‌های مولد تقابلی رابطه‌ی واضحی با مدل‌های پخشی وجود ندارد اما این دو مدل می‌توانند در صورت ترکیب با یکدیگر مشکلات و ضعف هایشان را برطرف کنند. به عنوان مثال یک مساله‌ی خیلی رایج در شبکه‌های مولد تقابلی، مساله‌ی عدم پایداری آنها است که عموماً به علت عدم همپوشانی بین توزیع داده‌ی واقعی و توزیع داده‌ی تولید شده ایجاد می‌شود. یک راه حل این مساله این است که با اضافه کردن نویز به ورودی تمایزگر فضای پشتیبان توزیع‌های تمایزگر و تولید کننده را افزایش دهیم. این راه حل می‌تواند توسط یک مدل پخشی صورت بگیرد. در [۱۸] نویسندگان از یک مدل پخشی با برنامه‌ی نویز تطبیقی برای اضافه کردن نویز به تمایزگر استفاده کردند. از سوی دیگر شبکه‌های مولد تقابلی نیز می‌توانند در افزایش سرعت نمونه برداری مدل‌های پخشی موثر واقع شوند. در [۱۹] برای کاهش تعداد گام‌های نویزدایی از یک مدل مولد تقابلی شرطی برای مدل کردن هر گام نویزدایی استفاده شده است.

^۵Saddle point

۵-۱-۳ مشکل سه‌گانه‌ی مدل‌های مولد

با وجود توسعه‌ی حجم زیادی از مدل‌های مولد در سال‌های اخیر برای انواعی از کاربردها از جمله تصویر، صدا، متن و گراف‌ها، همه‌ی مدل‌های موجود فعلی قابلیت ارضای سه نیاز اصلی یک مدل مولد به کار گرفته شده در مسائل واقعی را ندارند [۱۹]. این سه نیاز اصلی عبارتند از:

۱. تولید نمونه‌های با کیفیت

۲. پوشش مد و تنوع نمونه‌های تولیدی

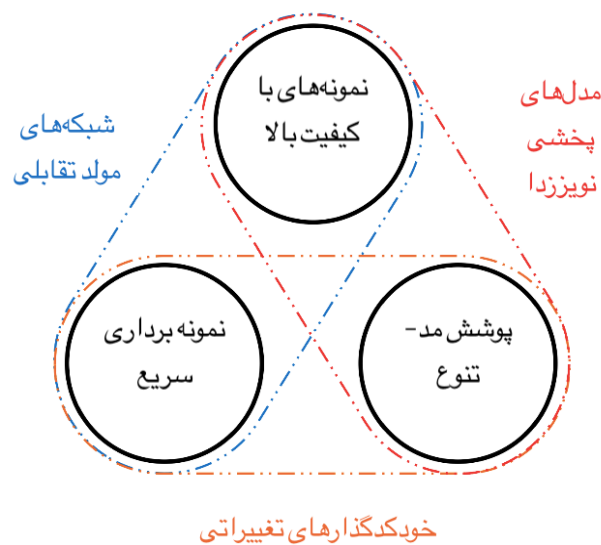
۳. نمونه برداری سریع و به صرفه از نظر محاسباتی

به عنوان مثال اکثر مدل‌های تولید تصویر فعلی روی تولید نمونه‌های با کیفیت تمرکز دارند اما تنوع داده‌ها نیز موضوع مهمی در کاهش تاثیرات اجتماعی منفی مدل‌های مولد است. همچنین برای کاربردهای تعاملی نیاز به نمونه برداری سریع یک مساله‌ی حائز اهمیت است. این چالش در ادبیات مدل‌ها مولد به عنوان مشکل سه‌گانه‌ی یادگیری مولد^۶ شناخته می‌شود. هر کدام از مدل‌های مولد در حل بخشی از مشکلات یادگیری مولد تبحر دارند و بخش‌هایی از این مشکل به عنوان ضعف آنها محسوب می‌شود. به عنوان مثال شبکه‌های مولد تقابلی قابلیت تولید نمونه‌های با کیفیت در زمان سریع را دارند اما پوشش مد ضعیفی دارند. از طرفی خودکدگذارهای تغییراتی به خوبی مدهای داده را پوشش می‌دهند و سرعت نمونه‌برداری مناسبی دارند اما کیفیت نمونه‌های تولیدی آنها کم است. برای مقایسه‌ی عملکرد مدل‌های پخشی با سایر مدل‌های مولد خوب است بدانیم که این مدل‌ها قابلیت رفع این مشکل سه‌گانه را دارند یا خیر. مدل‌های پخشی با ظهور خود توانستند از شبکه‌های مولد تقابلی در تولید نمونه‌های با کیفیت پیشی بگیرند و همچنین با تغییر و تنظیم این مدل‌ها و دستیابی به درست‌نمایی بهینه، قابلیت پوشش مدهای مختلف داده را دارند اما نمونه برداری از این مدل‌ها نیازمند صدها بار استفاده از شبکه می‌باشد و همین امر سرعت نمونه برداری را کاهش می‌دهد و به کارگیری آنها در کاربرد را هزینه‌بر و دشوار می‌سازد. شکل ۵-۲ این مشکل سه‌گانه و حوزه‌ی عملکرد مناسب هر مدل را نشان می‌دهد.

همچنین جدول ۵-۱-۳ نیز مقایسه‌ای از ویژگی‌های کلی این سه مدل مولد ارائه می‌دهد. ویژگی تعیین یکتا^۷ نشان دهنده‌ی این است که با تغییر معماری و مقداردهی اولیه تا چه حد خروجی مدل و کدگذاری آن تغییر می‌کند. در مدل‌های پخشی به شرط داشتن داده‌ی کافی، دقت بهینه‌سازی مناسب و ظرفیت لازم مدل، برخلاف مدل‌های مولد تقابلی و خودکدگذارهای تغییراتی با تغییر معماری و مقداردهی اولیه همچنان به یک خروجی و کدگذاری میرسیم [۱۶].

^۶Generative learning trilemma

^۷Unique identifiability



شکل ۵-۲: مشکل سه‌گانه‌ی مدل‌های مولد

سه نیاز اساسی در توسعه‌ی مدل‌های مولد و حوزه‌ی عملکرد موثر هر مدل [۱۹].

جدول ۵-۱: مقایسه‌ی مدل‌های مولد با یکدیگر

مدل	مدل‌های پخشی	خودکدگذارهای تغییراتی	شبکه‌های مولد تقابلی
تعداد شبکه	۱	۲	۲
نوع آموزش	نویرزدایی تکراری	استنتاج تغییراتی	تقابلی
سرعت آموزش	سریع	آهسته	آهسته
کیفیت تصاویر	بالا	متوسط	بالا
تعیین یکتا	بله	خیر	خیر

۵-۲ مقایسه‌ی انواع روش‌های پخشی

در سه فصل گذشته سه مدل از انواع روش‌های پخشی نام برده شد که هرکدام می‌توانستند هم ارز یا در برگیرنده‌ی روش دیگری باشند. در حالت کلی می‌توان گفت که مدل‌های پخشی مبتنی بر معادلات دیفرانسیل تصادفی در برگیرنده‌ی دو مدل دیگر هستند و این دو مدل با وجود اشتراک و هم ارزی با یکدیگر تحت شرایط خاص (تخمین نویر و واریانس ثابت) در حالت کلی دو فرمولاسیون متفاوت از معادلات دیفرانسیل تصادفی هستند. در این قسمت می‌خواهیم عملکرد و نتایج، معماری این مدل‌ها و برخی از ویژگی‌های کلی آنها را با یکدیگر مقایسه کنیم.

جدول ۵-۲: مقایسه‌ی کیفیت تصاویر تولیدی با سه مدل پخش‌ی روی مجموعه داده‌ی $CIFAR-10$

مدل	FID
مدل‌های مبتنی بر امتیاز	۱۰.۸۷
مپنا	۳.۱۷
مدل مبتنی بر معادله دیفرانسیل (۴-۵)	۲.۲۰
مدل مبتنی بر معادله دیفرانسیل (۴-۷)	۲.۴۱

۵-۲-۱ نتایج و عملکرد

برای مقایسه‌ی نتایج هر مدل از معیار کیفیت فاصله‌ی اکتسابی فرچت^۸ (FID) [۵] استفاده می‌کنیم و نتایج را بر روی مجموعه داده‌ی $CIFAR-10$ برای مقالات اولیه‌ی هر کدام که [۱۴] برای تطبیق امتیاز نویزدا با دینامیک لانژوین، [۶] برای مپنا و [۱۶] برای مدل‌های پخش‌ی مبتنی بر معادلات دیفرانسیل تصادفی است، بررسی می‌کنیم. فاصله‌ی اکتسابی فرچت معیاری است که فاصله‌ی بین بردار ویژگی تصاویر اصلی و تصاویر تولید شده را محاسبه می‌کند. هر چقدر که این معیار کمتر باشد نشان می‌دهد که تصاویر دو گروه مشابهت آماری بیشتری دارند و همچنین کمتر بودن این معیار با تصاویر پایه‌ی مجموعه داده همبستگی خوبی با کیفیت تصاویر تولیدی دارد.

جدول ۵-۲-۱ معیار FID را برای سه روش معرفی شده بر روی مجموعه‌ی $CIFAR-10$ نشان می‌دهد.

۵-۲-۲ معماری

معماری شبکه‌ی یو [۱۲] رایج‌ترین معماری استفاده شده برای مدل‌های پخش‌ی است. علت محبوبیت این معماری این است که یک ساختار متقارن دارد که بازنمایی هم ویژگی‌های سطح پایین و هم ویژگی‌های سطح بالا را یاد می‌گیرد و به نسبت، ساختاری ساده و قابل تغییر دارد. نوع شبکه‌ی یو به کار رفته در هر یک از مدل‌ها متفاوت است. معماری شبکه بسته به رزولوشن تصویر، وظیفه‌ی مورد نظر و نوع مدل متفاوت است و وجه مشترک همه‌ی آنها تنها در استفاده از انواعی از شبکه‌ی یو است. معماری‌های ذکر شده برای تصاویر با رزولوشن تا 64×64 هستند.

مدل تطبیق امتیاز با دینامیک لانژوین در شبکه‌ی امتیاز شرطی شده با نویز خود از یک شبکه‌ی یو با ۶ بلوک باقی مانده‌ای^۹ در هر سمت کدگذار و کدگشا استفاده می‌کند. مپناها از یک شبکه‌ی یو با ۸ بلوک باقی مانده‌ای در هر سمت کدگذار و کدگشا استفاده می‌کنند که جمعاً شامل ۱۵ بلوک باقی مانده‌ای است. مدل‌های پخش‌ی مبتنی بر معادله دیفرانسیل تصادفی نیز از یک شبکه‌ی یو با ۱۶ بلوک باقی مانده‌ای در هر سمت کدگذار و کدگشا استفاده می‌کنند.

^۸Frechet inception distance

^۹Residual block

مزیت اصلی شبکه‌ی یو به کارگیری اتصالات پرشی^{۱۰} است. اتصالات پرشی اتصالاتی هستند که از لایه‌های مختلف شبکه عبور می‌کنند و اجازه‌ی حفظ ویژگی‌های سطح پایین داده را می‌دهند که امری بسیار مهم برای تولید عکس‌هایی شبیه به واقعیت است.

انتخاب نوع معماری و تعداد پارامترها بسته به رویکرد و وظیفه‌ی مدنظر می‌تواند متفاوت باشد. به عنوان مثال اگر هدف تولید تصاویر با کیفیت بالا باشد ممکن است نیاز باشد که ساختار پیچیده‌تر با تعداد پارامتر بیشتر به کار گرفته شود و از طرفی اگر به دنبال سرعت در تولید تصاویر باشیم، نیاز داریم از یک شبکه‌ی ساده‌تر با تعداد پارامتر کمتر استفاده کنیم.

۵-۲-۳ ویژگی‌های کلی

با توجه به همه‌ی مواردی که پیش‌تر گفته شد و آنچه در مقالات مطالعه شده است، می‌توان چند نکته‌ی کلی درخصوص شباهت‌ها و تفاوت‌های سه مدل پخشی معرفی شده بیان کرد:

- هر سه مدل برای آموزش خود به منابع محاسباتی زیادی نیاز دارند و از نظر هزینه‌ی محاسباتی و زمان مورد نیاز برای یادگیری بسیار هزینه‌بر و زمان‌بر هستند.
- در مدل مپنا و مبتنی بر معادلات دیفرانسیل تصادفی، حساسیت زیادی نسبت به تعیین ابرپارامترها، تابع نویز زدا و معماری مدل وجود دارد و عملکرد مناسب آنها به انتخاب صحیح این پارامترها بستگی دارد.
- روش‌های مبتنی بر معادلات دیفرانسیل تصادفی انعطاف زیادی در تغییر فرآیند نویزافزایی دارند و می‌توان از آنها برای کاربردهایی که فرآیند تخریب نامشخص یا نیازمند تغییر است استفاده کرد.
- مدل تطبیق امتیاز با دینامیک لانژوین به دلیل معضل بعد^{۱۱} در داده‌های با بعد بالا عملکرد مناسبی ندارد.

¹⁰Skip connections

¹¹curse of dimensionality

فصل ششم

جمع‌بندی و کارهای آتی

۱-۶ جمع‌بندی

در این گزارش تلاش شد تا اطلاعات اولیه‌ای برای آشنایی و یافتن یک دیدگاه کلی در خصوص مدل‌های پخش ا ارائه شود. هدف کلی از این گزارش بررسی مقالات پایه و مبنای تئوری این مدل‌ها برای درک ساختار و سازو کار تولید داده بود و از پرداختن به حجم زیادی از کاربردها و بهبودها که با توجه به رشد روز افزون به کارگیری مدل‌های پخش رو به فزونی است خودداری شده است.

در فصل اول به معرفی مدل‌های پخش مبتنی بر امتیاز پرداختیم. این مدل‌ها در رویکرد اولیه خود تماماً پخش نیستند و تنها با استفاده از یک تطبیق امتیاز و دینامیک لانژوین قابل انجام هستند. اما رویکرد اولیه این روش‌ها با محدودیت‌هایی از جمله مشکل رویه‌ها و نواحی کمتر چگال توزیع داده روبرو بود که این محدودیت‌ها در داده‌های واقعی اکثر اوقات برقرار هستند. به همین دلیل در مدل ارائه شده که تطبیق امتیاز نویزدا با دینامیک لانژوین تابکاری شده نام دارد، به دنبال برطرف کردن این مشکل هستیم و این روش مانند تمام رویکردهای پخش به صورت مرحله به مرحله داده را با نویز تخریب و سپس در فرآیند معکوس نمونه‌هایی از داده را تولید می‌کند.

در فصل دوم به معرفی مپنا پرداختیم، مدلی که رایج‌ترین نوع مدل‌های پخش است و در کاربردهای تولید تصویر مطرح مانند $Dall - E$ به کار گرفته شده است. این مدل به صورت احتمالاتی و با یک زنجیره ی مارکوفی فرآیند پخش را مدل می‌کند. برای آموزش این شبکه رویکرد های مختلفی وجود دارد. در یک رویکرد می‌توانیم از تخمین میانگین پسین فرآیند رو به جلو برای آموزش شبکه استفاده کنیم و در رویکرد دوم می‌توانیم مستقیماً نویز اضافه شده به تصویر در یک گام از فرآیند پخش را تخمین بزنیم. می‌توان نشان داد که با تغییر پارامتر در رویکرد دوم و ثابت در نظر گرفتن واریانس این روش هم ارز روش تطبیق امتیاز با دینامیک لانژوین است.

در فصل سوم چهارچوبی کلی برای فرآیند پخش معرفی شد که قابلیت دربرگیری روش‌های پیش‌تر معرفی شده را دارد. در این روش نیز به صورت مرحله به مرحله داده را با نویز تخریب می‌کنیم اما به جای به کارگیری تعداد گام‌های محدود از تعداد گام نامتناهی پیوسته در یک بازه استفاده می‌کنیم. برای مدل کردن فرآیند پخش در چنین حالتی از یک معادله دیفرانسیل تصادفی استفاده می‌کنیم که معکوس چنین معادله‌ای را می‌توانیم تنها با داشتن تابع امتیاز بسازیم. برای ساخت نمونه از چنین معادله‌ای تنها کافی است از یک حل کننده استفاده کنیم و آن را گسسته کنیم. همچنین در این فصل نشان دادیم که هر دو روش پیش‌تر معرفی شده با تعیین معادلات متناسب با آنها می‌توانند در چهارچوب روش‌های مبتنی بر معادله دیفرانسیل تصادفی قرار گیرند.

در فصل چهارم پس از معرفی سه دسته‌ی اصلی مدل‌های پخش به مقایسه‌ی آنها هم با سایر مدل‌های مولد و هم با یکدیگر پرداختیم. نکته‌ی قابل تاکید در این فصل این است که برای مقایسه‌ی روش‌های معرفی شده تنها به مقالات به کار برده شده در سراسر گزارش اتکا شده است و با توجه به استقبال روزافزون از توسعه و به کارگیری این روش‌ها ممکن است به نتایج و بهبودهای به‌ترو یا رفع مشکلات مطرح شده پرداخته شده باشد که شرح آنها در این گزارش نیامده است.

۲-۶ کارهای آتی

با توجه به مطالعات انجام شده و با توجه به نوظهور بودن این دسته از مدل‌ها قابلیت بهبود و توسعه برای آنها در انواعی از جنبه‌ها وجود دارد که در ادامه آنها را نام می‌بریم.

- تئوری

از آنجاییکه مدل‌های پخشی اساساً بر پایه‌ی روابط ریاضی توسعه داده شده‌اند، امکان بهبود آنها از جنبه‌ی ریاضی و تئوری وجود دارد. به عنوان مثال ارائه‌ی فرمولاسیون جدیدی برای تخریب داده‌ها، استفاده از انواعی از حل‌کننده‌های معادله دیفرانسیل و بررسی تأثیرات آنها و همچنین به کارگیری مباحثی از فیزیک ترمودینامیک از جمله روش‌های قابل اعمال هستند.

- معماری

معماری فعلی تمام مدل‌های پخشی از هسته‌ی شبکه‌ی یو بهره می‌برد. با وجود عملکرد فوق‌العاده‌ی این شبکه‌ها همچنان مشکلاتی در استفاده از آنها در بعضی از کاربردهای خاص وجود دارد و علاوه بر این تعداد پارامترهای این شبکه‌ها آنها را از لحاظ محاسباتی بسیار سنگین می‌کند. پیاده‌سازی و استفاده از سایر شبکه‌ها برای شبکه‌ی نوین زدا می‌تواند پیشرفت قابل توجهی در به کارگیری این مدل‌ها در همه‌ی کاربردها فراهم آورد.

- داده

آموزش و توسعه‌ی مدل‌های پخشی با میزان داده‌ی کارا یک چالش باز در بحث توسعه‌ی مدل‌های پخشی است. این مدل‌ها برای تولید داده‌ی با کیفیت نیاز به دریافت داده‌ی با کیفیت دارند و از این رو کار روی یادگیری محدود و یادگیری با داده‌ی کم کیفیت می‌تواند یک رویکرد قابل پیگیری باشد.

- کاربرد

همانطور که در این گزارش نیز مشاهده شد، اکثر توسعه‌های مدل‌های پخشی بر روی داده‌های تصویری انجام شده است. اما تحقیقات متعددی روی پیاده‌سازی و تطبیق این مدل‌ها در سایر حوزه‌های کاربرد مانند صوت، متن، یادگیری تقویتی و گراف نیز در حال انجام است. تطبیق مدل‌های پخشی با انواع اشکال داده‌گونه‌ای دیگر از مسائلی است که می‌تواند مورد بررسی و تغییر قرار گیرد.

کتاب نامه

- [1] Anderson, Brian DO. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- [2] Feller, William. Retracted chapter: On the theory of stochastic processes, with particular reference to applications. In *Selected Papers I*, pages 769–798. Springer, 2015.
- [3] Goodfellow, Ian, Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair, Sherjil, Courville, Aaron, and Bengio, Yoshua. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- [4] Henighan, Tom, Kaplan, Jared, Katz, Mor, Chen, Mark, Hesse, Christopher, Jackson, Jacob, Jun, Heewoo, Brown, Tom B, Dhariwal, Prafulla, Gray, Scott, et al. Scaling laws for autoregressive generative modeling. *arXiv preprint arXiv:2010.14701*, 2020.
- [5] Heusel, Martin, Ramsauer, Hubert, Unterthiner, Thomas, Nessler, Bernhard, and Hochreiter, Sepp. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [6] Ho, Jonathan, Jain, Ajay, and Abbeel, Pieter. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

- [7] Hyvärinen, Aapo and Dayan, Peter. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- [8] Kingma, Diederik P, Welling, Max, et al. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019.
- [9] Kreis, Karsten, Gao, Ruiqi, and Vahdat, Arash. Denoising diffusion-based generative modeling: Foundations and applications, 2022.
- [10] Nichol, Alexander Quinn and Dhariwal, Prafulla. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- [11] Razavi, Ali, Van den Oord, Aaron, and Vinyals, Oriol. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32, 2019.
- [12] Ronneberger, Olaf, Fischer, Philipp, and Brox, Thomas. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015.
- [13] Salimans, Tim, Karpathy, Andrej, Chen, Xi, and Kingma, Diederik P. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. *arXiv preprint arXiv:1701.05517*, 2017.
- [14] Song, Yang and Ermon, Stefano. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- [15] Song, Yang and Ermon, Stefano. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.

-
- [16] Song, Yang, Sohl-Dickstein, Jascha, Kingma, Diederik P, Kumar, Abhishek, Ermon, Stefano, and Poole, Ben. Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456, 2020.
- [17] Vincent, Pascal. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
- [18] Wang, Zhendong, Zheng, Huangjie, He, Pengcheng, Chen, Weizhu, and Zhou, Mingyuan. Diffusion-gan: Training gans with diffusion. arXiv preprint arXiv:2206.02262, 2022.
- [19] Xiao, Zhisheng, Kreis, Karsten, and Vahdat, Arash. Tackling the generative learning trilemma with denoising diffusion gans. arXiv preprint arXiv:2112.07804, 2021.
- [20] Yang, Ling, Zhang, Zhilong, Song, Yang, Hong, Shenda, Xu, Runsheng, Zhao, Yue, Shao, Yingxia, Zhang, Wentao, Cui, Bin, and Yang, Ming-Hsuan. Diffusion models: A comprehensive survey of methods and applications. arXiv preprint arXiv:2209.00796, 2022.

واژه‌نامه‌ی فارسی به انگلیسی

Denoising score . . . تطبیق امتیاز نویزدا . . . matching	۱
Discrimination تمایزگر	Skip connections اتصالات پرشی
Prior distribution توزیع اولیه	Variance exploding انفجار واریانس
Generator تولید کننده	ب
ج	Residual block بلوک باقی‌مانده‌ای
Normalizing flows جریان‌های نرمال‌ساز	ت
ح	Score function تابع امتیاز
Evidence lower bound حد پایین مشاهده‌شده	Simulated شبیه‌سازی شده annealing
Variational bound حد تغییراتی	Mixing ترکیب
Euler-maruyama حل کننده‌ی اوایلر ماریاما solver	Image inpainting ترمیم تصویر
خ	Embedding تعبیه
Variational . . . خودکدگذارهای تغییراتی	Unique identifiability تعیین یکتا
autoencoders	Score matching تطبیق امتیاز

Deterministic process . . . فرآیند قطعی	خودکدگذارهای تغییراتی سلسله مراتبی
Forward process فرآیند رو به جلو	Hierarchical variational autoencoders
Standard wiener فرآیند وینر استاندارد	د
process	Likelihood درست‌نمایی
Support فضای پشتیبان	Langevin dynamics دینامیک لانژوین
ق	Annealed دینامیک لانژوین تابکاری شده
Semantic قطعه‌بندی معنایی	langevin dynamics
segmentation	KL-divergence دیورجنس KL
ک	ر
Decoder کدگشا	Colorization رنگ‌آمیزی
Encoder کدگذار	ش
ل	Generative شبکه‌های مولد تقابلی
Log-likelihood لگاریتم درست‌نمایی	adversarial networks
م	شبکه‌ی یو
Mode مد	U-net ض
مدل پخشی نويززدای احتمالاتی	Diffusion coefficient ضریب پخش
Denoising diffusion probabilistic models	Drift coefficient ضریب رانش
مدل‌های مولد مبتنی بر امتیاز	غ
generative models	Intractable غیرقابل محاسبه
مدل‌های متغیر پنهان	ف
models	Frechet فاصله‌ی اکتسابی فرچت
	inception distance

Diffusion models مدل پخشی

Forward trajectory مسیر روبه جلو

Generative model یادگیری مولد
learning trilemma

Ordinary معادلات دیفرانسیل ساده
differential equations

Posterior mean میانگین پسین

Vector field میدان برداری

ن

Anomaly ناهنجاری

Diffusion rate نرخ پخش

Saddle point نقطه‌ی زینی

Importance نمونه‌برداری اهمیت
sampling

ه

Kernel هسته

Binomial kernel . . . هسته‌ی دوجمله‌ای

Transition kernel هسته‌ی گذر

واژه‌نامه‌ی انگلیسی به فارسی

A	ضریب پخش Diffusion coefficient
Anomaly	مدل پخشی Diffusion models
Annealed	نرخ پخش Diffusion rate
دینامیک لانژوین تابکاری شده . langevin dynamics	تمایزگر Discriminator
B	ضریب رانش Drift coefficient
Binomial kernel . . .	E
C	کدگذار Encoder
Colorization	تعبیه Embedding
D	حد پایین مشاهده Evidence lower bound
Decoder	حل‌کننده‌ی اویلر ماریاما Euler-maruyama solver
Denoising score . . .	F
تطبیق امتیاز نویزدا . . . matching	فرآیند رو به جلو Forward process
مدل پخشی نویززدای احتمالاتی Denoising diffusion probabilistic models	مسیر روبه جلو Forward trajectory
Deterministic process . . .	

Frechet فاصله‌ی اکتسابی فرچت inception distance	Latent variable . . مدل‌های متغیر پنهان models
G	Likelihood درست‌نمایی
Generative شبکه‌های مولد تقابلی adversarial networks	Log-likelihood . . . لگاریتم درست‌نمایی
Generative مشکل سه‌گانه‌ی یادگیری مولد learning trilemma	M
Generator تولید کننده	Mixing ترکیب
H	Mode مد
خودکدگذارهای تغییراتی سلسله مراتبی Hierarchical variational autoencoders	N
I	Normalizing flows جریان‌های نرمال‌ساز
Image inpainting ترمیم تصویر	O
Importance نمونه‌برداری اهمیت sampling	Ordinary معادلات دیفرانسیل ساده differential equations
Intractable غیرقابل محاسبه	P
K	Posterior mean میانگین پسین
Kernel هسته	Prior distribution توزیع اولیه
KL-divergence دیورجنس KL	R
L	Residual block بلوک باقی‌مانده‌ای
Langevin dynamics . . دینامیک لانژوین	S
	Saddle point نقطه‌ی زینی
	Score-based مدل‌های مولد مبتنی بر امتیاز generative models
	Score function تابع امتیاز

Score matching تطبیق امتیاز	Transition kernel هسته‌ی گذر
Semantic قطعه‌بندی معنایی segmentation	U
Semi-supervised . یادگیری نیمه نظارتی . learning	U-net شبکه‌ی یو
Simulated تابکاری شبیه‌سازی شده annealing	Unique identifiability تعیین یکتا
Skip connections اتصالات پرشی	V
Standard wiener . . فرآیند وینر استاندارد process	Variance exploding انفجار واریانس
Support فضای پشتیبان	Variational . . . خودکدگذارهای تغییراتی autoencoders
T	Variational bound حد تغییراتی
	Vector field میدان برداری