

① ابتدا فرضیات و املاعاتی که داریم را می نویسیم

$P_t(x) \rightarrow$ توزیع x در زمان t

$P_{ot}(x'|x) \rightarrow$ احتمال رفتن از x_0 به x_t

forward:

$$dx = f(x, t) dt + g(t) d\bar{w}$$

Inverse:

$$dx = [f(x, t) - g(t)^T \nabla_x \log P_t(x)] dt + g(t) d\bar{w}$$

$$J_{SM}(\theta; \lambda(\cdot)) = \frac{1}{T} \int_0^T E_{P_t(x)} [\lambda(t) \|\nabla_x \log P_t(x) - S_\theta(x, t)\|_2^2] dt$$

$$J_{DSM}(\theta; \lambda(\cdot)) = \frac{1}{T} \int_0^T E_{P(x) P_{ot}(x'|x)} [\lambda(t) \|\nabla_x \log P_{ot}(x'|x) - S_\theta(x', t)\|_2^2] dt$$

انتظار داریم در نهایت: $P_T(x) \approx \pi(x)$

پس از انتشار از $\hat{x}_0(T) \sim \pi$ شروع می کنیم و به $\hat{x}_0(0) \sim P_0^{SDE}$ می رسیم

$$d\hat{x} = [f(\hat{x}, t) - g(t)^T S_\theta(\hat{x}, t)] dt + g(t) d\bar{w}, \quad \hat{x}_0(T) \sim \pi$$

$$\frac{dx}{dt} = f(x, t) - \frac{1}{T} g(t)^T \nabla_x \log P_t(x) \leftarrow \text{ODE}$$

$$\frac{d\tilde{x}}{dt} = f(\tilde{x}, t) - \frac{1}{T} g(t)^T S_\theta(\tilde{x}, t) \leftarrow \text{تقریب}$$

Subject

Date : Year:

Month:

Day:

$$D_{KL}(P \parallel P_{\theta}^{SDE}) \leq J_{SM}(\theta; g(\cdot)) + D_{KL}(P_T \parallel \pi) \quad : \text{نیمه}$$

$$H(P) = H(P_T(m)) - \frac{1}{T} \int_0^T E_{P_t(m)} \left[\gamma \nabla \cdot f(m_t) + g(t)' \nabla_m \log P_t(m) \right] dt$$

5

10

15

20

الف) KL بین p و p_{θ}^{SDE} به ما نشان می‌دهد که مدل ما چقدر از توزیع داده واقعی فاصله دارد.

ایده آل مان هر این است که $p_{\theta}^{SDE} \approx p$. یعنی واگرایی مان کمینه شود.

J_{SM} هزینه اصلی تقنین score است و هرچه آنرا کوچکتر کنیم یعنی بهتر score

را یاد گرفته ایم.

KL بین p_T و p در زمان نهایی T ، توزیع نویز شده p_T چقدر با توزیع هدف p تفاوت دارد.

10

این مکان می‌گوید که: اگر می‌خواهید کل واگرایی KL بین p و p_{θ}^{SDE} کم شود، کافی است که:

① تابع score را خوب یاد گرفته باشید تا J_{SM} کوچک شود، ② مطمئن باشید که فرایند

پس رو در زمان T به توزیع p_T که از نظر KL فاصله جتنانی با p نداشته باشد، می‌رسد. 15

در مقاله هر داریم که اگر p را به fixed prior در نظر بگیریم:

$$-E_{p(w)} [\log p_{\theta}^{SDE}] \leq J_{SM}(\theta; g(\cdot)^2) + C_1$$

20

حنگامی که اسکور به خوبی یاد گرفته شود، فرایند معکوس می‌تواند نمونه‌ها را بصیر واقعی و

قابل قبولی بسازد، پس از نظر کیفیت بهبود داریم، از طرفی کوچک شدن J_{SM} باعث

تردیک شدن P و P_{θ}^{SDE} در فضای توزیع‌های می‌شود. بنابراین اگر با معیارهایی مثل

KL , FID یا حتی تست شرایلی $likelihood$ ارزیابی کنیم، بهبود در J_{SM}

باعث بهبود در این معیارها نیز خواهد شد.

← J_{SM} و در کل این نگران بالا از نظر کیفی کمی حساس عمل می‌کند. وقتی این عبارت کوچک

شود، هم از بعد بسیر و هم از بعد آماری، احتمالاتی نمونه‌ها را تردیک به واقعیت و توزیع

واقعی خواص دست است.

(ب) حکم:

$$-E_{P(n)}[\log P_{\theta}^{SDE}(n)] \leq -E_{P_T(n)}[\log \pi(n)] + \frac{1}{T} \int_0^T E_{P_{\theta,t}(n'|n)} P(n) [g(t)^2$$

$$\|s_{\theta}(n',t) - \nabla_{n'} \log P_{\theta,t}(n'|n)\|_2^2 - g(t)^2 \|\nabla_{n'} \log P_{\theta,t}(n'|n)\|_2^2 - 2 \nabla \cdot f(n',t)] dt$$

$$D_{KL}(P \parallel P_{\theta}^{SDE}) = E_{P(n)} \left[\log \frac{P(n)}{P_{\theta}^{SDE}(n)} \right] = E_{P(n)}[\log P(n)] - E_{P(n)}[\log P_{\theta}^{SDE}(n)]$$

$$= -H(P) - E_{P(n)}[\log P_{\theta}^{SDE}(n)] \Rightarrow$$

$$-E_{P(n)}[\log P_{\theta}^{SDE}(n)] \leq H(P) + J_{SM}(\theta; g(\cdot)) + D_{KL}(P \parallel \pi)$$

$$D_{KL}(P_T \parallel \pi) = E_{P_T(n)} [\log P_T(n)] - E_{P_T(n)} [\log \pi(n)] =$$

$$-H(P_T(n)) - E_{P_T(n)} [\log \pi(n)]$$

$$5 \rightarrow H(P) + D_{KL}(P_T \parallel \pi) = \cancel{H(P_T(n))} - \frac{1}{T} \int_0^T E_{P_t(n)} [\gamma \nabla \cdot \dot{F}(n,t) + g(t)^T$$

$$\|\nabla_n \log P_t(n)\|_r^2] dt$$

$$- \cancel{H(P_T(n))} - E_{P_T(n)} [\log \pi(n)]$$

$$10 \rightarrow -E_{P(n)} [\log P_{\theta}^{SDE}(n)] \leq -E_{P_T(n)} [\log \pi(n)] - \frac{1}{T} \int_0^T E_{P_t(n)} [\gamma \nabla \cdot \dot{F}(n,t)$$

$$+ g(t)^T \|\nabla_n \log P_t(n)\|_r^2] dt + \frac{1}{T} \int_0^T E_{P_t(n)} [g(t)^T \|\nabla_n \log P_t(n) - S_{\theta}(n,t)\|_r^2] dt$$

$$\rightarrow -E_{P(n)} [\log P_{\theta}^{SDE}(n)] \leq -E_{P_T(n)} [\log \pi(n)] + \frac{1}{T} \int_0^T E_{P_t(n)} [g(t)^T \|S_{\theta}(n,t)$$

$$15 \quad -\nabla_n \log P_t(n)\|_r - g(t)^T \|\nabla_n \log P_t(n)\|_r - \gamma \nabla \cdot \dot{F}(n,t)] dt \quad ***$$

$$\int_0^T E_{P_{\theta t}(x'|x)P(n)} [g(t)^T \|S_{\theta}(n',t) - \nabla_{n'} \log P_{\theta t}(n'|n)\|_r - g(t)^T \|\nabla_{n'} \log$$

$$20 \quad P_{\theta t}(n'|n)\|_r^2] dt = \int_0^T E_{P_{\theta t}(n'|n)P(n)} [g(t)^T (\|S_{\theta}(n',t)\|_2^2 - \gamma \nabla_{n'} \log P_{\theta t}(n'|n)^T$$

$$S_{\theta}(n',t))\|_r^2] dt = \int_0^T \int \int P_{\theta t}(n'|n) P(n) g(t)^T (\|S_{\theta}(n',t)\|_r^2 - \gamma \nabla_{n'} \log P_{\theta t}(n'|n)^T$$

$$S_{\theta}(n',t)) dx' dn dt = \int_0^T g(t)^T \int_{x'} \|S_{\theta}(n',t)\|_2^2 \underbrace{\int_n P(n) P_{\theta t}(n'|n) dn}_{P_t(x')} dn' dt$$

parsian

$$- \gamma \int_0^T g(t)^r \int_{\mathbf{x}'} s_{\theta}(\mathbf{x}', t)^T \int_{\mathbf{n}} p(\mathbf{n}) p_t(\mathbf{n}'|\mathbf{n}) \nabla_{\mathbf{n}'} \log p_t(\mathbf{n}'|\mathbf{n}) d\mathbf{n} d\mathbf{n}' dt$$

$$= \int_0^T g(t)^r E_{p_t(\mathbf{n}')} [\|s_{\theta}(\mathbf{x}', t)\|_2^2] dt - \gamma \int_0^T g(t)^r \int_{\mathbf{x}'} s_{\theta}(\mathbf{n}', t)^T \int_{\mathbf{n}} p(\mathbf{n}) p_t(\mathbf{n}'|\mathbf{n})$$

$$\frac{\nabla_{\mathbf{n}'} p_t(\mathbf{n}'|\mathbf{n})}{p_t(\mathbf{n}'|\mathbf{n})} d\mathbf{n} d\mathbf{n}' dt = \frac{\nabla_{\mathbf{n}'} \int_{\mathbf{n}} p(\mathbf{n}) p_t(\mathbf{n}'|\mathbf{n}) d\mathbf{n}}{\nabla_{\mathbf{n}'} p_t(\mathbf{n}')} =$$

$$\int_0^T g(t)^r E_{p_t(\mathbf{n}')} [\|s_{\theta}(\mathbf{n}', t)\|_r^r] dt - \gamma \int_0^T g(t)^r \int_{\mathbf{x}'} s_{\theta}(\mathbf{n}', t)^T \nabla_{\mathbf{n}'} p_t(\mathbf{n}') d\mathbf{n}' dt$$

$p_t(\mathbf{n}') \nabla_{\mathbf{n}'} \log p_t(\mathbf{n}') \rightarrow E_{p_t(\mathbf{n}')} [---]$

$$= \int_0^T g(t)^r E_{p_t(\mathbf{n}')} [\|s_{\theta}(\mathbf{n}', t)\|_r^r - \gamma s_{\theta}(\mathbf{n}', t)^T \nabla_{\mathbf{n}'} \log p_t(\mathbf{n}')] dt \quad 10$$

$$= \int_0^T \overset{g(t)^r}{\sqrt{E_{p_t(\mathbf{n}')}}} [\|s_{\theta}(\mathbf{n}', t) - \nabla_{\mathbf{n}'} \log p_t(\mathbf{n}')\|_r^r - \|\nabla_{\mathbf{n}'} \log p_t(\mathbf{n}')\|_r^r] dt$$

$$= \int_0^T E_{p_t(\mathbf{n}')} [g(t)^r \|s_{\theta}(\mathbf{n}', t) - \nabla_{\mathbf{n}'} \log p_t(\mathbf{n}')\|_r^r - g(t)^r \|\nabla_{\mathbf{n}'} \log p_t(\mathbf{n}')\|_r^r] dt \quad 15$$

از طرفی داریم:

$$\int_0^T E_{p_t(\mathbf{n}'|\mathbf{n})} p(\mathbf{n}) [-\gamma \nabla \cdot \mathbf{f}(\mathbf{n}', t)] dt = \int_0^T -\gamma \nabla \cdot \mathbf{f}(\mathbf{n}', t) \left(\int p_t(\mathbf{n}'|\mathbf{n}) p(\mathbf{n}) d\mathbf{n} \right) d\mathbf{n}' dt$$

$$= \int_0^T E_{p_t(\mathbf{n}')} [-\gamma \nabla \cdot \mathbf{f}(\mathbf{n}', t)] dt \Rightarrow$$

$$\frac{1}{\gamma} \int_0^T E_{p_t(\mathbf{n}'|\mathbf{n})} p(\mathbf{n}) [g(t)^r \|s_{\theta}(\mathbf{n}', t) - \nabla_{\mathbf{n}'} \log p_t(\mathbf{n}'|\mathbf{n})\|_r^r - g(t)^r \|\nabla_{\mathbf{n}'} \log p_t(\mathbf{n}')\|_r^r - \gamma \nabla \cdot \mathbf{f}(\mathbf{n}', t)] dt \quad \text{parsian}$$

این عبارت هر معادل *** است.

پس حکم بخش ب اثبات می شود.

ب) می دانیم که حکم این بخش همان حکم بخش ب است که یک Expected روی $P(n)$

از آن حذف شده است، حالا اگر یک SDE فیلتر را در نظر بگیریم، یعنی تمامی

transition kernel های $P_{0t}(n'|n)$ فیلتر باشند، قضیه قسمت ب برابر تمام توزیع های

$P(n)$ برقرار است. حالا با استفاده از اثبات با تناقض بسمت می آید که می توانیم $E_P(n)$ را

از دو طرف رابطه حذف کنیم و به این برسیم:

$$-\log P_0^{SDE}(n) \leq L_0^{DSM}(n)$$

ابتدا عبارت بخش ب را بازنویسی می کنیم:

$$E_{P(n)}[-\log P_0^{SDE}(n)] \leq -E_{P_T(n)}[\log \pi(n)] + \frac{1}{T} \int_0^T E_{P_{0t}(n'|n)} P(n) [$$

$$g(t)^2 \|s_0(n', t) - \nabla_{n'} \log P_{0t}(n'|n)\|_2^2 - g(t)^2 \|\nabla_{n'} \log P_{0t}(n'|n)\|_2^2 - 2 \nabla f(n', t)] dt$$

عبارت سمت راست راست را به این صورت می نویسیم:

$$\int_0^T \int_{\mathbf{x}'} \int_{\mathbf{n}} P_{\theta t}(\mathbf{n}'|\mathbf{n}) P(\mathbf{n}) [g(t)^r \| S_{\theta}(\mathbf{n}', t) - \nabla_{\mathbf{n}'} \log P_{\theta t}(\mathbf{n}'|\mathbf{n}) \|_r^r - g(t)^r \| \nabla_{\mathbf{n}'} \log P_{\theta t}(\mathbf{n}') \|_r^r - \gamma \nabla \cdot \mathbf{f}(\mathbf{n}', t)] d\mathbf{n} d\mathbf{n}' dt = \int_{\mathbf{n}} P(\mathbf{n}) \int_0^T \int_{\mathbf{n}'} P_{\theta t}(\mathbf{n}'|\mathbf{n}) [g(t)^r \| S_{\theta}(\mathbf{n}', t) - \nabla_{\mathbf{n}'} \log P_{\theta t}(\mathbf{n}'|\mathbf{n}) \|_r^r - g(t)^r \| \nabla_{\mathbf{n}'} \log P_{\theta t}(\mathbf{n}') \|_r^r - \gamma \nabla \cdot \mathbf{f}(\mathbf{n}', t)] d\mathbf{x}' dt d\mathbf{n} =$$

$$E_P(\mathbf{n}) \left[\int_0^T P_{\theta t}(\mathbf{n}'|\mathbf{n}) [g(t)^r \| S_{\theta}(\mathbf{n}', t) - \nabla_{\mathbf{n}'} \log P_{\theta t}(\mathbf{n}'|\mathbf{n}) \|_r^r - g(t)^r \| \nabla_{\mathbf{n}'} \log P_{\theta t}(\mathbf{n}') \|_r^r - \gamma \nabla \cdot \mathbf{f}(\mathbf{n}', t)] d\mathbf{x}' dt \right] =$$

$$E_P(\mathbf{n}) \left[\int_0^T E_{P_{\theta t}(\mathbf{n}'|\mathbf{n})} [g(t)^r \| S_{\theta}(\mathbf{n}', t) - \nabla_{\mathbf{n}'} \log P_{\theta t}(\mathbf{n}'|\mathbf{n}) \|_r^r - g(t)^r \| \nabla_{\mathbf{n}'} \log P_{\theta t}(\mathbf{n}') \|_r^r - \gamma \nabla \cdot \mathbf{f}(\mathbf{n}', t)] dt \right]$$

از طرفی برابر $E_{P_T(\mathbf{n})} [\log \pi(\mathbf{n})]$ هم داریم که

$$E_{P_T(\mathbf{n})} [\log \pi(\mathbf{n})] = \int P_T(\mathbf{n}) \log \pi(\mathbf{n}) d\mathbf{n} = \int \log \pi(\mathbf{n}) \int P(\mathbf{n}') P_{\theta T}(\mathbf{n}|\mathbf{n}') d\mathbf{n}' d\mathbf{n}$$

$$= \iint P(\mathbf{n}') P_{\theta T}(\mathbf{n}|\mathbf{n}') \log \pi(\mathbf{n}) d\mathbf{n} d\mathbf{n}' = E_{P(\mathbf{n}')} \left[\int P_{\theta T}(\mathbf{n}|\mathbf{n}') \log \pi(\mathbf{n}) d\mathbf{n} \right]$$

$$= E_{P(\mathbf{n}')} [E_{P_{\theta T}(\mathbf{n}|\mathbf{n}')} [\log \pi(\mathbf{n})]] = E_{P(\mathbf{n})} [E_{P_{\theta T}(\mathbf{n}'|\mathbf{n})} [\log \pi(\mathbf{n}')]]$$

پس داریم که

$$E_{P(n)} [-\log P_{\theta}^{SDE}(n)] \leq E_{P(n)} [-E_{P_{\theta T}(n'|n)} [\log \pi(n')]]$$

$$+ E_{P(n)} \left[\frac{1}{T} \int_0^T E_{P_{\theta t}(n'|n)} [g(t)^T \| \nabla_{n'} \log P_{\theta t}(n'|n) \|^2] dt \right]$$

$$- \frac{1}{T} \int_0^T E_{P_{\theta t}(n'|n)} [g(t)^T \| \nabla_{n'} \log P_{\theta t}(n'|n) \|^2 + 2 \nabla_{n'} \cdot f(n', t)] dt$$

که یعنی داریم:

$$* E_{P(n)} [-\log P_{\theta}^{SDE}(n)] \leq E_{P(n)} [\mathcal{L}_{\theta}^{DSM}(n)]$$

دقت شود که اگر همه را به یک $E_{P(n)}$ ببریم رابطه بالا بدست می آید.

حالا می دانیم که به ازای هر توزیع $P(n)$ ای این رابطه برقرار است به ازای ساختار

فیلکس SDE، اثبات با تناقض.

فرض کنیم که نقطه n^* ای داریم که در آن $\mathcal{L}_{\theta}^{DSM}(n^*) > -\log P_{\theta}^{SDE}(n^*)$ باشد.

آنگاه خواهیم داشت که: به ازای هر توزیع $P(n)$ رابطه $*$ برقرار است، پس

برای توزیع $P(n)$ به صورتی که $P(n^*) = 1$ ، در این نقطه صفر باید این رابطه

برقرار باشد که به وضوح برقرار نیست ← همواره باید داشته باشیم که:

$$-\log P_{\theta}^{SDE}(n) \leq \mathcal{L}_{\theta}^{DSM}(n) \rightarrow$$

بعضی ب هم اثبات می شود.

الف) در سری های زمانی ممکن است به سری از دیتاها یک شده باشند یا دچار خطا شده باشند، از این رو

ما دچار مشکل می شویم در یادگیری آنها، بدین منظور از مدل های CSDI استفاده می کنیم، ایده آن

هم این است که به جای تولید داده های جدید (مانند تقاضا)، از این مدل ها برابر بازسازی مقادیر گمشده

با استفاده از اطلاعات موجود در داده های مشاهده شده بهره بگیرند.

روش آموزش:

در این روش، داده های مشاهده شده به دو بخش تقسیم می شوند. هجین به بخشی را برابر

جایگزینی به بخشی را برابر اطلاعات شرطی در نظر می گیریم. هجین به mask تعریف می کنیم

که برابر دیتار از دست رفته ی $mask = 0$ است وگرنه 1 است. گفتیم که بخشی را به عنوان

جایگزینی target: x_0^{tel} و بخشی به عنوان اطلاعات شرطی: x_0^c استفاده می شود.

در واقع از x_0^c به عنوان راهنما استفاده می کنیم، احتمالات به شرط آنرا حساب می کنیم.

سعی می کنیم x_0^{tel} را پس بینی کنیم. در گام t فرایند هر نسخه ای از x_0^{tel} با نویز ترکیب

$$x_t^{tel} = \sqrt{\alpha_t} x_0^{tel} + \sqrt{1 - \alpha_t} \epsilon$$

شده است.

برای آموزش مدل، مدل باید یاد بگیرد نویز x_t^{tot} را با توجه به x_0^c پیش‌بینی کند ←

تابع لاسی با: $\mathcal{L}(\theta) = \mathbb{E} \|\epsilon - \epsilon_\theta(x_t^{tot}, t | x_0^c)\|^2$ تعریف می‌کنیم و سعی

می‌کنیم آنرا کمینه کنیم.

مزایای معکوس هر به این صورت است که: در زمان استنتاج، مدل از نویز تصادفی شروع

می‌کند و تکرار حذف نویز، مقادیر کم‌شده را بازسازی می‌کند.

ب) استراتژی‌ها عبارتند از:

① رندم: در هر داده‌ها مشاهده شده به صورت تصادفی به عنوان هدف انتخاب می‌شوند

مزایای انعطاف پذیری برابر نسبت به هر کم‌شده است، معایب این است که ممکن است

الگوها را به‌خاطر یافته را یاد نکند.

② تارگتی: از الگوها کم‌شده در داده‌ها آموزشی استفاده می‌کند، این استراتژی برابر داده‌هایی

با الگوها را به‌خاطر یافته (مانند مقادیر کم‌شده/موتالی) مناسب است.

مزایا: برابر داده‌ها تکرار شوند مناسب است. معایب: خطر $overfitting$ ممکن است

۳) ترکیبی: ترکیبی از استراتژی‌های رندوم و تاریخی است تا هر از تغییرپذیر و هر از

الگوها را اختار یافته بهره ببرد. این استراتژی ترکیبی تعادل مناسبی بین

تعمیرپذیر و یادگیری الگوها را اختار یافته ایجاد می‌کند و احتمالاً بهبود بیشتر در عملکرد

نهایی مدل دارد.

۴) الگو آزمون: زمانی استفاده می‌شود که الگوها را کمینه در داده‌های آزمون از قبل

10

مستخرج باشند (مثلاً پیش‌بینی آینده) این استراتژی هر فقط زمانی کاربرد دارد که

الگوها را کمینه / آزمون مستخرج باشند.

15

ج) معیارها عبارتند از:

CRPS (Continuous Ranked Probability Score):

هدفش ارزیابی دقت پیش‌بینی‌های احتمالاتی است. هرچه کمتر CRPS کمتر باشد.

20

یعنی عملکرد مدل بهتر است.

CSDI، در حدود ۴۵-۴۰٪ مقدار CRPS را بهبود داده است نسبت به روش‌های

احتمالاتی موجود.

MAE (Mean absolute Error).

هدفش ارزیابی دقت پیش بینی هار قطعی است.

مکانیزم آن به این صورت است که میانگین اختلاف قدر مطلق بین مقادیر واقعی

و پیش بینی شده را محاسبه می کند.

CSDI مقدار MAE را بین ۵-۲۰٪ نسبت به state-of-the-art deterministic method

کاهش داده است.

با پیاده سازی