

Most times Are used multiple times (MIASM)

Sylvain Le Corff, Sorbonne Université, LPSM

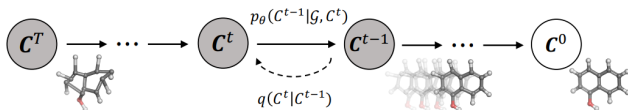
Based on joint works with
Stanislas Strasman, Claire Boyer, Vincent Lemaire (LPSM)
Antonio Ocello, Yazid Janati, Gabriel Cardoso, Eric Moulines (CMAP)

Generative modeling

Assumption: (X^1, \dots, X^N) in \mathbb{R}^{d_x} are samples from some **unknown** distribution π_{data} .

What: generate synthetic instances of a target distribution π_{data}

Probabilistic model for generating molecular conformations (GeoDiff, Xu et al., 2022)

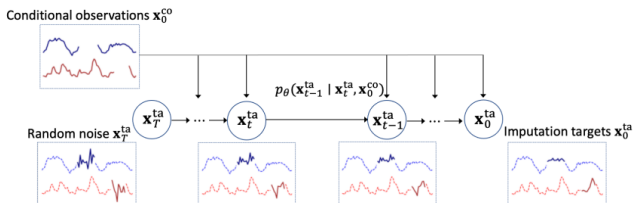


GEOM (37 million molecular conformations annotated by energy):
generates new structures + chemical toolkit to calculate conformation energy.



Generative modeling

Probabilistic time series imputation (CSDI, Tashiro et al., 2021)



Healthcare dataset in PhysioNet Challenge 2012 (4000 clinical time series with 35 variables for 48 hours from intensive care unit).

Synthetic Data Generation for Privacy and Security (TabDDPM, Kotelnikov et al., 2023), etc.

most Times ARE used multiple times (TARDigrade)



Generative modeling

PierrE pictURes (PEUR) dataset (3 pictures, dataset being built slowly but surely).

Improving corrupted data



Generative modeling

PierrE pictURes (PEUR) dataset (3 pictures, dataset being built slowly but surely).

Improving corrupted data



Generative modeling

PierrE pictURes (PEUR) dataset (3 pictures, dataset being built slowly but surely).

Improving corrupted data



most times are used multiple times (juste pour le plaisir)

Generative modeling

① Estimate π_{data} with a **parametric** probability distribution p_{θ} .

1. Choose a **suitable parametric form** for p_{θ} .

↪ p_{θ} is parameterized using a **Neural Network**.

↪ $p_{\theta} \geq 0$, $\int p_{\theta} = 1 \rightarrow$ **constraints** on the modelization.

2. Train p_{θ} **to approximate** π_{data} using the samples
 $(X^1, \dots, X^N) \sim \pi_{\text{data}}$: $\mathcal{L}(\theta) = \sum_{i=1}^N -\log p_{\theta}(X^i)$.

↪ Minimize $\mathcal{L}(\theta)$ or an **upperbound** \rightarrow find optimal parameter θ_* .



mOst times Used aRe multiple timeS (OURS)

Controlled generation

② Perform controlled generation using p_{θ_*} .

↪ Target distribution: weight p_{θ_*} with a function $x \mapsto g(x)$

$$\phi(dx) = \frac{g(x)p_{\theta_*}(dx)}{\int g(z)p_{\theta_*}(dz)},$$

↪ Posterior sampling: $g(x) = p(y|x)$.

mosT timEs aRe used Multiple TimES (TERMITES)



Score-Based Generative Models (SGMs)

What: generate synthetic instances of a target distribution π_{data}

Why: challenges in modeling the complexity of real data, preventing conventional parametric modeling or traditional maximum likelihood methods.

Creating noise from data is easy; creating data from noise is generative modeling. (Song et al., *Score-Based Generative Modeling through Stochastic Differential Equations*)

Who: SGMs address this by

1. (forward phase) introducing progressively noise into the samples,
2. (backward phase) reversing the noising dynamics, with the help of a score function usually learned using deep neural networks.



mOst times aRe usEd multiple times (OREortyx)

Noise scheduling tuning

Practitioners' corner: time-inhomogeneous SGMs - the noise schedule has central role, as exhibited in numerical experiments.

[Chen et al., 2023]: performance of SGMs relies on the chosen noise scheduling and optimal strategy varies depending on the task, *e.g.*, image sizes.

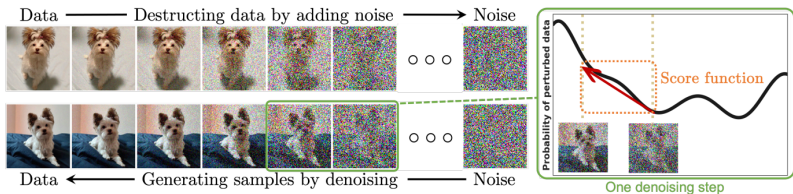


Figure: From Yang et al., 2023

MOSt times are used multiple Times (MOTh)



Estimating π_{data} with diffusion models

Diffusion models: the forward process - DDPM [Ho et al., 2020]

Consider the *forward noising* process

$$X_k = \sqrt{1 - \beta_k} X_{k-1} + \sqrt{\beta_k} Z_k, \quad \beta_k \in [0, 1], \quad X_0 \sim \pi_{\text{data}},$$

where $Z_k \sim \mathcal{N}(\mathbf{0}_{d_x}, \mathbf{I}_{d_x})$.

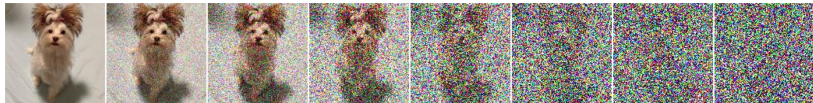


Figure: One sample $X_{0:n}$.

$X_k \sim \pi_k$ where $\pi_k(dx_k) := \int \pi_{\text{data}}(dx_0) \mathcal{N}(dx_k; \sqrt{\bar{\alpha}_k} x_0, (1 - \bar{\alpha}_k) \mathbf{I}_{d_x})$.

$(X_k)_{k \geq 0}$ is a discrete-time OU process.



$${}^1 \bar{\alpha}_k := \prod_{j=1}^k (1 - \beta_j).$$

Most times are used multiple times (MOLA)

Diffusion models: the backward process

Note that $\pi_{1:n|0}(x_{1:n}|x_0) = \pi_{n|0}(x_n|x_0) \prod_{k=2}^n \pi_{k-1|0,k}(x_{k-1}|x_0, x_k)$,
where $\pi_{n|0}(x_n|x_0) = \mathcal{N}(x_n; \bar{\alpha}_n^{1/2}x_0, (1 - \bar{\alpha}_n)\mathbf{I})$ and

$$\pi_{k-1|0,k}(x_{k-1}|x_0, x_k) = \mathcal{N}(x_{k-1}; \boldsymbol{\mu}_k(x_0, x_k), \sigma_k^2 \mathbf{I}_d) ,$$

with

$$\boldsymbol{\mu}_k(x_0, x_k) = \bar{\alpha}_{k-1}^{1/2}x_0 + (1 - \bar{\alpha}_{k-1} - \sigma_k^2)^{1/2}(x_k - \bar{\alpha}_k^{1/2}x_0)/(1 - \bar{\alpha}_k)^{1/2}.$$

↪ We know how to write **the joint distribution of $X_{1:n}$ given X_0** .

↪ Use this decomposition to turn **noise into samples from π_0** .

$$p_{0:n}^\theta(dx_{0:n}) = p_n(dx_n) \prod_{k=0}^{n-1} p_k^\theta(dx_k|x_{k+1}).$$



Most times are Used multiple times (MULE)

Diffusion models: the backward process

↔ Use this decomposition to turn **noise into samples from π_0** .

$$p_{0:n}^{\theta}(dx_{0:n}) = p_n(dx_n) \prod_{k=0}^{n-1} p_k^{\theta}(dx_k|x_{k+1}),$$

where p_n is a std Gaussian and

$$p_k^{\theta}(dx_k|x_{k+1}) = \mathcal{N}(dx_k; \mu_{k+1}^{\theta}(x_{k+1}), \beta_{k+1}I_{d_x})$$

with $\mu_{k+1}^{\theta}(x_{k+1})$ obtained by replacing x_0 in $\mu_{k+1}(x_0, x_{k+1})$ with a prediction

$$\hat{x}_{0|k,\theta}(x_{k+1}) := \bar{\alpha}_{k+1}^{-1/2} \left(x_{k+1} - (1 - \bar{\alpha}_{k+1})^{1/2} e^{\theta}(x_{k+1}, k+1) \right).$$



Most times are Used multiple times (MUSkox)

Diffusion models: the backward process

We use

$$p_k^\theta(dx_k|x_{k+1}) = \mathcal{N}(dx_k; \mu_{k+1}^\theta(x_{k+1}), \beta_{k+1}I_{d_x})$$

with $\mu_{k+1}^\theta(x_{k+1})$ obtained by replacing x_0 in $\mu_{k+1}(x_0, x_{k+1})$ with a prediction

$$\hat{x}_{0|k,\theta}(x_{k+1}) := \bar{\alpha}_{k+1}^{-1/2} \left(x_{k+1} - (1 - \bar{\alpha}_{k+1})^{1/2} e^\theta(x_{k+1}, k+1) \right).$$

$e^{\theta^*}(X_t, t)$ might be seen as the predictor of the noise added to X_0 to obtain X_t (in the forward pass) and justifies the *prediction* terminology.

The parameter θ is obtained by minimizing a variational loss between the forward and backward joint distributions.



moSt timEs Are used muLtiPle times (SEA Lion)

Diffusion models: an illustration

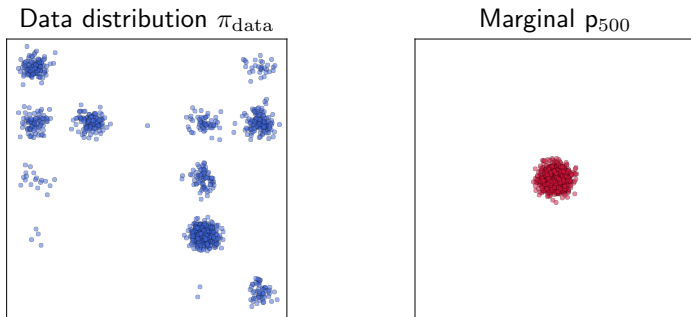


Figure: Samples from p_t for some time steps ranging from $n := 500$ to 1. π_0 is a Gaussian mixture.



Diffusion models: an illustration

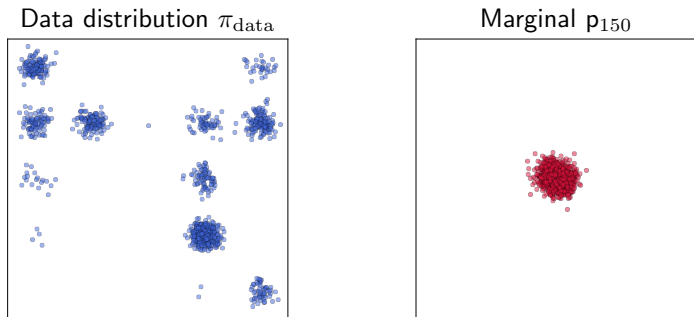


Figure: Samples from p_t for some time steps ranging from $n := 500$ to 1. π_0 is a Gaussian mixture.



Diffusion models: an illustration

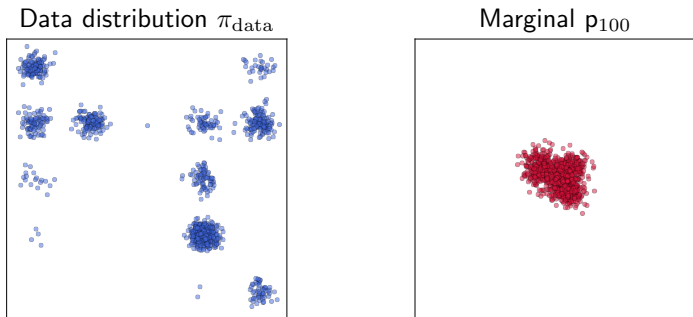


Figure: Samples from p_t for some time steps ranging from $n := 500$ to 1. π_0 is a Gaussian mixture.



Diffusion models: an illustration

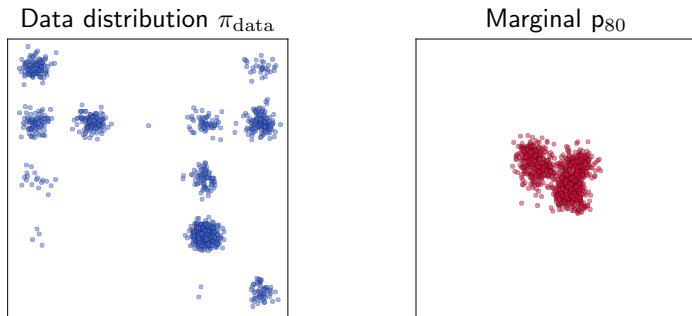


Figure: Samples from p_t for some time steps ranging from $n := 500$ to 1. π_0 is a Gaussian mixture.



Diffusion models: an illustration

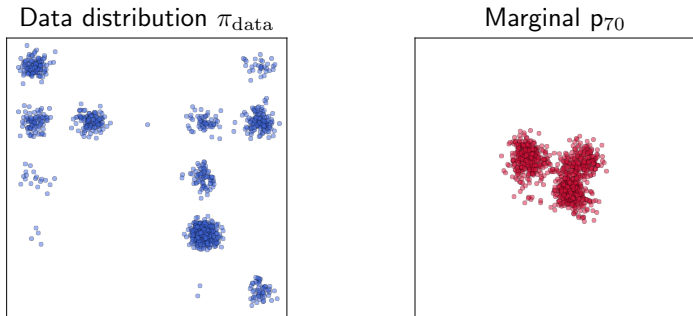


Figure: Samples from p_t for some time steps ranging from $n := 500$ to 1. π_0 is a Gaussian mixture.



Diffusion models: an illustration

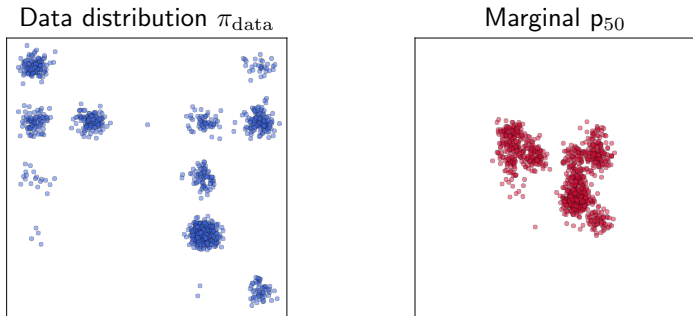


Figure: Samples from p_t for some time steps ranging from $n := 500$ to 1. π_0 is a Gaussian mixture.



Diffusion models: an illustration

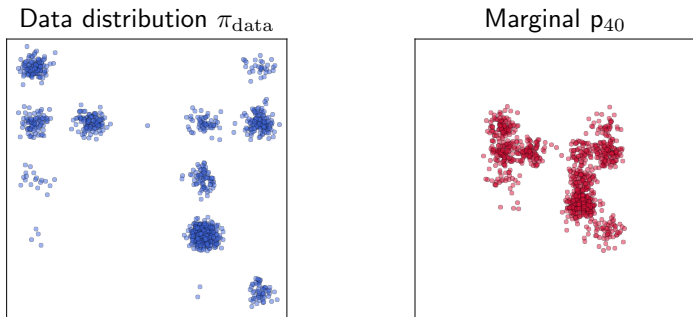


Figure: Samples from p_t for some time steps ranging from $n := 500$ to 1. π_0 is a Gaussian mixture.



Diffusion models: an illustration

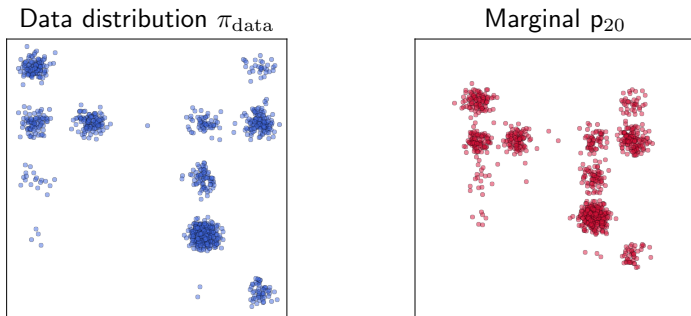


Figure: Samples from p_t for some time steps ranging from $n := 500$ to 1. π_0 is a Gaussian mixture.



Diffusion models: an illustration

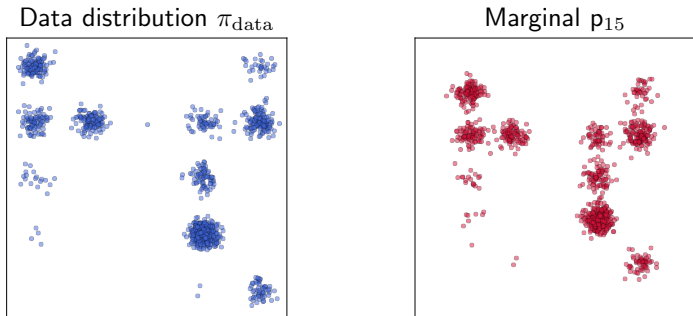


Figure: Samples from p_t for some time steps ranging from $n := 500$ to 1. π_0 is a Gaussian mixture.



Diffusion models: an illustration

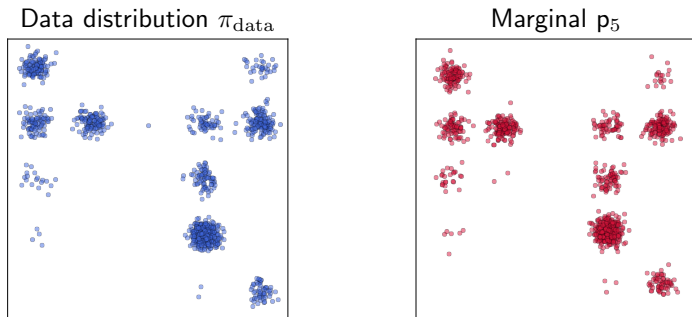


Figure: Samples from p_t for some time steps ranging from $n := 500$ to 1. π_0 is a Gaussian mixture.



Diffusion models: convergence

In the first theoretical results: **strong assumptions on the data distribution** (polynomial growth of the score), (de Bortoli et al., 2021):

$$\|\pi_{\text{data}} - p_0^\theta\|_{\text{tv}} \leq c_0 M \exp(c_1 n) + c_2 \left(n^{-1} + n^{-1/2} \right),$$

where M **quantifies the quality of the score approximation**.

In most recent works, we only require π_{data} to have a **finite relative Fisher information w.r.t the standard Gaussian distribution**, (Conforti et al., 2023):

$$\text{KL}(\pi_{\text{data}}, p_0^\theta) \leq \exp(-c_0 n) \text{KL}(\mathcal{N}(0, I_d), \pi_0) + Mn + c_1 h,$$

Assuming that $\mathbb{E}_{\pi_{\text{data}}} [\|\nabla \log(d\pi_{\text{data}}/d\gamma_d)\|^2] < \infty$.



moSt times Are used multiple TimEs (SATanic leaf-tailed gEcko)

Guarantees on the approximation of π_{data}
Score-based training procedures

SGM - Forward phase

Data noised using the Ornstein–Uhlenbeck (OU) process:

$$d\vec{X}_t = -\frac{\beta(t)}{2\sigma^2}\vec{X}_t dt + \sqrt{\beta(t)}dB_t, \quad \vec{X}_0 \sim \pi_{\text{data}}.$$

Fix $T > 0$, then,

$$\text{KL} \left(\mathcal{L} \left(\vec{X}_T \right), \pi_{\infty} \right) \lesssim \exp \left(-\frac{1}{2\sigma^2} \int_0^T \beta(s) ds \right) \text{KL} \left(\pi_{\text{data}}, \pi_{\infty} \right)$$

with

$$\text{KL} (\mu, \nu) := \int \log \left(\frac{d\mu(x)}{d\nu(x)} \right) \mu(dx).$$

Fokker-Planck for $(p_t)_{0 \leq t \leq T}$ + logarithmic Sobolev inequality + Gronwall's inequality



moSt tlmeS Are uSed multiple Times (SciSSor-tailed flycaTcher)

Time reversal

For $p_t := \mathcal{L}(\vec{X}_t)$, $(t, x) \mapsto \nabla \log p_t(x)$ is the **score function**. We consider the **time reversal** of the forward process, *i.e.*, the process satisfying

$$d\overleftarrow{X}_t = \left(\frac{\bar{\beta}(t)}{2\sigma^2} \overleftarrow{X}_t + \bar{\beta}(t) \nabla \log p_{T-t}(\overleftarrow{X}_t) \right) dt + \sqrt{\bar{\beta}(t)} dB_t, \quad \overleftarrow{X}_0 \sim p_T,$$

with $\bar{\beta}(t) := \beta(T - t)$. It satisfies

$$\left(\vec{X}_t \right)_{t \in [0, T]} = \left(\overleftarrow{X}_{T-t} \right)_{t \in [0, T]}.$$

In our setting,

$$\mathcal{L}(\vec{X}_T) = \mathcal{L}(\overleftarrow{X}_0) \approx \pi_\infty, \quad \mathcal{L}(\overleftarrow{X}_T) = \mathcal{L}(\vec{X}_0) \approx \pi_{\text{data}}.$$

moSt timEs Are uSed multiPle times (SEA SPIder)



Modified score

Let $\tilde{p}_t := p_t / \varphi_{\sigma^2}$, with φ_{σ^2} the density of π_∞ . With respect to the **modified score function**, the backward dynamics is

$$d\overleftarrow{X}_t = \left(-\frac{\bar{\beta}(t)}{2\sigma^2} \overleftarrow{X}_t + \bar{\beta}(t) \nabla \log \tilde{p}_{T-t}(\overleftarrow{X}_t) \right) dt + \sqrt{\bar{\beta}(t)} dB_t, \quad \overleftarrow{X}_0 \sim p_T.$$

If the score (or the modified score) is known, we can (in theory) simulate the backward process and get *data from noise*.

Let $s_\theta : [0, T] \times \mathbb{R}^d \mapsto \mathbb{R}^d$ be such that

$$\mathcal{L}_{\text{score}}(\theta) = \mathbb{E} \left[\left\| s_\theta \left(\tau, \overrightarrow{X}_\tau \right) - \nabla \log p_\tau \left(\overrightarrow{X}_\tau \right) \right\|^2 \right],$$

with $\tau \sim \mathcal{U}(0, T)$ independent of the forward process $(\overrightarrow{X}_t)_{t \geq 0}$.



mosT times aRE usEd multiple times (TREE creeper)

Discretization scheme

As the linear part can be simulated exactly, we consider the **Exponential Integrator scheme**:

Let $0 =: t_0 \leq t_1 \leq \dots \leq t_N := T$. Consider

$$d\overleftarrow{X}_t^\theta = \bar{\beta}(t) \left(-\frac{1}{2\sigma^2} \overleftarrow{X}_t^\theta + \tilde{s}_\theta \left(T - t_k, \overleftarrow{X}_{t_k}^\theta \right) \right) dt + \sqrt{\bar{\beta}(t)} dB_t,$$

for $t \in [t_k, t_{k+1})$, with $\tilde{s}_\theta(t, x) := s_\theta(t, x) + x/\sigma^2$, and $\overleftarrow{X}_0^\theta \sim \pi_\infty$.

We denote $\hat{\pi}_N^{(\beta, \theta)}$ the marginal probability density of $\overleftarrow{X}_T^\theta$.

The **loss function** is built using the conditional score:

$$\mathcal{L}(\theta) = \mathbb{E} \left[\alpha_\tau \left\| \nabla \log p_{\tau|0}(X_\tau | X_0) - s_\theta(\tau, X_\tau) \right\|^2 \right].$$



most times aRe used multiple times (River otter)

Assumptions

H1 The noise schedule is continuous, non decreasing and such that

$$\int_0^{\infty} \beta(t) dt = \infty.$$

H2 The data distribution has **finite Fisher information** w.r.t. the normal distribution, *i.e.*,

$$\mathcal{I}(\pi_{\text{data}}|\pi_{\infty}) := \int \left\| \nabla \log \left(\frac{d\pi_{\text{data}}}{d\pi_{\infty}} \right) \right\|^2 d\pi_{\text{data}} < \infty.$$

H3 The parameter θ and the schedule β satisfy

$$\mathbb{E} \left[\exp \left\{ \frac{1}{2} \int_0^T \bar{\beta}(t) \left\| \left(\tilde{s}(T-t, \bar{X}_t) - \tilde{s}_{\theta}(T-t_k, \bar{X}_{t_k}) \right) \right\|^2 dt \right\} \right] < \infty.$$



most times are used multiple times (River otter)

Why it works

$\bar{\mathbb{Q}}_N^{\beta, \theta} \in \mathcal{P}(C([0, T], \mathbb{R}^d))$: **path measure** (backward diffusion). $\hat{\pi}_N^{(\beta, \theta)}$ the **marginal probability density** of $\overleftarrow{X}_T^\theta$.

By the **data processing inequality**,

$$\text{KL} \left(\pi_{\text{data}} \left\| \hat{\pi}_N^{(\beta, \theta)} \right. \right) = \text{KL} \left(p_T \mathbb{Q}_T \left\| \hat{\pi}_N^{(\beta, \theta)} \right. \right) \leq \text{KL} \left(p_T \mathbb{Q}_T \left\| \pi_\infty \bar{\mathbb{Q}}_N^{\beta, \theta} \right. \right).$$

By applying **Girsanov theorem**,

$$\begin{aligned} \text{KL} \left(\pi_{\text{data}} \left\| \hat{\pi}_N^{(\beta, \theta)} \right. \right) &\leq \text{KL} (p_T \left\| \varphi_{\sigma^2} \right.) \\ &+ \frac{1}{2} \int_0^T \bar{\beta}(t) \mathbb{E} \left[\left\| \nabla \log \tilde{p}_{\tau_t} \left(\overleftarrow{X}_t \right) - \tilde{s}_\theta \left(\tau_k, \overleftarrow{X}_{t_k} \right) \right\|^2 \right] dt. \end{aligned}$$

most times are used multiple times (River otter)



Upper bound

Theorem

Assume that H1, H2 and H3 hold. Then,

$$\text{KL} \left(\pi_{\text{data}} \left\| \hat{\pi}_N^{(\beta, \theta)} \right. \right) \leq \mathcal{E}_1(\beta) + \mathcal{E}_2(\theta, \beta) + \mathcal{E}_3(\beta),$$

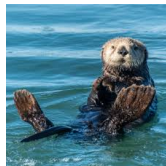
where

$$\mathcal{E}_1(\beta) := \text{KL}(\pi_{\text{data}} \|\pi_{\infty}) \exp \left\{ -\frac{1}{\sigma^2} \int_0^T \beta(s) ds \right\},$$

$$\mathcal{E}_2(\theta, \beta) := \sum_{k=1}^N \int_{t_k}^{t_{k+1}} \beta(t) dt \mathbb{E} \left[\left\| \nabla \log \tilde{p}_{T-t_k} \left(\vec{X}_{T-t_k} \right) - \tilde{s}_{\theta} \left(T - t_k, \vec{X}_{T-t_k} \right) \right\|^2 \right],$$

$$\mathcal{E}_3(\beta) := 2h\beta(T) \mathcal{I}(\pi_{\text{data}} \|\pi_{\infty}),$$

with $h := \sup_{k \in \{1, \dots, N\}} (t_k - t_{k-1})$ and $t_0 := 0$.



Gaussian case

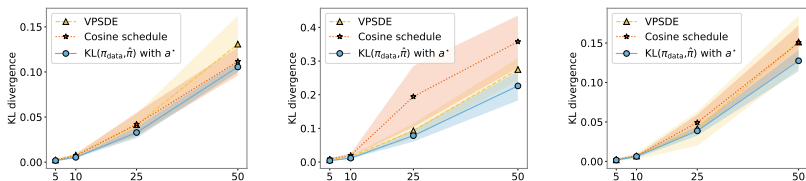
Let the true distribution be Gaussian in dimension $d = 50$ with mean $\mathbf{1}_d$ and different choices of covariance structure.

1. **(Isotropic)** $\Sigma^{(\text{iso})} = 0.5\mathbf{I}_d$.
2. **(Heteroscedastic)** $\Sigma^{(\text{heterosc})} \in \mathbb{R}^{d \times d}$ is a diagonal matrix such that $\Sigma_{jj}^{(\text{heterosc})} = 10$ for $1 \leq j \leq 5$, and $\Sigma_{jj}^{(\text{heterosc})} = 0.1$ otherwise.
3. **(Correlated)** $\Sigma^{(\text{corr})} \in \mathbb{R}^{d \times d}$ is a full matrix whose diagonal entries are equal to one and the off-diagonal terms are $\Sigma_{jj'}^{(\text{corr})} = 1/\sqrt{|j - j'|}$ for $1 \leq j \neq j' \leq d$.



most times are used multiple times (River otter)

Comparison with existing literature



(a) Isotropic setting (b) Heteroscedastic setting (c) Correlated setting

Figure: Comparison of the empirical KL divergence between π_{data} and the generative distribution $\hat{\pi}_N^{(\beta, \theta)}$ w.r.t. SGM for β_{a^*} , the VPSDE model and the one with a cosine schedule, presented in Chen et al. (2023).

- ▶ Optimizing the noise schedule has an impact even with simple parametrization of the β scheduling.



most times aRe used multiple times (River otter)

Conclusion and perspectives

- ↪ Optimizing the noise schedule allows to optimize virtually all score-based samplers.
- ↪ Constants in the upper bound are crucial and should be taken care of.
- ↪ Important open problem: design π_{data} tailored (Swift) to π_{∞} (then maybe use Langevin for the forward).

most times are used multiple times (River otter)

