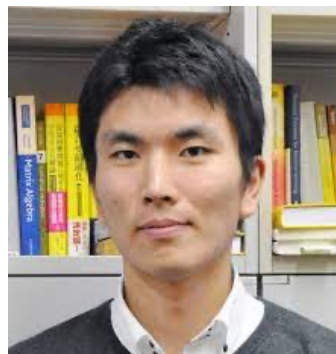


Diffusion Models are Minimax Optimal Distribution Estimators

Kazusato Oko (The University of Tokyo / AIP RIKEN)

Joint work with Shunta Akiyama (The University of Tokyo)

Taiji Suzuki (The University of Tokyo / AIP RIKEN)



THE UNIVERSITY OF TOKYO



May 4, 2023
ME-FoMo spotlight talk
ICLR 2023 Workshop

Success of diffusion models

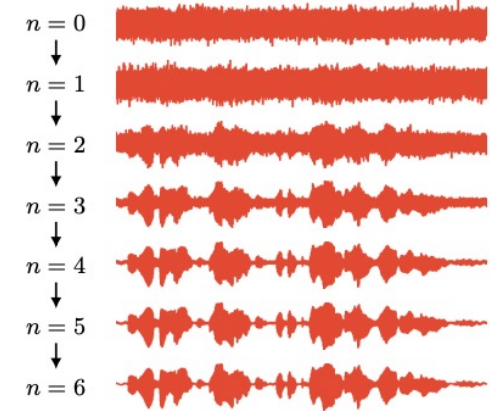
- Image, video, audio, etc...



Image generated by DALL·E2



Video generated by Video Diffusion Models



Visualization of WaveGrad

- Several pictures I made with Stable Diffusion:



A mecha robot in a favela in an expressionist style



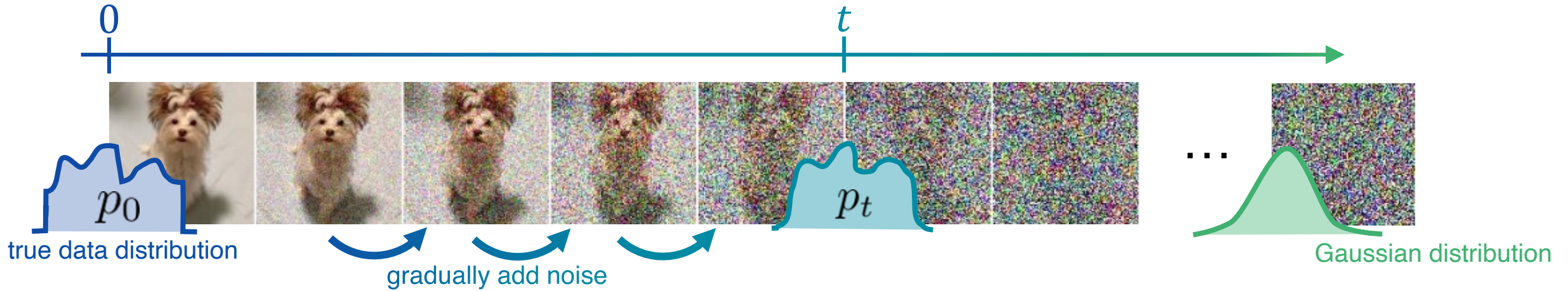
A high tech solarpunk utopia in the Amazon forest



Pikachu in the style of Munch's "The Scream"

Formulation as SDE (Song et al., 2020): forward process

3



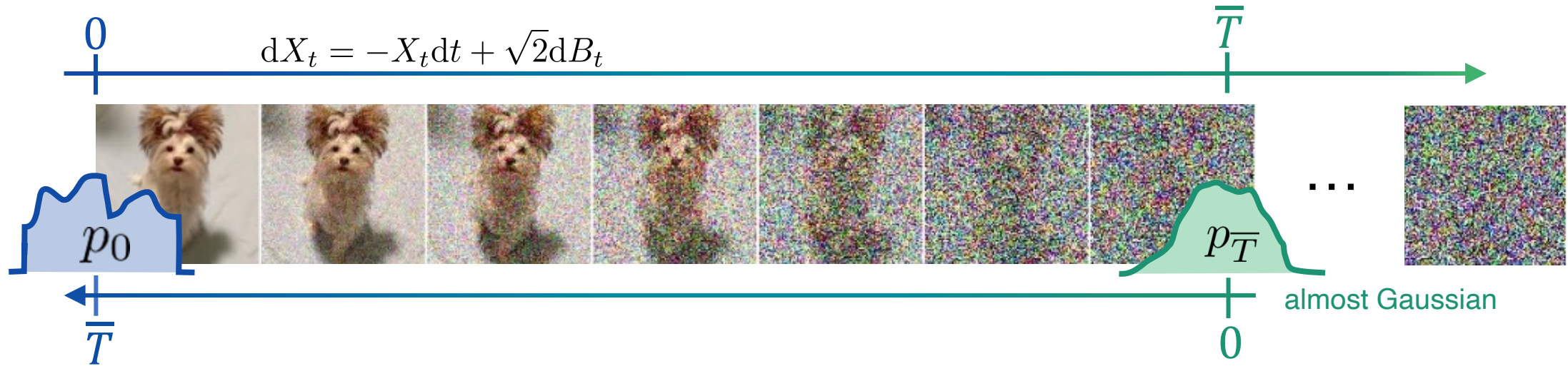
$$X_0 \sim p_0, \quad dX_t = -X_t dt + \sqrt{2} dB_t \quad (\text{OU process})$$

Brownian motion

- The distribution of X_t at time t :

$$p_t(x) = \int p_0(y) \frac{1}{\sigma_t^d (2\pi)^{\frac{d}{2}}} \exp\left(-\frac{\|x - \mu_t y\|^2}{2\sigma_t^2}\right) dy$$

Formulation as SDE (Song et al., 2020): backward process

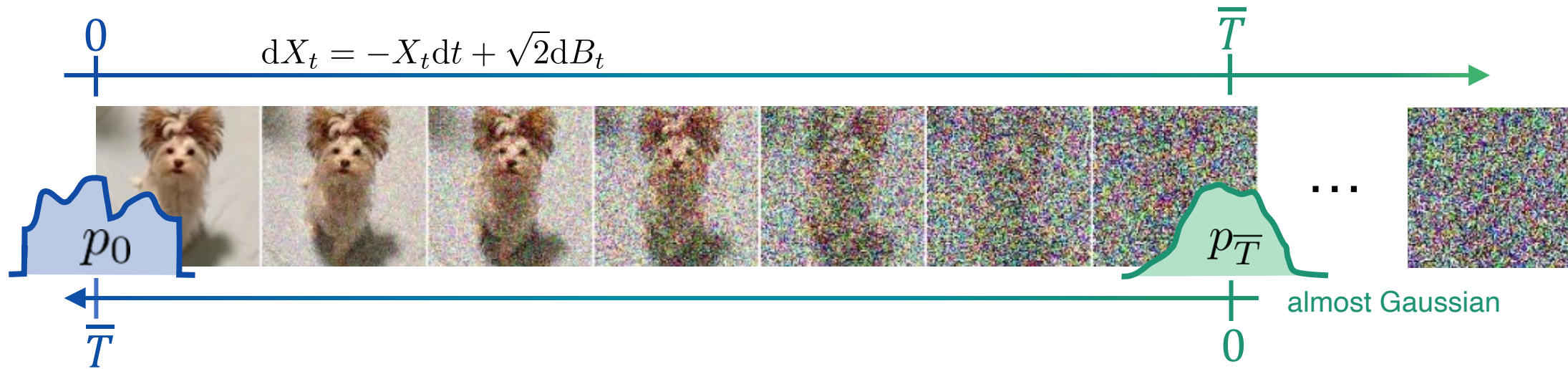


$$Y_0 \sim p_{\bar{T}}, \quad dY_t = (Y_t + 2\nabla \log p_{\bar{T}-t}(Y_t))dt + \sqrt{2}dB_t$$

Brownian motion

→ $Y_{\bar{T}} \sim p_0$ (recovers the true data distribution)

Formulation as SDE (Song et al., 2020): backward process



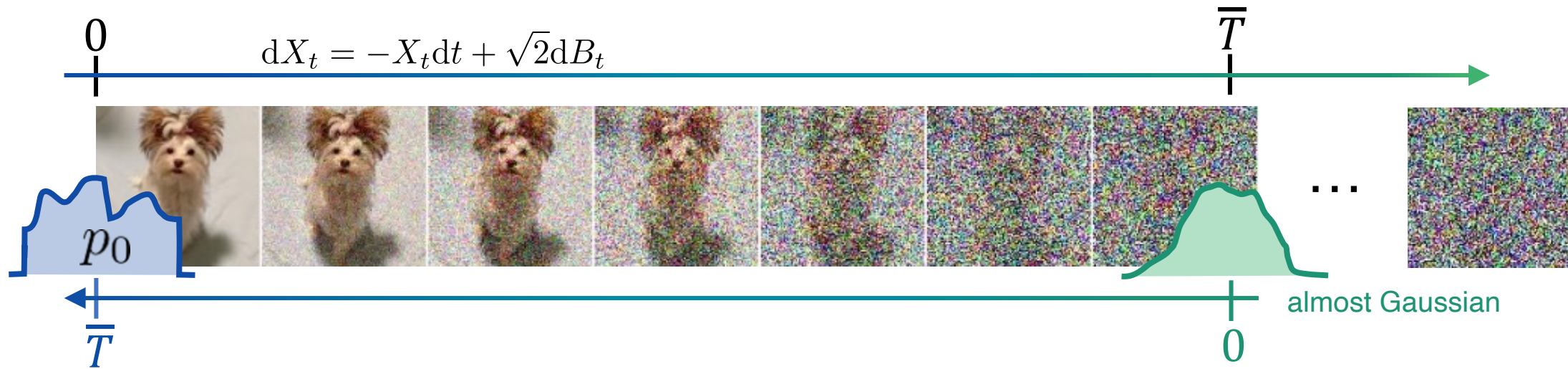
$$Y_0 \sim p_{\bar{T}},$$

$$dY_t = (Y_t + 2 \nabla \log p_{\bar{T}-t}(Y_t)) dt + \sqrt{2} dB_t$$

Brownian motion

The exact value of **the score** $\nabla \log p_t(x)$ cannot be obtained because it **depends on** p_0

Formulation as SDE (Song et al., 2020): backward process



$$Y_0 \sim p_{\bar{T}},$$

$$dY_t = (Y_t + 2 \nabla \log p_{\bar{T}-t}(Y_t)) dt + \sqrt{2} dB_t$$

Brownian motion

The exact value of the score $\nabla \log p_t(x)$ cannot be obtained because it depends on p_t

$$\hat{Y}_0 \sim \mathcal{N}(0, I),$$

$$d\hat{Y}_t = (\hat{Y}_t + 2 \hat{s}(\hat{Y}_t, \bar{T} - t)) dt + \sqrt{2} dB_t$$

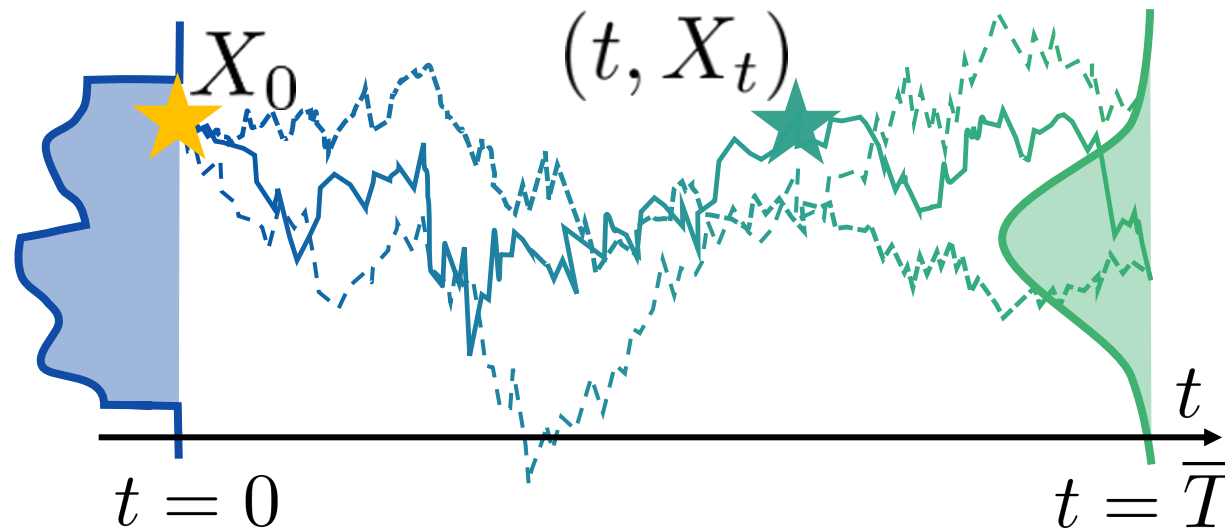
the score network, trained with finite sample

Score matching to train the score network

- The true score minimizes the following loss:

$$\mathbb{E}_{X_0 \sim p_0} \left[\underbrace{\int_{t=0}^{\bar{T}} \mathbb{E}_{X_t | X_0 \sim \mathcal{N}(m_t X_0, \sigma_t^2)} [\|s(X_t | X_0, t) - \nabla \log p_t(X_t | X_0)\|^2] dt}_{\text{Computed by sampling } (t, X_t) \sim \text{Unif}[0, \bar{T}] \times \mathcal{N}(m_t X_0, \sigma_t^2)} \right]$$

Computed by sampling $(t, X_t) \sim \text{Unif}[0, \bar{T}] \times \mathcal{N}(m_t X_0, \sigma_t^2)$



Score matching to train the score network

- The true score minimizes the following loss:

$$\mathbb{E}_{X_0 \sim p_0} \left[\int_{t=0}^{\bar{T}} \mathbb{E}_{X_t | X_0 \sim \mathcal{N}(m_t X_0, \sigma_t^2)} [\|s(X_t | X_0, t) - \nabla \log p_t(X_t | X_0)\|^2] dt \right]$$

- Replace the expectation w.r.t. p_0 by finite sample $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} p_0$

$$\operatorname{argmin}_{s \in \mathcal{S}: \text{DNNs}} \frac{1}{n} \sum_{i=1}^n \left[\int_{t=0}^{\bar{T}} \mathbb{E}_{X_t | x_i \sim \mathcal{N}(m_t x_i, \sigma_t^2)} [\|s(X_t | x_i, t) - \nabla \log p_t(X_t | x_i)\|^2] dt \right]$$

empirical score matching loss

How close the generated distribution is to p_0 ?

Score matching to train the score network

- The true score minimizes the following loss:

$$\mathbb{E}_{X_0 \sim p_0} \left[\int_{t=0}^{\bar{T}} \mathbb{E}_{X_t | X_0 \sim \mathcal{N}(m_t X_0, \sigma_t^2)} [\|s(X_t | X_0, t) - \nabla \log p_t(X_t | X_0)\|^2] dt \right]$$

- Replace the expectation w.r.t. p_0 by finite sample $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} p_0$

$$\operatorname{argmin}_{s \in \mathcal{S}: \text{DNNs}} \frac{1}{n} \sum_{i=1}^n \left[\int_{t=0}^{\bar{T}} \mathbb{E}_{X_t | x_i \sim \mathcal{N}(m_t x_i, \sigma_t^2)} [\|s(X_t | x_i, t) - \nabla \log p_t(X_t | x_i)\|^2] dt \right]$$

empirical score matching loss

**How good is the diffusion modeling
as a distribution estimator based on n data?**

Existing work on error analysis

- Most works analyzed how the score estimation error, **given as an assumption**, affects the generated data distribution
 - ❖ **Continuous time**: Song et al. (2021); De Bortoli et al. (2022a)
 - ❖ **Discrete time**: Lee et al. (2022a;b); Chen et al. (2023)
 - ❖ Non-quantitative bound under manifold assumption: Pidstrigach (2022)
- Sample complexity bounds
 - ❖ **W1 bound of $n^{-1/d}$ with manifold assumption**: De Bortoli et al. (2021)
 - * **Not considering “generalization” and unimprovable**
 - * Based on the convergence of the empirical measure $W_1(\frac{1}{n} \sum_{i=1}^n \delta_{x_i}, p_0) \simeq n^{-1/d}$ (Weed and Bach, 2019)
- **Concurrent work** (appeared after the submission of this work): Chen et al. (2023)

Analysis of diffusion models

from statistical learning theory

A1 p_0 is supported on $[-1,1]^d$, upper and lower bounded in the support, and

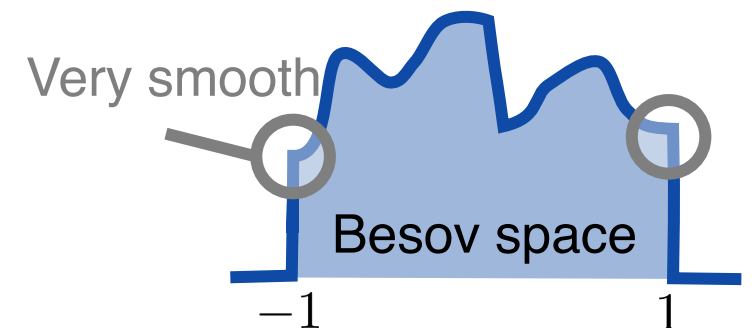
$$p_0 \in B_{p,q,C}^s$$

with $s > (1/p - 1/2)_+$ as a density function on $[-1,1]^d$.

• $B_{p,q,C}^s$: Besov space $B_{p,q}^s$ with the norm bounded by C (some constant)

❖ Intuition: $\|f\|_{B_{p,q}^s(\Omega)} = \|f\|_{L^p(\Omega)} + \|D^s f\|_{L^p(\Omega)}$

A2 p_0 is sufficiently smooth on the edge of the support $[-1,1]^d \setminus [-1 + n^{-\frac{1-\delta}{d}}, 1 - n^{-\frac{1-\delta}{d}}]^d$.



Analysis of diffusion models from statistical learning theory

- Hypothesis network class: sparsity-constrained deep ReLU networks

(Schmidt-Hieber, 2020; Suzuki, 2019)

$\mathcal{S}(L \text{ (depth)}, W \text{ (width)}, S \text{ (sparsity-constraint; num. of non-zero params)}, B \text{ (magnitude)})$

$:= \{(A^L \text{ReLU}(\cdot) + b^L) \circ \dots \circ (A^1 x + b^1) \mid A^i \in \mathbb{R}^{w_i \times w_{i+1}}, b^i \in \mathbb{R}^{w_{i+1}}, \|w\|_\infty \leq W,$

$$\left. \sum_{i=1}^L (\|A^i\|_0 + \|b^i\|_0) \leq S, \max \|A^i\|_\infty \vee \|b^i\|_\infty \leq B \right\}$$

- Sparsity-constraint yields tighter generalization error bounds

Main result ①: minimax optimality in TV

Theorem 1

The generated data distribution by using the score network \hat{S} that minimizes the empirical score matching loss over $\mathcal{S}(L, W, S, B)$ yields that

$$\mathbb{E}_{\{x_i\}_{i=1}^n} \left[\text{TV}(\hat{Y}_{\bar{T}}, X_0) \right] \lesssim n^{-\frac{s}{2s+d}} \log^9(n),$$

under the appropriate choice of $\bar{T}, \underline{T}, L, W, S$, and B .

This rate is **the minimax optimal** (up to polylog), because it also holds that

$$n^{-\frac{s}{2s+d}} \lesssim \inf_{\hat{\mu}: \text{estimator}} \sup_{p_0 \in B_{p,q,C}^s} \mathbb{E}_{\{x_i\}_{i=1}^n} [\text{TV}(\hat{\mu}, X_0)].$$

Basis decomposition tailored for score approximation 14

- B-spline basis decomposition of $p_0 (\in B_{p,q,C}^s)$: $p_0(x) \approx \sum_{j=1}^N \alpha_j M_{a^j, b^j}^d(x)$
(Devore & Popov, 1988) B-spline basis

- Approximation of $p_t(x)$:

$$p_t(x) = \int p_0(y) \underbrace{\frac{1}{\sigma_t^d (2\pi)^{\frac{d}{2}}} \exp\left(-\frac{\|x - \mu_t y\|^2}{2\sigma_t^2}\right)}_{=: K_t(x|y)} dy$$

approximation via
B-spline basis

$$\approx \sum_{j=1}^N \alpha_j \int M_{a^j, b^j}^d(y) K_t(x|y) dy$$

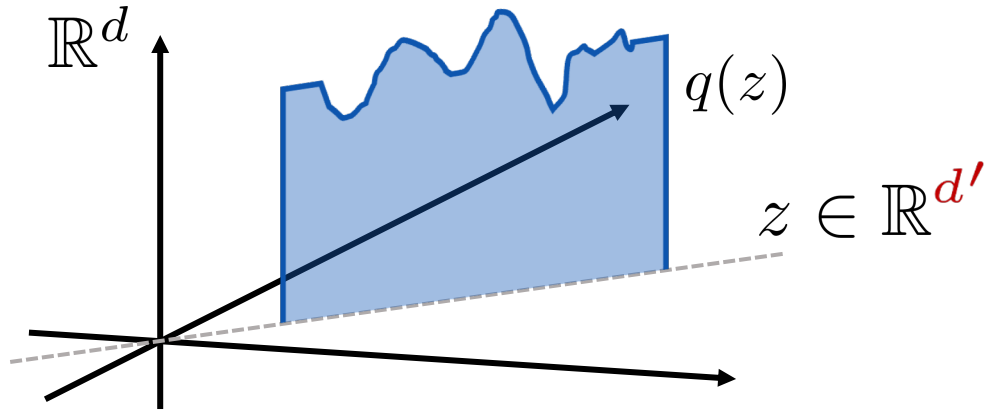
$=: E_{a^j, b^j}(x, t)$ **diffused B-spline basis**

Approximated by NN very efficiently (polylog size)



- ❖ Approximate $\nabla p_t(x)$ in the same way and use $\nabla \log p_t(x) = \frac{\nabla p_t(x)}{p_t(x)}$

Main result ②: manifold hypothesis



- p_0 lies on a d' -dimensional plane ($d' \leq d$)
- Density function on the canonical coordinate system on the plane q belongs to $B_{p,q,C}^s$

Theorem 2

Based on $\{x_i\}_{i=1}^n$, we can train the score network \hat{s} that satisfies

$$\mathbb{E}_{\{x_i\}_{i=1}^n} \left[\mathbf{W}_1(\hat{Y}_{\bar{T}}, X_0) \right] \lesssim n^{-\frac{s+1-\delta}{2s+d'}}.$$

($\delta(> 0)$): arbitrarily fixed constant)

Diffusion models can avoid the curse of dimensionality!