

AI Applications Lecture 14

Image Generation AI 4: Goals and Scheduling of Diffusion Processes

SUZUKI, Atsushi

Jing WANG

Contents

1	Introduction	3
1.1	Recap	3
1.2	Learning Outcomes	3
2	Preparation: Mathematical Notations	3
3	Revisiting the Goal and Steps of the Reverse Diffusion Process	4
3.1	Goal: Sampling from a Conditional Target Distribution	4
3.2	Invalidity of the Naive Method and Motivation	4
3.3	Review of the Role of Each Step	5
3.4	Discrete-Time Reverse Diffusion and Stability	5
4	The Non-Triviality of Learning and Forward Noising	5
4.1	What is Observable and What is Missing	5
4.2	Two Simultaneous Challenges	6
5	Available Probabilistic Tools and Construction of Artificial Distribution Sequence	6
5.1	Constraints on High-Dimensional Distributions	6
5.2	Constructing Random Variables by Linear Combination	6
5.3	Designing a Smooth Distribution Sequence q_k	6
5.4	Global Schedule and Training Data Sequence	7
5.5	Noise Estimator Learning Objective and x Reconstruction	7

6 Goal and Tools: Connection to Scheduler Design	8
6.1 Scheduler Design Goal	8
6.2 Available Tool	8
7 Overall Strategy for Denoising Scheduler Construction	8
8 Deterministic Scheduling	10
8.1 Revisiting the Goal	10
8.2 Abstract Definition of One-Step Method	10
8.3 First-Order Expansion of KL Divergence for One-Step Method	10
8.4 Optimization of One-Step Method Coefficients (DDIM/Euler)	12
8.5 Abstract Definition of Two-Step Method	14
9 Markovian Scheduling	16
9.1 Revisiting the Goal	16
9.2 Local Linear Noising and Gaussian Approximation of Reverse Conditional . .	16
9.3 \tilde{q} as a Discrete-Time Diffusion Process	17
9.4 Abstract Definition of Markovian Update and Divergence Expansion	18
9.5 Optimal Coefficients and DDPM/Euler a	19
10 Summary and Next Time	19
10.1 Summary Corresponding to Learning Outcomes	19
10.2 Next Time	20

1 Introduction

1.1 Recap

Last time, we learned that we can generate low-resolution **latent images** through continuous update steps using **pseudo-random numbers** and a **denoising scheduler**. However, what each update step **aims for** was largely based on intuition. In this lecture, we will explain the **meaning of the scheduler's update equations** and the **design of convergence to the target distribution** by formulating them mathematically and rigorously.

1.2 Learning Outcomes

By the end of this lecture, students should be able to:

- Explain how to obtain **data points** using a **diffusion process** for **implicit distribution learning** in the context of **sampling** from a **target distribution**.
- Explain how training a **noise estimator** using **data point pairs generated by adding noise** simultaneously solves **two difficulties**: (i) learning from **realistically obtainable data** and (ii) **realizing the target distribution** through a **reverse diffusion process**.
- Explain, through theorems and calculations, in what sense a **denoising scheduler approaches the target distribution with each update**.

2 Preparation: Mathematical Notations

- **Definition:**

- (LHS) := (RHS): The left-hand side is defined by the right-hand side.

- **Set:**

- Sets are often denoted by uppercase calligraphic letters. E.g.: \mathcal{A} .
 - $x \in \mathcal{A}$: The element x belongs to the set \mathcal{A} .
 - $\{\}$: The empty set.
 - $\{a, b, c\}$: The set consisting of elements a, b, c (extensional notation).
 - $\{x \in \mathcal{A} \mid P(x)\}$: The set of elements for which the proposition $P(x)$ is true (intensional notation).
 - $|\mathcal{A}|$: The number of elements in set \mathcal{A} (finite sets are assumed in this lecture).
 - $\mathbb{R}, \mathbb{R}_{>0}, \mathbb{R}_{\geq 0}, \mathbb{Z}, \mathbb{Z}_{>0}, \mathbb{Z}_{\geq 0}$ have their standard meanings.

- $[1, k]_{\mathbb{Z}}$: Integer interval for $k \in \mathbb{Z}_{>0} \cup \{+\infty\}$.

- **Function:**

- Notation $f : \mathcal{X} \rightarrow \mathcal{Y}$, $y = f(x)$.

- **Vector:**

- Vectors are column vectors, denoted by bold italic lowercase \mathbf{v} .
- $\mathbf{v} = \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix}^T \in \mathbb{R}^n$.
- Standard inner product $\langle \mathbf{u}, \mathbf{v} \rangle := \sum_{i=1}^n u_i v_i$.

- **Sequence:**

- $\mathbf{a} : [1, n]_{\mathbb{Z}} \rightarrow \mathcal{A}$ is called a sequence of length n .

- **Matrix:**

- Matrices are bold italic uppercase $\mathbf{A} \in \mathbb{R}^{m,n}$. Transpose $\mathbf{A}^T \in \mathbb{R}^{n,m}$.

- **Tensor:**

- Tensors as multidimensional arrays are written as $\underline{\mathbf{A}}$.

3 Revisiting the Goal and Steps of the Reverse Diffusion Process

3.1 Goal: Sampling from a Conditional Target Distribution

The **goal** is to achieve **sampling** from a **target distribution** P_c that depends on a text condition $c \in \mathbb{R}^{d_{\text{AllText}}}$. Here,

$$c = (c^{[j]})_{j=1}^n \quad (1)$$

is an output vector sequence from **text encoders**, and although it generally consists of multiple vectors, in this lecture, we treat it as **concatenated into a single vector**.

3.2 Invalidity of the Naive Method and Motivation

The most naive method considered is to **directly output** the **low-resolution latent images that appeared in the training data** corresponding to c , or to add **Gaussian noise** to them. However, this is **practically meaningless**. The c used in inference is unlikely to **reappear** from training, making the direct output of training data not useful, and minor modifications with Gaussian noise tend to be **unnatural as natural images**. What we want is a **distribution** that is **continuous** with respect to c (allowing interpolation to unseen c) and corresponds to **natural images**. For this reason, we employ the **reverse diffusion process**,

which uses the **continuity** and **nonlinearity** of **neural networks** to achieve sampling from a **non-Gaussian** distribution by **push-forward** from a **simple base distribution**.

3.3 Review of the Role of Each Step

In general, the **composition of functions**, such as neural networks, only provides a **point-to-point correspondence between input and output**, not a **distribution** directly. Therefore, we achieve **sampling** by providing **(pseudo-)random numbers** to the **input side** to obtain a **push-forward distribution**. In other words, it is a **trick** to "sample from a distribution using an input-output relationship."

3.4 Discrete-Time Reverse Diffusion and Stability

We take a discrete time sequence $T = t_0 > t_1 > \dots > t_K = 0$. In a one-step reverse diffusion process, we use a function

$$g_{\theta} : \mathbb{R}^{d_{\text{Latent}}} \times \mathbb{R}^{d_{\text{Latent}}} \times \mathbb{R}^{d_{\text{AllText}}} \times \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^{d_{\text{Latent}}} \quad (2)$$

(containing a neural network internally), start from a standard multivariate normal

$$\mathbf{x}_{t_0} \sim \text{StdNormal}_{d_{\text{Latent}}} \quad (3)$$

and update as

$$\mathbf{x}_{t_{k+1}} = g_{\theta}\left(\mathbf{x}_{t_k}, \boldsymbol{\epsilon}^k, \mathbf{c}, t_k, t_k - t_{k+1}\right), \quad k = 0, 1, \dots, K-1 \quad (4)$$

For **inference stability**, it is designed such that at each step

$$\|\mathbf{x}_{t_{k+1}} - \mathbf{x}_{t_k}\|_2 \text{ is sufficiently small.} \quad (5)$$

If we write the distribution of \mathbf{x}_{t_k} as P_{t_k} , the ideal is

$$P_{t_0} = \text{StdNormal} \Rightarrow P_{t_1} \Rightarrow \dots \Rightarrow P_{t_K} \approx P_c. \quad (6)$$

4 The Non-Triviality of Learning and Forward Noising

4.1 What is Observable and What is Missing

What we can actually obtain is only the **original data** x , which can be regarded as generated from $P_{t_K} = P_c$. The corresponding $\mathbf{x}_{t_{K-1}}, \dots, \mathbf{x}_{t_0}$ are **not uniquely given**. Therefore, it is necessary to **artificially** construct $\mathbf{x}_{t_{K-1}}, \dots, \mathbf{x}_{t_0}$. The distributions these follow correspond to $P_{t_{K-1}}, \dots, P_{t_0}$.

4.2 Two Simultaneous Challenges

Surveying the situation, we need to satisfy the following **two challenges simultaneously**:

- Making the final distribution P_{t_0} match (or sufficiently approximate) the **target distribution** P_c .
- Being **learnable** through **realistic operations** (possible with available computational resources and data).

5 Available Probabilistic Tools and Construction of Artificial Distribution Sequence

5.1 Constraints on High-Dimensional Distributions

The only practical distributions that are **directly realizable** in high dimensions and defined by a **small number of parameters** are the **isotropic Gaussian distribution**, the **Cauchy distribution**, and the **uniform distribution within a hypersphere**. However, the Cauchy distribution **lacks an expected value** and is difficult to handle, and the uniform distribution within a hypersphere lacks a **reproductive property with respect to sums** and is also difficult to handle. Therefore, the **isotropic Gaussian distribution** is **effectively the only tool**. An isotropic Gaussian distribution is determined by a **mean vector** and a **scalar standard deviation**.

5.2 Constructing Random Variables by Linear Combination

Using a sample x from the original data distribution P_c and an independent $\epsilon \sim \text{StdNormal}_d$, we construct

$$\zeta_{\lambda_{\text{signal}}, \lambda_{\text{noise}}} := \lambda_{\text{signal}}x + \lambda_{\text{noise}}\epsilon. \quad (7)$$

From a computational cost perspective, if we limit the use to one random vector per data point, the attainable random variables are effectively limited to the form (7). The distribution it follows is written as

$$p_{\zeta_{\lambda_{\text{signal}}, \lambda_{\text{noise}}}}. \quad (8)$$

5.3 Designing a Smooth Distribution Sequence q_k

To transition **smoothly** from q_0 to q_K , so that the mean and variance move smoothly, we define

$$0 = \bar{\alpha}_0 < \bar{\alpha}_1 < \dots < \bar{\alpha}_{K-1} < \bar{\alpha}_K = 1 \quad (9)$$

and set

$$q_k := p_{\zeta_{\sqrt{\alpha_k}, \sqrt{1-\alpha_k}}}, \quad (10)$$

$$\zeta_k := \sqrt{\alpha_k} \mathbf{x} + \sqrt{1 - \alpha_k} \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \text{StdNormal}_d. \quad (11)$$

Remark 5.1. Let the mean and covariance of \mathbf{x} be \mathbf{m} and \mathbf{V} , respectively. Then

$$\mathbb{E}[\zeta_k] = \sqrt{\alpha_k} \mathbf{m}, \quad (12)$$

$$\text{Cov}(\zeta_k) = \alpha_k \mathbf{V} + (1 - \alpha_k) \mathbf{I}. \quad (13)$$

Thus, as k increases, the mean **monotonically** approaches \mathbf{m} from 0, and the covariance **monotonically** approaches \mathbf{V} from \mathbf{I} .

5.4 Global Schedule and Training Data Sequence

The increasing sequence $\{\bar{\alpha}_k\}_{k=0}^K$ used during inference may not be known during training. Therefore, to cover as wide a range as possible, we choose a sufficiently large $T \in \mathbb{Z}_{>0}$ and fix a

$$1 = \bar{\alpha}_0 > \bar{\alpha}_1 > \dots > \bar{\alpha}_T = 0 \quad (14)$$

decreasing sequence (adopting the decreasing direction by convention; in principle, increasing is also fine). If $\bar{\alpha}_t \approx \bar{\alpha}_k$, then t roughly corresponds to k . Furthermore,

$\mathbf{x}^{(i)}$ is assumed to be independently generated from $q_K = P_{\mathbf{c}^{(i)}}$, (15)

$t^{(i)} \sim \text{Unif}\{0, 1, \dots, T - 1\}$, $\boldsymbol{\epsilon}^{(i)} \sim \text{StdNormal}_d$, (16)

$\zeta^{(i)} = \sqrt{\bar{\alpha}_{t^{(i)}}} \mathbf{x}^{(i)} + \sqrt{1 - \bar{\alpha}_{t^{(i)}}} \boldsymbol{\epsilon}^{(i)}$. (17)

Thus, a **sequence of points** $(\zeta^{(i)}, t^{(i)}, \mathbf{x}^{(i)})_{i=1}^m$ is obtained.

5.5 Noise Estimator Learning Objective and x Reconstruction

Using a neural network

$$\hat{\boldsymbol{\epsilon}}_{\theta} : \mathbb{R}^{d_{\text{Latent}}} \times \mathbb{R}^{d_{\text{AllText}}} \times [0, T] \rightarrow \mathbb{R}^{d_{\text{Latent}}} \quad (18)$$

we minimize the following objective function:

$$\min_{\theta} \sum_{i=1}^m \left\| \boldsymbol{\epsilon}^{(i)} - \hat{\boldsymbol{\epsilon}}_{\theta}(\zeta^{(i)}, \mathbf{c}^{(i)}, t^{(i)}) \right\|_2^2. \quad (19)$$

Here, in implementation, $\hat{\epsilon}$ is estimating the **noise**, but due to the **linear relationship**

$$x = \frac{1}{\sqrt{a_t}} \zeta - \frac{\sqrt{1 - \bar{a}_t}}{\sqrt{a_t}} \epsilon, \quad (20)$$

if $\hat{\epsilon}$ is obtained, an estimate of x is also **immediately** obtained.

Remark 5.2. When t is large (strong noise), $\zeta \approx \epsilon$, and information about x is scarce. In particular, at $t = T$, $\zeta = \epsilon$. Therefore, naively generating using $\hat{\epsilon}_\theta(\epsilon, c, T)$ is not appropriate.

6 Goal and Tools: Connection to Scheduler Design

6.1 Scheduler Design Goal

The **goal** is to construct an appropriate function (update equation) and provide a **scheduler** such that its **push-forward distribution**

$$q_0, q_1, \dots, q_{K-1}, q_K \quad (21)$$

gradually approaches the target distribution q_K (i.e., the ideal distribution matching P_c) at each step. Here, q_k is the distribution that

$$\zeta_{\sqrt{\alpha_k}, \sqrt{1-\bar{\alpha}_k}} = \sqrt{\alpha_k} x + \sqrt{1 - \bar{\alpha}_k} \epsilon \quad (22)$$

follows.

6.2 Available Tool

The main tool available is the **noise estimator** $\hat{\epsilon}_\theta$ ((18)), and the objective function to obtain it is (19). This is expected to **indirectly** contain information about q_0, \dots, q_K . Therefore, the challenge is how to construct **sampling** from (21) using this $\hat{\epsilon}_\theta$.

Remark 6.1. Note that when t is large (k is small), $\hat{\epsilon}_\theta$ is **almost useless**. In particular, at $t = T$, $\zeta = \epsilon$ and contains no x information, so one should not naively generate using $\hat{\epsilon}_\theta(\epsilon, c, T)$.

7 Overall Strategy for Denoising Scheduler Construction

Finally, we construct the **denoising scheduler**. The **goal** is to use the trained noise estimator $\hat{\epsilon}_\theta$ to construct a random variable sequence

$$z_0, z_1, z_2, \dots, z_{K-1}, z_K \quad (23)$$

such that the distribution of each z_k satisfies

$$\text{Law}(z_k) \approx q_k \quad (k = 0, 1, \dots, K) \quad (24)$$

(where q_k is the artificial distribution sequence defined in (10)–(11)). To achieve this goal, we consider the following two **main strategies**:

- **Transformation by constructing deterministic functions (deterministic push-forward):** Construct a **deterministic** recurrence relation such as

$$z_{k+1} = \mathcal{F}_k(z_0, z_1, \dots, z_k; \hat{\epsilon}_\theta, c, t_k, h_k) \quad (25)$$

and design the coefficients and step sizes so that the distribution of z_{k+1} **approaches** q_{k+1} .

Remark 7.1. In this case, note that the origin of stochasticity comes only from the initial value $z_0 \sim \text{StdNormal}_d$, and the **stochastic behavior of $\{z_k\}$ depends only on the behavior of z_0** .

- **Generation by constructing Markovian conditional probability densities (Markovian generative transitions):** Choose a **joint probability density**

$$\tilde{q}_{0:K}(\zeta_0, \dots, \zeta_K) \quad (26)$$

such that each marginal \tilde{q}_k **matches** q_k , and the **conditional** $\tilde{q}_{k+1|k}$ is **easy to construct** (specifically, using a **discrete-time diffusion process**). Then, introduce a standard multivariate normal $u_k \sim \text{StdNormal}_d$,

$$z_{k+1} = \mathcal{G}_k(z_k, \hat{\epsilon}_\theta, c, t_k, h_k, u_k) \quad (27)$$

and design the coefficients and variance terms so that the probability distribution of z_{k+1} (given z_k) **matches** the **designed** $\tilde{q}_{k+1|k}$.

These two strategies are complementary. The former provides **deterministic scheduling** based on moment matching and local error analysis from the **push-forward** perspective, while the latter systematically derives **stochastic updates** (such as ancestral steps) that match the **conditional distribution**. In the following, we first provide the concretization (coefficient design) of the deterministic update, and then proceed to the design of the Markovian update.

8 Deterministic Scheduling

8.1 Revisiting the Goal

The **goal** of this section is to use the trained noise estimator $\hat{\epsilon}_\theta$ to construct a **deterministic** recurrence relation

$$z_{k+1} = \mathcal{F}_k(z_0, \dots, z_k; \hat{\epsilon}_\theta, c, t_k, h_k) \quad (28)$$

and design the coefficients and discrete width h_k such that its push-forward distribution $\mathcal{L}(z_{k+1})$ **approaches** q_{k+1} . Here q_k is the artificial distribution sequence defined in (10)–(11).

8.2 Abstract Definition of One-Step Method

We define the **log-SNR** as

$$\lambda_k := \log\left(\frac{\sqrt{\alpha_k}}{\sqrt{1 - \bar{\alpha}_k}}\right), \quad h_k := \lambda_{k+1} - \lambda_k \quad (29)$$

and call h_k the **discrete width**.

Definition 8.1 (One-Step Method Update Upd1 _{$\hat{\epsilon}$,coeff}). Given an arbitrary sequence of coefficients $\{a_k, b_k\}_{k=0}^{K-1} \subset \mathbb{R}$. The **one-step method** is a deterministic recurrence relation determined by a linear combination of z_k and the noise estimate $\hat{\epsilon}_\theta(z_k, c, t_k)$ as

$$z_{k+1} := a_k z_k + b_k \hat{\epsilon}_\theta(z_k, c, t_k). \quad (30)$$

Remark 8.1. Through training, $\hat{\epsilon}_\theta(z, c, t_k)$ is proportional to the **score** of q_k :

$$\nabla_z \log q_k(z) \approx -\frac{1}{\sqrt{1 - \bar{\alpha}_k}} \hat{\epsilon}_\theta(z, c, t_k) \quad (31)$$

(see [5]).

8.3 First-Order Expansion of KL Divergence for One-Step Method

First, let's set up the notation. We write (30) as

$$z_{k+1} = z_k + h_k v_{\text{upd1}}(z_k, \lambda_k), \quad v_{\text{upd1}}(z, \lambda_k) := \frac{a_k - 1}{h_k} z + \frac{b_k}{h_k} \hat{\epsilon}_\theta(z, c, t_k). \quad (32)$$

The **ideal** continuous-time connection q_λ between q_k and q_{k+1} follows the **continuity equation**:

$$\partial_\lambda q_\lambda(z) = -\nabla_z \cdot (q_\lambda(z) v_\star(z, \lambda)), \quad (33)$$

where $v_\star(z, \lambda) = A(\lambda)z + B(\lambda)\nabla_z \log q_\lambda(z)$ is the deterministic drift consistent with q_λ ([2, 5]).

Definition 8.2 (Updated Distribution and KL Divergence). We denote the distribution of z_{k+1} resulting from the update in Definition 8.1 as $p_{k+1}^{(1)}$, and its divergence as

$$D_{\text{KL}}\left(p_{k+1}^{(1)} \parallel q_{k+1}\right). \quad (34)$$

The following proposition rigorously expands (34) in the limit of small h_k .

Proposition 8.1 (First-Order Expansion of KL Divergence for One-Step Method). Assumptions:

- q_λ is C^2 with respect to λ and C^3 with respect to z , and $\|\nabla_z^j \log q_\lambda(z)\|$ is $L^1(\mathbb{R}^d)$ integrable and decays sufficiently for $j \leq 3$.
- $\hat{\epsilon}_\theta(z, c, t)$ is C^2 with respect to z and has at most polynomial growth.
- h_k is sufficiently small.

Then

$$D_{\text{KL}}\left(p_{k+1}^{(1)} \parallel q_{k+1}\right) = \frac{h_k^2}{2} \mathbb{E}_{z \sim q_k} \left[\|\mathbf{v}_{\text{upd1}}(z, \lambda_k) - \mathbf{v}_*(z, \lambda_k)\|_2^2 \right] + O(h_k^3). \quad (35)$$

Proof. (1) **First-order expansion of q_{k+1} :** Discretizing (33) with forward Euler in λ ,

$$q_{k+1}(z) = q_k(z) - h_k \nabla_z \cdot (q_k(z) \mathbf{v}_*(z, \lambda_k)) + O(h_k^2). \quad (36)$$

(2) **First-order expansion of the push-forward distribution $p_{k+1}^{(1)}$:** Consider the transformation $\Phi(z) = z + h_k \mathbf{v}_{\text{upd1}}(z, \lambda_k)$. For small h_k , Φ is a diffeomorphism, and its inverse map Φ^{-1} can be Taylor expanded as

$$\Phi^{-1}(\mathbf{y}) = \mathbf{y} - h_k \mathbf{v}_{\text{upd1}}(\mathbf{y}, \lambda_k) + O(h_k^2). \quad (37)$$

The Jacobian is

$$\begin{aligned} \det(\nabla \Phi^{-1}(\mathbf{y})) &= \det(\mathbf{I} - h_k \nabla_{\mathbf{y}} \mathbf{v}_{\text{upd1}}(\mathbf{y}, \lambda_k)) + O(h_k^2) \\ &= 1 - h_k \text{tr}(\nabla_{\mathbf{y}} \mathbf{v}_{\text{upd1}}(\mathbf{y}, \lambda_k)) + O(h_k^2). \end{aligned} \quad (38)$$

Therefore

$$\begin{aligned} p_{k+1}^{(1)}(\mathbf{y}) &= q_k(\Phi^{-1}(\mathbf{y})) \det(\nabla \Phi^{-1}(\mathbf{y})) \\ &= \left(q_k(\mathbf{y}) - h_k \nabla_{\mathbf{y}} q_k(\mathbf{y}) \cdot \mathbf{v}_{\text{upd1}}(\mathbf{y}, \lambda_k) + O(h_k^2) \right) \\ &\quad \times \left(1 - h_k \nabla_{\mathbf{y}} \cdot \mathbf{v}_{\text{upd1}}(\mathbf{y}, \lambda_k) + O(h_k^2) \right) \\ &= q_k(\mathbf{y}) - h_k (\nabla_{\mathbf{y}} q_k \cdot \mathbf{v}_{\text{upd1}} + q_k \nabla_{\mathbf{y}} \cdot \mathbf{v}_{\text{upd1}}) + O(h_k^2) \\ &= q_k(\mathbf{y}) - h_k \nabla_{\mathbf{y}} \cdot (q_k(\mathbf{y}) \mathbf{v}_{\text{upd1}}(\mathbf{y}, \lambda_k)) + O(h_k^2), \end{aligned} \quad (39)$$

where the last equality uses the product rule for derivatives $\nabla(q_k) \cdot \mathbf{v} + q_k \nabla \cdot \mathbf{v} = \nabla \cdot (q_k \mathbf{v})$.

(3) First-order expression of the difference: The difference between (36) and (39)

$$\delta(\mathbf{y}) := p_{k+1}^{(1)}(\mathbf{y}) - q_{k+1}(\mathbf{y}) = -h_k \nabla_{\mathbf{y}} \cdot (q_k(\mathbf{y}) [\mathbf{v}_{\text{upd1}}(\mathbf{y}) - \mathbf{v}_*(\mathbf{y})]) + O(h_k^2). \quad (40)$$

(4) Quadratic expansion of KL: For $p = q + \delta$,

$$\begin{aligned} D_{\text{KL}}(p\|q) &= \int_{\mathbb{R}^d} (q + \delta) \log\left(1 + \frac{\delta}{q}\right) d\mathbf{y} \\ &= \int (q + \delta) \left(\frac{\delta}{q} - \frac{1}{2} \frac{\delta^2}{q^2}\right) d\mathbf{y} + O\left(\int \frac{|\delta|^3}{q^2} d\mathbf{y}\right) \\ &= \int \frac{\delta^2}{2q} d\mathbf{y} + O(\|\delta\|_{L^3}^3). \end{aligned} \quad (41)$$

By assumption $\delta = O(h_k)$, so $O(\|\delta\|_{L^3}^3) = O(h_k^3)$.

(5) Evaluation of the principal term by integration by parts: Substituting (40) into (41) and sequentially applying integration by parts, using the boundary term being 0 (decay assumption),

$$\begin{aligned} \int \frac{\delta^2}{2q_{k+1}} d\mathbf{y} &= \frac{h_k^2}{2} \int \frac{[\nabla \cdot (q_k(\mathbf{v}_{\text{upd1}} - \mathbf{v}_*))]^2}{q_k} d\mathbf{y} + O(h_k^3) \\ &= \frac{h_k^2}{2} \int q_k \|\mathbf{v}_{\text{upd1}} - \mathbf{v}_*\|_2^2 d\mathbf{y} + O(h_k^3), \end{aligned} \quad (42)$$

where the last equality is obtained by successively applying integration by parts identities (Fisher information type reduction) that pair gradient and divergence terms, in line with the equality condition of the Picard inequality in the $L^2(q_k)$ inner product (intermediate calculations involve terms of the form $\int \partial_i(q_k w_i) \partial_j(q_k w_j)/q_k$, which sequentially reduce to $\int q_k w_i w_i$). Substituting this into (41) yields (35). \square

8.4 Optimization of One-Step Method Coefficients (DDIM/Euler)

Theorem 8.1 (Optimal Coefficients for One-Step Method by Minimizing Divergence). Under the assumptions of Proposition 8.1,

$$a_k^* = \sqrt{\frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} + O(h_k^2), \quad (43)$$

$$b_k^* = -\sqrt{\bar{\alpha}_{k+1}} \frac{\sqrt{1 - \bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} + \sqrt{1 - \bar{\alpha}_{k+1}} + O(h_k^2), \quad (44)$$

minimize the principal term of $D_{\text{KL}}(p_{k+1}^{(1)} \| q_{k+1})$.

Proof. **(1) Form of the ideal drift:** From the structure of (33) and (11), the drift consistent

with q_k is

$$\nu_\star(z, \lambda_k) = \alpha'_k z + \beta'_k \nabla_z \log q_k(z), \quad (45)$$

where α'_k, β'_k are smooth functions of $\bar{\alpha}_k$, and to first-order accuracy

$$\alpha'_k = \frac{1}{h_k} \left(\frac{\sqrt{\bar{\alpha}_{k+1}} - \sqrt{\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} \right) + O(h_k), \quad (46)$$

$$\beta'_k = \sqrt{\bar{\alpha}_{k+1}} - \sqrt{\bar{\alpha}_k} + O(h_k^2), \quad (47)$$

(can be derived from the first-order difference of λ).

(2) Substitution of ϵ -approximation for the score: Substituting (31) into (45),

$$\nu_\star(z, \lambda_k) = \alpha'_k z - \frac{\beta'_k}{\sqrt{1 - \bar{\alpha}_k}} \hat{\epsilon}_\theta(z, c, t_k). \quad (48)$$

(3) Separation of the minimization problem: From Proposition 8.1, the minimization of the principal term is equivalent to

$$\min_{a_k, b_k} \mathbb{E}_{q_k} \left[\left\| \frac{a_k - 1}{h_k} z + \frac{b_k}{h_k} \hat{\epsilon}_\theta - \alpha'_k z + \frac{\beta'_k}{\sqrt{1 - \bar{\alpha}_k}} \hat{\epsilon}_\theta \right\|_2^2 \right]. \quad (49)$$

Although z and $\hat{\epsilon}_\theta$ are not independent in $L^2(q_k)$, from coefficient comparison (uniqueness of least squares in any vector space), the minimum is achieved by matching the coefficients of each basis direction:

$$\frac{a_k - 1}{h_k} = \alpha'_k + O(h_k), \quad (50)$$

$$\frac{b_k}{h_k} = -\frac{\beta'_k}{\sqrt{1 - \bar{\alpha}_k}} + O(h_k). \quad (51)$$

(4) Substitution and simplification of (46)–(47): From (50)

$$\begin{aligned} a_k &= 1 + h_k \alpha'_k + O(h_k^2) = 1 + \left(\frac{\sqrt{\bar{\alpha}_{k+1}} - \sqrt{\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} \right) + O(h_k^2) \\ &= \frac{\sqrt{\bar{\alpha}_{k+1}}}{\sqrt{\bar{\alpha}_k}} + O(h_k^2) = \sqrt{\frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} + O(h_k^2), \end{aligned} \quad (52)$$

From (51)

$$\begin{aligned}
b_k &= -h_k \frac{\beta'_k}{\sqrt{1-\bar{\alpha}_k}} + O(h_k^2) = -\frac{\sqrt{\bar{\alpha}_{k+1}} - \sqrt{\bar{\alpha}_k}}{\sqrt{1-\bar{\alpha}_k}} + O(h_k^2) \\
&= -\sqrt{\bar{\alpha}_{k+1}} \frac{1}{\sqrt{1-\bar{\alpha}_k}} + \frac{\sqrt{\bar{\alpha}_k}}{\sqrt{1-\bar{\alpha}_k}} + O(h_k^2) \\
&= -\sqrt{\bar{\alpha}_{k+1}} \frac{\sqrt{1-\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} \cdot \frac{1}{1-\bar{\alpha}_k} + \sqrt{1-\bar{\alpha}_k} \cdot \frac{1}{1-\bar{\alpha}_k} + O(h_k^2) \\
&= -\sqrt{\bar{\alpha}_{k+1}} \frac{\sqrt{1-\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} + \sqrt{1-\bar{\alpha}_k} \frac{1}{1-\bar{\alpha}_k} (1 - \bar{\alpha}_{k+1} + \bar{\alpha}_{k+1} - (1 - \bar{\alpha}_k)) \\
&= -\sqrt{\bar{\alpha}_{k+1}} \frac{\sqrt{1-\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} + \sqrt{1-\bar{\alpha}_{k+1}} + O(h_k^2),
\end{aligned} \tag{53}$$

where the inserted term in the fourth line is an equivalent transformation to match up to the first order of h_k using the identity $1 = (1 - \bar{\alpha}_{k+1}) + (\bar{\alpha}_{k+1} - (1 - \bar{\alpha}_k))$. (52) and (53) give (43)–(44). \square

Remark 8.2. Substituting (43), (44) back into (30) gives

$$z_{k+1} = \sqrt{\frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} z_k + \left(-\sqrt{\bar{\alpha}_{k+1}} \frac{\sqrt{1-\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} + \sqrt{1-\bar{\alpha}_{k+1}} \right) \hat{\epsilon}_\theta(z_k, c, t_k), \tag{54}$$

which matches the deterministic update equation of **DDIM** widely used in implementation.

Example 8.1 (Euler ODE Scheduler). Discretizing (48) with forward Euler in λ gives

$$z_{k+1} = z_k + h_k \left(\alpha'_k z_k - \frac{\beta'_k}{\sqrt{1-\bar{\alpha}_k}} \hat{\epsilon}_\theta(z_k, c, t_k) \right). \tag{55}$$

This matches the coefficients of (32) and satisfies the principal term minimization condition of Proposition 8.1.

8.5 Abstract Definition of Two-Step Method

Definition 8.3 (Two-Step Method Update Upd2 _{$\hat{\epsilon}$,coeff}). Given arbitrary coefficients $\{a_k, b_k^{(0)}, b_k^{(1)}\} \subset \mathbb{R}$, we define

$$z_{k+1} := a_k z_k + b_k^{(0)} \hat{\epsilon}_\theta(z_k, c, t_k) + b_k^{(1)} \left(\hat{\epsilon}_\theta(z_k, c, t_k) - \hat{\epsilon}_\theta(z_{k-1}, c, t_{k-1}) \right). \tag{56}$$

Proposition 8.2 (Second-Order Accuracy of KL Divergence for Two-Step Method). Under similar regularity assumptions as the one-step method,

$$D_{KL}\left(p_{k+1}^{(2)} \| q_{k+1}\right) = \frac{h_k^3}{2} \mathbb{E}_{q_k} [\|\mathbf{a}_2(z)\|_2^2] + O(h_k^4), \tag{57}$$

holds. Here $p_{k+1}^{(2)}$ is the push-forward distribution of (56), and $\mathbf{a}_2(z)$ is the difference in the **second-order local coefficients** of $\mathbf{v}_{\text{upd2}} - \mathbf{v}_\star$.

Proof. (1) **Expand the two-step method into a second-order local map:** View (56) as

$$z_{k+1} = z_k + h_k \mathbf{v}_{\text{upd2}}(z_k) + \frac{h_k^2}{2} \mathbf{w}_{\text{upd2}}(z_k) + O(h_k^3) \quad (58)$$

(approximating the second-order term of the finite difference with the difference of $\hat{\epsilon}$), and on the q_{k+1} side, expand (33) to second order (equivalent to Heun/Adams–Bashforth)

$$q_{k+1} = q_k - h_k \nabla \cdot (q_k \mathbf{v}_\star) + \frac{h_k^2}{2} \Xi(q_k, \mathbf{v}_\star) + O(h_k^3), \quad (59)$$

where Ξ is a linear operator for the second-order term involving the Jacobian and density gradient.

(2) **Second-order Taylor of change of variables:** Let $\Phi(z) = z + h_k \mathbf{v}_{\text{upd2}} + \frac{h_k^2}{2} \mathbf{w}_{\text{upd2}}$, and expand Φ^{-1} and $\det(\nabla \Phi^{-1})$ to second order. Simplifying similarly to the one-step method,

$$p_{k+1}^{(2)} = q_k - h_k \nabla \cdot (q_k \mathbf{v}_{\text{upd2}}) + \frac{h_k^2}{2} \left(\Xi(q_k, \mathbf{v}_{\text{upd2}}) - \nabla \cdot (q_k \mathbf{w}_{\text{upd2}}) \right) + O(h_k^3). \quad (60)$$

(3) **Difference and KL:** Write the difference between (59) and (60) as $\delta = -h_k \nabla \cdot (q_k (\mathbf{v}_{\text{upd2}} - \mathbf{v}_\star)) + \frac{h_k^2}{2} (\Xi(q_k, \mathbf{v}_{\text{upd2}}) - \Xi(q_k, \mathbf{v}_\star) - \nabla \cdot (q_k \mathbf{w}_{\text{upd2}})) + O(h_k^3)$, and apply a quadratic expansion similar to (41). Due to the matching of coefficients up to second order, the first-order contribution vanishes, and the principal term is of order h_k^3 , $\propto \int q_k \|\mathbf{a}_2(z)\|^2$. This gives (57). \square

Theorem 8.2 (Optimal Coefficients for Two-Step Method by Minimizing Divergence (DPM++ 2M form)). The coefficients that minimize D_{KL} to second-order accuracy are given by

$$a_k^\star = \sqrt{\frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} + O(h_k^3), \quad (61)$$

$$b_k^{(0)\star} = -\frac{\sqrt{\bar{\alpha}_{k+1}}}{\sqrt{\bar{\alpha}_k}} \phi_1(h_k) + \sqrt{1 - \bar{\alpha}_{k+1}}, \quad \phi_1(h) := \frac{e^h - 1}{h}, \quad (62)$$

$$b_k^{(1)\star} = -\frac{\sqrt{\bar{\alpha}_{k+1}}}{\sqrt{\bar{\alpha}_k}} \phi_2(h_k), \quad \phi_2(h) := \frac{e^h - 1 - h}{h^2}, \quad (63)$$

Proof. (1) **Exponential integrator form of the probability flow ODE:** Using (48), $\partial_\lambda z = \alpha'(\lambda)z - \gamma(\lambda)\hat{\epsilon}_\theta(z, \mathbf{c}, t(\lambda))$ ($\gamma(\lambda) := \beta'(\lambda)/\sqrt{1 - \bar{\alpha}(\lambda)}$). Approximating this with a two-step exponential integrator ϕ -function, the local transition is

$$z_{k+1} = e^{h_k} z_k - e^{h_k} \phi_1(h_k) \hat{\gamma}_k \hat{\epsilon}_\theta(z_k) - e^{h_k} \phi_2(h_k) \hat{\gamma}_k \left(\hat{\epsilon}_\theta(\tilde{z}_{k+1}) - \hat{\epsilon}_\theta(z_k) \right) + O(h_k^3), \quad (64)$$

where $\hat{\gamma}_k \approx \frac{\sqrt{\bar{\alpha}_{k+1}}}{\sqrt{\bar{\alpha}_k}} - 1 = O(h_k)$, $e^{h_k} = \sqrt{\bar{\alpha}_{k+1}/\bar{\alpha}_k}$ (from the definition of λ).

(2) Coefficient identification: Comparing coefficients between (64) and (56), and substituting $e^{h_k} = \sqrt{\bar{\alpha}_{k+1}/\bar{\alpha}_k}$, we obtain (61)–(63). The cancellation of the second-order error satisfies the condition $\mathbf{a}_2(z) \equiv 0$ of Proposition 8.2. \square

Example 8.2 (Explicit formula for DPM++ 2M Karras). The implementation is given in the following two stages:

$$\tilde{z}_{k+1} = \sqrt{\frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} z_k - \frac{\sqrt{\bar{\alpha}_{k+1}}}{\sqrt{\bar{\alpha}_k}} \phi_1(h_k) \hat{\epsilon}_\theta(z_k, \mathbf{c}, t_k) + \sqrt{1 - \bar{\alpha}_{k+1}} \hat{\epsilon}_\theta(z_k, \mathbf{c}, t_k), \quad (65)$$

$$\begin{aligned} z_{k+1} &= \sqrt{\frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} z_k - \frac{\sqrt{\bar{\alpha}_{k+1}}}{\sqrt{\bar{\alpha}_k}} \left(\phi_1(h_k) \hat{\epsilon}_\theta(z_k, \mathbf{c}, t_k) + \phi_2(h_k) [\hat{\epsilon}_\theta(\tilde{z}_{k+1}, \mathbf{c}, t_{k+1}) - \hat{\epsilon}_\theta(z_k, \mathbf{c}, t_k)] \right) \\ &\quad + \sqrt{1 - \bar{\alpha}_{k+1}} \hat{\epsilon}_\theta(\tilde{z}_{k+1}, \mathbf{c}, t_{k+1}), \end{aligned} \quad (66)$$

$h_k = \lambda_{k+1} - \lambda_k$. This matches the coefficients of Theorem 8.2 and achieves a local error of $O(h_k^3)$ [2, 3].

9 Markovian Scheduling

9.1 Revisiting the Goal

The **goal** of this section is to choose a **joint probability density** $\tilde{q}_{0:K}$ such that each marginal \tilde{q}_k matches q_k , and the conditional $\tilde{q}_{k+1|k}$ is easy to construct (discrete-time diffusion process), and to design

$$z_{k+1} = \mathcal{G}_k(z_k, \hat{\epsilon}_\theta, \mathbf{c}, t_k, h_k, \mathbf{u}_k), \quad \mathbf{u}_k \sim \text{StdNormal}_d, \quad (67)$$

such that the distribution of z_{k+1} given z_k **matches** $\tilde{q}_{k+1|k}$.

9.2 Local Linear Noising and Gaussian Approximation of Reverse Conditional

Proposition 9.1 (Gaussian Approximation of Reverse Conditional in Local Linear Noising). Define a sequence of random variables $\{\xi_k\}$ by

$$\xi_k = \lambda_{\text{signal}} \xi_{k+1} + \delta \epsilon_k, \quad \epsilon_k \sim \text{StdNormal}_d, \text{ independent.} \quad (68)$$

When $\delta > 0$ is sufficiently small, even if ξ_k is non-Gaussian, the reverse conditional distribution $p(\xi_{k+1} | \xi_k)$ is

$$p(\xi_{k+1} | \xi_k) = \mathcal{N}(\mathbf{M}\xi_k, \Sigma) + O(\delta^2), \quad (69)$$

(where \mathbf{M}, Σ are first-order functions of $\lambda_{\text{signal}}, \delta$).

Proof. (1) **Application of Bayes' rule:** $p(\xi_{k+1} \mid \xi_k) \propto p(\xi_k \mid \xi_{k+1}) p(\xi_{k+1})$. Here $p(\xi_k \mid \xi_{k+1}) = \mathcal{N}(\lambda_{\text{signal}} \xi_{k+1}, \delta^2 \mathbf{I})$ is exactly Gaussian.

(2) **Laplace approximation:** The log-posterior $\ell(\mathbf{u}) := \log p(\mathbf{u} \mid \xi_k)$ with $\mathbf{u} = \xi_{k+1}$ is

$$\ell(\mathbf{u}) = -\frac{1}{2\delta^2} \|\xi_k - \lambda_{\text{signal}} \mathbf{u}\|_2^2 + \log p(\mathbf{u}) + \text{const.} \quad (70)$$

The gradient and Hessian are

$$\nabla_{\mathbf{u}} \ell(\mathbf{u}) = \frac{\lambda_{\text{signal}}}{\delta^2} (\xi_k - \lambda_{\text{signal}} \mathbf{u}) + \nabla_{\mathbf{u}} \log p(\mathbf{u}), \quad (71)$$

$$\nabla_{\mathbf{u}}^2 \ell(\mathbf{u}) = -\frac{\lambda_{\text{signal}}^2}{\delta^2} \mathbf{I} + \nabla_{\mathbf{u}}^2 \log p(\mathbf{u}). \quad (72)$$

As $\delta \rightarrow 0$, the first term dominates.

(3) **First-order approximation of the mode:** From $\nabla \ell(\mathbf{u}^*) = 0$, $\frac{\lambda_{\text{signal}}}{\delta^2} (\xi_k - \lambda_{\text{signal}} \mathbf{u}^*) = -\nabla \log p(\mathbf{u}^*)$. The right side is $\mathcal{O}(1)$, the left side is $\mathcal{O}(\delta^{-2})$, so $\mathbf{u}^* = \lambda_{\text{signal}}^{-1} \xi_k + \mathcal{O}(\delta^2)$.

(4) **Quadratic approximation and covariance:** The value of (72) at the mode is $-\nabla^2 \ell(\mathbf{u}^*) = \frac{\lambda_{\text{signal}}^2}{\delta^2} \mathbf{I} - \nabla^2 \log p(\mathbf{u}^*)$. The principal term is $\frac{\lambda_{\text{signal}}^2}{\delta^2} \mathbf{I}$, so $\Sigma = \left(\frac{\lambda_{\text{signal}}^2}{\delta^2} \mathbf{I} \right)^{-1} + \mathcal{O}(\delta^2) = \frac{\delta^2}{\lambda_{\text{signal}}^2} \mathbf{I} + \mathcal{O}(\delta^4)$. The mean is $\mathbf{M} \xi_k = \mathbf{u}^* + \mathcal{O}(\delta^2) = \lambda_{\text{signal}}^{-1} \xi_k + \mathcal{O}(\delta^2)$. \square

9.3 \tilde{q} as a Discrete-Time Diffusion Process

Definition 9.1 (Joint Distribution of Forward Discrete Diffusion). In d dimensions, let

$$\xi_K \sim q_K, \quad \epsilon_k \sim \text{StdNormal}_d \text{ independent}, \quad (73)$$

$$\alpha_k := \frac{\bar{\alpha}_{k-1}}{\alpha_k}, \quad \xi_{k-1} = \sqrt{\alpha_k} \xi_k + \sqrt{1 - \alpha_k} \epsilon_k, \quad (k = 1, \dots, K), \quad (74)$$

define the joint distribution $\tilde{q}_{0:K}$.

Theorem 9.1 (Identity $\tilde{q}_k = q_k$). Under Definition 9.1, for all k ,

$$\tilde{q}_k = q_k \quad (75)$$

holds.

Proof. (1) $k = K$: By definition $\tilde{q}_K = q_K$.

(2) **Induction:** Assume $\tilde{q}_k = q_k$ and show $\tilde{q}_{k-1} = q_{k-1}$. From (74) $\sqrt{\bar{\alpha}_{k-1}} \xi_{k-1} = \sqrt{\alpha_k} \xi_k + \sqrt{\bar{\alpha}_{k-1} - \bar{\alpha}_k} \epsilon_k$. By assumption $\xi_k \stackrel{d}{=} \sqrt{\alpha_k} \mathbf{x} + \sqrt{1 - \alpha_k} \epsilon$. By addition of independent normals $\sqrt{\alpha_k} \xi_k + \sqrt{\bar{\alpha}_{k-1} - \bar{\alpha}_k} \epsilon_k \stackrel{d}{=} \sqrt{\alpha_k} \mathbf{x} + \sqrt{1 - \alpha_{k-1}} \epsilon'$ (additivity of variance and independence). Thus $\xi_{k-1} \stackrel{d}{=} \sqrt{\bar{\alpha}_{k-1}} \mathbf{x} + \sqrt{1 - \bar{\alpha}_{k-1}} \epsilon'$, i.e., $\tilde{q}_{k-1} = q_{k-1}$. \square

9.4 Abstract Definition of Markovian Update and Divergence Expansion

Definition 9.2 (Markovian Update MUpd _{$\hat{\epsilon}$,coeff}). Given arbitrary coefficients $\{A_k, B_k, C_k\} \subset \mathbb{R}$, define

$$z_{k+1} := A_k z_k + B_k \hat{\epsilon}_{\theta}(z_k, c, t_k) + C_k u_k, \quad u_k \sim \text{StdNormal}_d. \quad (76)$$

Proposition 9.2 (First-Order Expansion of Conditional Divergence). Let $\tilde{q}_{k+1|k}$ be the exact posterior (Gaussian) from Definition 9.1. Let the distribution of z_{k+1} conditioned on $z_k = z$ be $p_{k+1|k}^{\text{cond}}(\cdot | z)$. Then

$$\begin{aligned} \mathbb{E}_{z \sim q_k} \left[\text{D}_{\text{KL}}(p_{k+1|k}^{\text{cond}}(\cdot | z) \| \tilde{q}_{k+1|k}(\cdot | z)) \right] &= \frac{1}{2} \mathbb{E}_{q_k} \left[\frac{\|\mu_{\text{mupd}}(z) - \mu_*(z)\|_2^2}{\sigma_*^2} + \frac{(\sigma_{\text{mupd}} - \sigma_*)^2}{\sigma_*^2} \right] \\ &\quad + O(h_k^2), \end{aligned} \quad (77)$$

where

$$\mu_{\text{mupd}}(z) = A_k z + B_k \hat{\epsilon}_{\theta}(z, c, t_k), \quad \sigma_{\text{mupd}}^2 = C_k^2, \quad (78)$$

$$\mu_*(z) = \sqrt{\frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} z - \sqrt{\bar{\alpha}_{k+1}} \frac{\sqrt{1-\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} \hat{\epsilon}_{\theta}(z, c, t_k), \quad (79)$$

$$\sigma_*^2 = 1 - \frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}. \quad (80)$$

Proof. **(1) Mean and variance of the posterior Gaussian:** In the linear Gaussian process of Definition 9.1, from standard conditional Gaussian formulas $\tilde{q}_{k+1|k} = \mathcal{N}(\mu_*(z), \sigma_*^2 I)$, and μ_*, σ_*^2 are given by (79)–(80) (identical to the derivation in [1]).

(2) Conditional of the generative side: From (76) $p_{k+1|k}^{\text{cond}}(\cdot | z) = \mathcal{N}(\mu_{\text{mupd}}(z), \sigma_{\text{mupd}}^2 I)$.

(3) Exact formula for KL between isotropic Gaussians: For isotropic Gaussians with means μ_1, μ_2 and variances σ_1^2, σ_2^2 ,

$$\text{D}_{\text{KL}}(\mathcal{N}(\mu_1, \sigma_1^2 I) \| \mathcal{N}(\mu_2, \sigma_2^2 I)) = \frac{1}{2} \left(\frac{\|\mu_1 - \mu_2\|_2^2}{\sigma_2^2} + d \left(\frac{\sigma_1^2}{\sigma_2^2} - 1 - \log \frac{\sigma_1^2}{\sigma_2^2} \right) \right). \quad (81)$$

Here d is the dimension.

(4) Small-step approximation: Since $\bar{\alpha}_{k+1}/\bar{\alpha}_k = 1 - O(h_k)$, $\sigma_*^2 = 1 - \bar{\alpha}_{k+1}/\bar{\alpha}_k = O(h_k)$. Thus $\frac{\sigma_{\text{mupd}}^2}{\sigma_*^2} - 1 - \log \frac{\sigma_{\text{mupd}}^2}{\sigma_*^2} = \frac{(\sigma_{\text{mupd}} - \sigma_*)^2}{\sigma_*^2} + O(h_k^2)$ quadratic expansion applies. Substitute $\mu_1 = \mu_{\text{mupd}}$, $\mu_2 = \mu_*$, $\sigma_1 = \sigma_{\text{mupd}}$, $\sigma_2 = \sigma_*$ into (81) for this problem, and take the q_k expectation to obtain (77). \square

9.5 Optimal Coefficients and DDPM/Euler a

Theorem 9.2 (Optimal Coefficients by Minimizing Divergence). The coefficients that minimize the principal term on the right-hand side of Proposition 9.2 are

$$A_k^* = \sqrt{\frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} + O(h_k^2), \quad (82)$$

$$B_k^* = -\sqrt{\bar{\alpha}_{k+1}} \frac{\sqrt{1-\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} + O(h_k^2), \quad (83)$$

$$C_k^* = \sqrt{1 - \frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} + O(h_k^2). \quad (84)$$

Proof. (1) **Point-wise minimization:** The integrand of (77) for each z is $\frac{\|\mu_{\text{mupd}}(z) - \mu_*(z)\|^2}{2\sigma_*^2} + \frac{(\sigma_{\text{mupd}} - \sigma_*)^2}{2\sigma_*^2}$. To minimize this, it suffices to satisfy $\mu_{\text{mupd}}(z) = \mu_*(z)$, $\sigma_{\text{mupd}} = \sigma_*$.

(2) **Coefficient matching:** Comparing (78) and (79), $A_k^* = \sqrt{\bar{\alpha}_{k+1}/\bar{\alpha}_k}$, $B_k^* = -\sqrt{\bar{\alpha}_{k+1}} \frac{\sqrt{1-\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}}$, $C_k^* = \sqrt{1 - \bar{\alpha}_{k+1}/\bar{\alpha}_k}$ are obtained. The $O(h_k^2)$ residuals are due to the first-order smoothness of $\bar{\alpha}$. \square

Example 9.1 (DDPM and Euler a). Substituting Theorem 9.2 into (76)

$$z_{k+1} = \sqrt{\frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} z_k - \sqrt{\bar{\alpha}_{k+1}} \frac{\sqrt{1-\bar{\alpha}_k}}{\sqrt{\bar{\alpha}_k}} \hat{\epsilon}_\theta(z_k, c, t_k) + \sqrt{1 - \frac{\bar{\alpha}_{k+1}}{\bar{\alpha}_k}} u_k, \quad (85)$$

is obtained. This matches the **DDPM ancestral sample (Euler a)** [1], and makes the principal term of Proposition 9.2 zero point-wise (i.e., first-order consistency).

10 Summary and Next Time

10.1 Summary Corresponding to Learning Outcomes

- **Data Acquisition for Implicit Distribution Learning:** We constructed an **artificial distribution sequence** q_0, \dots, q_K and corresponding data through forward noising.
- **Role of Noise Estimator (Reconciliation):** Learning via the objective function achieved the **reconciliation** of two difficult conditions: being learnable from real data and realizing the target distribution via reverse diffusion.
- **Proximity of Scheduler:** We showed that by determining the coefficients of deterministic (DDIM/Euler) and stochastic (DDPM/Euler a) updates through moment matching, we **approach** q_{k+1} at each update.

10.2 Next Time

Next time, we will discuss **convolutional neural networks** used in image generation AI from an implementation and design perspective (centering on U-Net [4]).

References

- [1] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In NeurIPS, 2020.
- [2] Tero Karras, Miika Aittala, Samuli Laine, Ari Herva, and Jaakko Lehtinen. Elucidating the design space of diffusion-based generative models. arXiv:2206.00364, 2022.
- [3] Cheng Lu and Yuhao Zhou. Dpm-solver++: Fast solver for diffusion odes with optimal error bounds. arXiv preprint arXiv:2211.01095, 2022.
- [4] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In MICCAI, 2015.
- [5] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456, 2021.