

Topics on Diffusion Generative Models

Original Text in Chinese by Jianlin Su¹ and Translated and Annotated by Jack Yang²

¹⁾ <https://kexue.fm>

²⁾ *School of Materials Science and Engineering, University of New South Wales, Sydney, New South Wales 2052, Australia^{a)}*

(Dated: 7 July 2025)

Diffusion generative models, originally developed for high-quality image synthesis, have gradually gained popularity in other scientific research domains, including materials science for innovative materials discoveries. Compared with other deep learning models, however, diffusion models are much more mathematically involved, thus more challenging to be understood by novelists. Nevertheless, the ideal of diffusion generative models have strong link with the diffusion processes, which can be described by stochastic diffusion equations and are very known across many fields of fundamental science. In light of these, the translator (JY) sees huge potential in the applications of both the theoretical frameworks and the numerical implementations of the diffusion generative models in the future. Throughout the past few years, the original author (JS) has written a series of blog articles on diffusion generative models in Chinese, providing a more approachable understanding on the mathematical motivations and derivations behind this family of generative models, of which JY believes an English translated version will bring greater benefit to the large community who are interested in utilising these models in their research and development. This document is, therefore, prepared for this purpose while JY was studying JS's blog articles. Some further annotations and notes are added by JY to facilitate the understanding of the original materials.

^{a)} Electronic mail: jianliang.yang1@unsw.edu.au

1. DDPM=Building Demolition and Rebuilding	4
1.1. A new starting point	4
1.2. Demolition and rebuilding	4
1.3. How should we demolish the building?	5
1.4. How to rebuild the building?	6
1.5. Reducing the square error	6
1.6. Regressive generations	7
1.7. Hyperparameter settings	7
2. DDPM=Autoregressive VAE	8
2.1. Solving the problem step-wise	8
2.2. The joint KL divergence	8
2.3. Divide-and-conquer	9
2.4. Reconstructing the diffusion model	9
2.5. Hyperparameter settings	10
3. DDPM=Bayes Plus Denoising	11
3.1. Recap on the model setup	11
3.2. Using the Bayes' theorem	11
3.3. The denoising process	12
3.4. Correcting the estimates	13
3.5. Remaining questions	13
4. DDIM=High Level View of DDPM	14
4.1. The thought process	14
4.2. Method of undetermined coefficients	14
4.3. Same procedure as before	15
4.4. A few examples	15
4.5. Accelerated generations	16
4.6. Differential equation	17
5. The General Framework Based on Stochastic Differential Equation	17
5.1. Stochastic differentiations	18
5.2. The reverse process	18
5.3. Score matching	19
5.4. The final results	20
6. The General Framework Based on Ordinary Differential Equation	22
6.1. Think again	22
6.2. The Dirac function	22
6.3. The Fokker-Planck equation	23
6.4. Equivalent transformation	23
6.5. Neural ODE	24
6.6. Revisiting DDIM	24
7. Optimum Estimation of the Variance (Part I)	25
7.1. Uncertainty in the reconstruction process	26
7.2. Optimisation for the mean	26
7.3. Estimation of the variance: part I	27
7.4. Estimation of the variance: part II	28
7.5. Further thoughts	28
8. Optimum Estimation of the Variance (Part II)	29
8.1. Revisiting the results	29
8.2. Different approaches to improve the results	29
8.3. Maximum likelihood	30
8.4. Conditional variance	30
8.5. Two stages	31
9. Conditional Generation	31
9.1. Theoretical analysis	31
9.2. Conditional input	32
9.3. Approximated distribution	32
9.4. Scaling the classifier gradient	33
9.5. Similarity control	33
9.6. Continuous scenarios	34
9.7. Classifier-free approach	35
10. Unified Diffusion Model: The Theoretical Framework	35
10.1. The forward process	36
10.2. The reverse process	36
10.3. Training target	36
10.4. Conditional probability	37

10.5. The thought process	38
11. Unified Diffusion Model: Applications	38
11.1. Recap on the general framework	38
11.2. Hot diffusion	39
11.3. Cold diffusion	40
11.4. The editing model	40
11.5. The masking model	41
11.6. The encoding model	42
12. Conquering the Diffusion ODE	42
12.1. Differential equation	42
12.2. The Jacobian	43
12.3. Taylor approximation	43
12.4. The equation of heat conduction	44
12.5. Solving for the probability distribution	44
12.6. Completing the model design	45
13. From the Law of Gravitational Force to Diffusion Model	45
13.1. The gravitational force	46
13.2. Following the lines of the gravitational field	46
13.3. Equivalent centre-of-mass	46
13.4. Mode Collapsing	46
13.5. Adding one extra dimension	47
13.6. Everything clicked	47
13.7. Training the field	48
14. The General Procedure of Constructing an ODE (Part I)	49
14.1. The basic conclusions from the previous sections	49
14.2. Simplifying the equation	50
14.3. The Green's function	50
14.4. The gravitational force	50
14.4.1. An ansatz	51
14.4.2. Imposing the constraints	51
14.4.3. Result analysis	51
14.4.4. Revisiting the problem	52
14.5. Space-time separation	53
14.5.1. An ansatz	53
14.5.2. Gaussian diffusion	54
14.6. Reverse engineering	54
14.7. Summary	55
15. The General Procedure of Constructing an ODE (Part II)	55
15.1. Revisiting the previous results	55
15.2. Geometric interpretations	56
15.3. The method of characteristics	57
15.4. The training target	57
15.5. A few examples	58
15.5.1. Linear trajectory	58
15.5.2. A further generalisation	58
15.5.3. Even more complicated?	59
15.6. Summary	59
16. Wasserstein Distance \leq Score-Matching	59
16.1. Result analysis	60
16.2. Preparatory works	61
16.3. The initial trial	61
16.4. Continuing the proof	63
16.5. Further generalisation and its difficulties	63
17. The General Procedure of Constructing an ODE (Part II)	64
17.1. The intuitive results	64
17.2. A simple example	65
17.3. Further proofs	66
18. Score-Matching = Conditional Score-Matching	67
18.1. Score-matching	67
18.2. Conditional score-matching	68
18.3. The inequality relationship	68
18.4. The equality relationship	68

1. DDPM=Building Demolition and Rebuilding

Generative Adversarial Networks (GAN) and Variational Autoencoders (VAE) are some of the well-known generative models in the deep learning community. Some of the less popular models, including the flow models¹, VQ-VAE² and so on, have also become increasingly attractive in the past few years. This is particularly true for VQ-VAE and its variant VQ-GAN, which have risen to become the image tokeniser, that can be adapted to (pre-)train the NLP models. Apart from these models, an even less known choice, namely, the **diffusion model**, has become increasingly popular in the area of artificial image generations. The two most successful text-to-image generation models to date, DALL-E2 from OpenAI and Imagen from Google, are both diffusion generative models.

Hereafter, we start to discuss the recent progresses in the development of the diffusion model. Generally, it was believed that the mathematical derivations behind the diffusion models were very complex thus lacking intuitions. Here, we shall try to unpick these complex mathematical derivations, aiming at obtaining an intuitively clearer understanding behind the diffusion model.

1.1. A new starting point

Generally, when discussing the diffusion model, most literatures will refer to concepts such as **the energy-based model**, **score-matching**, **Langevin dynamics**, and so on. In a simple word, what diffusion model is trying to do is to *utilise the technique of score-matching to train an energy model, and then sample new items from this trained model based on the Langevin equation*.

Theoretically, this is a very mature proposal, whereby we can, in principle, apply it to generate and sample any continuous objects such as images and audio. Practically, however, it is rather difficult to train an energy function, especially when the dimensionality of the data is very high (*e.g.* high-resolution images). The main challenge is to obtain a fully converged energy function. On the other hand, there exist lots of uncertainties when we perform samplings based on the energy model using the Langevin equation, which results in noisy samples. As such, for a long time, the traditional ideas behind the diffusion model had only been experimented on generating low-resolution images.

The popularity of modern diffusion generative models started from the **DDPM (Denoising Diffusion Probabilistic Model)**³ first proposed in 2020. Although it had been named as a diffusion model, the formalism of DDPM is completely different from the traditional diffusion model that was based on the Langevin equation. More precisely, this new type of diffusion model should really be named as **models of progressive changes**, whereas calling it the diffusion model can actually cause lots of misunderstanding on the essence of this model. The concepts in the traditional diffusion model, such as the energy function, score matching and Langevin dynamics actually *have nothing to do with* DDPM and its subsequent variants. Interestingly, before 2020, the fundamental mathematical framework for DDPM had already been published in an ICML2015 article⁴ entitled “*Deep Unsupervised Learning Using Nonequilibrium Thermodynamics*”, but the model has only become popular after its capacity in generating high-resolution images was demonstrated. This shows that the birth and popularisation of a new model require both time and opportunity.

1.2. Demolition and rebuilding

When introducing the DDPM, many literatures adapt the probabilistic interpretation and dive straight into the change in the probability distribution that is followed by variational inference, which scares many readers with the mathematical equations. This is further compounded by the perceptions from the traditional diffusion models, thinking that the maths involved are very complicated. As a matter of fact, DDPM can be more intuitively explained, in a way that is no more complicated than the idea of ‘fake it and then discriminate it’ behind GAN (Generative Adversarial Network).

Firstly, let’s think of a GAN-like model, it is actually a model that converts a random noise z into an actual sample x , as shown in Fig. 1.

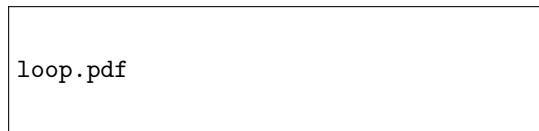


FIG. 1. Building construction analogy to the DDPM model.

Original blog post see <https://www.spaces.ac.cn/archives/9119>

June 2022, at the time of writing.

In statistical thermodynamics, it means we need to sample all possible states $\{x_i\}$ to obtain converged configurational partition function $Z = \sum_{i=1}^N \exp[-U(x_i)/k_B T]$, where $U(x)$ is the energy function, from which we can then gain access to the probability of individual sample $P(x_i) = \exp[-U(x_i)/k_B T]/Z$. However, it is generally impossible to obtain the converged Z through sampling!

Although, on the other hand, this proves that DDPM is a variational autoencoder (VAE) rather than a diffusion model.

We can rethink of the process as a building process, where the noises z are the raw materials, in this case, the bricks and mortar, and the samples x are the final buildings being built. In this case, the generative model becomes the construction team.

The reconstruction process is bound to be difficult, and this is why there are so many researches involved in the development of generative models. On the other hand, the demolition process is very easy. Considering a **step-wise process** that gradually demolishing a building into bricks, in this case, let x_0 be the fully-constructed building (initial sample) and x_T being the fully deconstructed bricks (random noises). Assuming the demolition takes T steps, then the whole demolition process can be written as

$$x = x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \cdots x_{T-1} \rightarrow x_T = z. \quad (1.1)$$

Within this framework, the difficulties in rebuilding a building is that, the jump from the raw materials x_T to the final building x_0 is too large for a novelist to comprehend how this process can be achieved. However, if we have the intermediate states from the demolition processes x_1, x_2 till x_{T-1} , and we know $x_{t-1} \rightarrow x_t$ represents one small demolition step, then why can't we utilise these information to learn the reverse process $x_t \rightarrow x_{t-1}$ to rebuild the building? This means that if we can learn the relationship between the two $\mu(x_t) = x_{t-1}$, then starting from x_T , by repetitively applying $x_{T-1} = \mu(x_T)$, $x_{T-2} = \mu(x_{T-1})$ and so on, wouldn't we be able to rebuild the whole building x_0 ?

1.3. How should we demolish the building?

Just like we need to eat our meals bite by bite, we also need to rebuild our buildings stepwise. The process of generation in DDPM is actually fully consistent with the the **analogy of "demolition followed by reconstruction"** for DDPM. It is a process that first gradually deconstruct the data samples into random noises, then considers a reverse process which is carried out repetitively to (re)generated the (original) new samples. As mentioned before, this is why DDPM is more like a model of **progressive change** rather than a diffusion model.

More specifically, the demolition process in DDPM takes the following form:

$$x_t = \alpha_t x_{t-1} + \beta_t \varepsilon_t, \quad \text{with } \varepsilon_t \sim \mathcal{N}(0, I), \quad (1.2)$$

in particular, α_t and $\beta_t > 0$, $\alpha_t^2 + \beta_t^2 = 1$. Usually, we chose $\beta_t \rightarrow 0$, it represents how much damage we will make to the building during each step of the demolition process. The introduction of the **noise** ε_t represents the destruction to the original signal. We can also think of it as the raw materials, such that in each step of the demolition process, we break x_{t-1} into $\alpha_t x_{t-1}$ part of the original building plus $\beta_t \varepsilon_t$ part of the raw materials.

If we repeat the above process, then we have

$$\begin{aligned} x_t &= \alpha_t x_{t-1} + \beta_t \varepsilon_t \\ &= \alpha_t (\alpha_{t-1} x_{t-2} + \beta_{t-1} \varepsilon_{t-1}) + \beta_t \varepsilon_t \\ &= \cdots \\ &= (\alpha_t \alpha_{t-1} \cdots \alpha_1) x_0 + (\alpha_t \alpha_{t-1} \cdots \alpha_2) \beta_1 \varepsilon_1 + (\alpha_t \cdots \alpha_3) \beta_2 \varepsilon_2 + \cdots + \beta_t \varepsilon_t. \end{aligned} \quad (1.3)$$

So why do we require the coefficients must satisfy the constraint $\alpha_t^2 + \beta_t^2 = 1$? Let us try to answer this question. In Eq. (1.3), the part that is highlighted in red represents the sum of multiple independent Gaussian noises with the mean of zero, and variances being $(\alpha_t \alpha_{t-1} \cdots \alpha_2) \beta_1$, $(\alpha_t \cdots \alpha_3) \beta_2, \dots, \alpha_t \beta_{t-1}^2$ and β_t^2 , respectively. Then, based on the superposition of Gaussians from the probability theory, the sum of Gaussians returns a new Gaussian, which, in the case is of mean zero and variance given by $(\alpha_t \alpha_{t-1} \cdots \alpha_2) \beta_1 + (\alpha_t \cdots \alpha_3) \beta_2 + \cdots + \beta_t$. Providing that $\alpha_t^2 + \beta_t^2 = 1$, it is not difficult to prove that

$$(\alpha_t \alpha_{t-1} \cdots \alpha_2) \beta_1 + (\alpha_t \cdots \alpha_3) \beta_2 + \cdots + \beta_t = 1. \quad (1.4)$$

This means that, effectively, we have

$$x_t = \underbrace{(\alpha_t \cdots \alpha_1)}_{\equiv \bar{\alpha}_t} x_0 + \underbrace{\sqrt{1 - (\alpha_t \cdots \alpha_1)^2}}_{\equiv \bar{\beta}_t} \bar{\varepsilon}_t, \quad \bar{\varepsilon}_t \sim \mathcal{N}(0, I), \quad (1.5)$$

which has now significantly simplified the computations of x_t . On the other hand, by choosing a suitable form of α_t in DDPM to ensure $\bar{\alpha}_t \approx 0$ will guarantee that, after T steps of

It is important to take a note of this at this point, as this is a recurrent trick in the diffusion generative model that tries to learn the reverse construction process. This can be parameterised by any major NN architectures, including convolutional NN, transformer, graph and message-passing NNs for learning atomistic structures. In this regard, diffusion model has very little to do with the designs of the NN but more with the manipulations in the probability distributions of the data.

This part is the sum of multiple independent Gaussian noises.

Take note of the definitions $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ and $\bar{\beta}_t = \sqrt{1 - \bar{\alpha}_t^2}$, which are commonly used in the diffusion models.

demolitions, the entire building will be fully deconstructed into its raw materials ϵ .

1.4. How to rebuild the building?

During the demolition process $\mathbf{x}_{t-1} \rightarrow \mathbf{x}_t$, we can obtain many pairs of data $(\mathbf{x}_{t-1}, \mathbf{x}_t)$, therefore, rebuilding is simply a process that learns a model to enable the reverse process $\mathbf{x}_t \rightarrow \mathbf{x}_{t-1}$ to occur. Suppose this model is $\boldsymbol{\mu}(\mathbf{x}_t)$, then the most straightforward learning protocol is to minimise the Euclidean distance

$$\|\boldsymbol{\mu}(\mathbf{x}_t) - \mathbf{x}_{t-1}\|_2^2. \quad (1.6)$$

This simple formulation is actually very close to the final DDPM model. Here, we can refine it further. If we first rearrange Eq. (1.2) from the demolition process into a reconstructing one by making \mathbf{x}_{t-1} the subject:

$$\mathbf{x}_{t-1} = \frac{1}{\alpha_t} (\mathbf{x}_t - \beta_t \epsilon_t),$$

this then inspires us to express the parameterised reconstruction process analogously as

$$\boldsymbol{\mu}(\mathbf{x}_t) = \frac{1}{\alpha_t} (\mathbf{x}_t - \beta_t \epsilon_\theta(\mathbf{x}_t, t)), \quad (1.7)$$

in which θ is the learnable parameter. Substituting Eq. (1.7) into the loss function Eq. (1.6), we get

$$\|\boldsymbol{\mu}(\mathbf{x}_t) - \mathbf{x}_{t-1}\|_2^2 = \frac{\beta_t^2}{\alpha_t^2} \|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)\|_2^2. \quad (1.8)$$

Here, β_t^2/α_t^2 represents the weight of the loss function, which we can ignore it for the time being. Substituting Eq. (1.2) and Eq. (1.5) for the expression of \mathbf{x}_t :

$$\mathbf{x}_t = \alpha_t \mathbf{x}_{t-1} + \beta_t \epsilon_t = \alpha_t (\bar{\alpha}_{t-1} \mathbf{x}_0 + \bar{\beta}_{t-1} \bar{\epsilon}_{t-1}) + \beta_t \epsilon_t = \alpha_t \mathbf{x}_0 + \alpha_t \bar{\beta}_{t-1} \bar{\epsilon}_{t-1} + \beta_t \epsilon_t, \quad (1.9)$$

from which we arrive at a new form of the loss function:

$$\|\epsilon_t - \epsilon_\theta(\alpha_t \mathbf{x}_0 + \alpha_t \bar{\beta}_{t-1} \bar{\epsilon}_{t-1} + \beta_t \epsilon_t, t)\|_2^2. \quad (1.10)$$

Some readers may wonder why, in Eq. (1.9), we have to go back on step to $t-1$, rather than just go from t to 0 directly? The answer is that, since we already sampled ϵ_t in Eq. (1.8), and $\bar{\epsilon}_t$ is not independent from ϵ_t , we cannot sample $\bar{\epsilon}_t$ completely as an independent variable, in which case we end up minimising towards a target that is sample dependent.

This process of stepping back one more step will occur recurrently in many derivations in the subsequent chapters.

1.5. Reducing the square error

In principle, we could train DDPM directly using the loss function Eq. (1.10). But in real practice, this is not very tractable, which can lead to slow convergence. This is because Eq. (1.10) involves four random variables that can only be obtained through samplings:

1. A random sample \mathbf{x}_0 from the inputs.
2. Sample $\bar{\epsilon}_{t-1}$ and $\bar{\epsilon}_t$ from a normal distribution $\mathcal{N}(0, \mathbf{I})$, giving extra two variables.
3. Choosing a t from 1 to T .

The more we need to sample, the more challenge it is to compute the loss function accurately! Fortunately, this problem can be partially mitigated by combining ϵ_t and $\bar{\epsilon}_{t-1}$ into a single random variable sampled from a normal distribution using the integration trick. Because the product of two normal distributions is still a normal distribution, so we have $\alpha_t \bar{\beta}_{t-1} \bar{\epsilon}_{t-1} + \beta_t \epsilon_t$ is equivalent to $\bar{\beta}_t \epsilon | \epsilon \sim \mathcal{N}(0, \mathbf{I})$. Equally, we have $\beta_t \bar{\epsilon}_{t-1} - \alpha_t \bar{\beta}_{t-1} \epsilon_t$ that is equivalent to $\bar{\beta}_t \mathbf{w} | \mathbf{w} \sim \mathcal{N}(0, \mathbf{I})$. Furthermore, it can be proven that $\mathbb{E}[\epsilon \mathbf{w}^T] = 0$, so they are two independent normal distributions. With this, we can express ϵ_t with ϵ and \mathbf{w} :

Not quite sure where this comes from?

$$\epsilon_t = \frac{(\beta_t \epsilon - \alpha_t \bar{\beta}_{t-1} \mathbf{w}) \bar{\beta}_t}{\beta_t^2 + \alpha_t^2 \bar{\beta}_{t-1}^2} = \frac{\beta_t \epsilon - \alpha_t \bar{\beta}_{t-1} \mathbf{w}}{\bar{\beta}_t}. \quad (1.11)$$

Substitute Eq. (1.11) into Eq. (1.10), we get

$$\begin{aligned} & \mathbb{E}_{\bar{\epsilon}_{t-1}, \epsilon_t \sim \mathcal{N}(0, \mathbf{I})} \left[\left\| \epsilon_t - \epsilon_\theta(\alpha_t \mathbf{x}_0 + \alpha_t \bar{\beta}_{t-1} \bar{\epsilon}_{t-1} + \beta_t \epsilon_t, t) \right\|_2^2 \right] \\ &= \mathbb{E}_{\mathbf{w}, \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\left\| \frac{\beta_t \epsilon - \alpha_t \bar{\beta}_{t-1} \mathbf{w}}{\bar{\beta}_t} - \epsilon_\theta(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \epsilon, t) \right\|_2^2 \right]. \end{aligned} \quad (1.12)$$

Note that the loss function is quadratic in \mathbf{w} , which we can expand the square and then determine the expectation value, this gives,

$$\frac{\beta_t^2}{\bar{\beta}_t^2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\left\| \epsilon - \frac{\bar{\beta}_t}{\beta_t} \epsilon_\theta(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \epsilon, t) \right\|_2^2 \right] + \mathcal{C}. \quad (1.13)$$

Ignoring the weight and constant again, we arrive at:

$$\mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\left\| \epsilon - \frac{\bar{\beta}_t}{\beta_t} \epsilon_\theta(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \epsilon, t) \right\|_2^2 \right], \quad (1.14)$$

which is the final form of the lost function for DDPM.

1.6. Regressive generations

Above, we have clarified the training procedures for DDPM. The derivation is a bit lengthy and not immediately trivial. While it is not easy to understand, it is not totally difficult either, since it has not applied the mathematical tools that are needed in the traditional diffusion models, but only an analogous model of building demolition and reconstruction plus some basic knowledge in the probability theory. As such, this newly formulated diffusion generated model, as represented by DDPM, is not as mathematically complicated as one may think.

After the model is being trained, we can perform samplings to generate new samples. Starting from a random noise $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$, and repeat the process for T steps according to Eq. (1.7):

$$\mathbf{x}_{t-1} = \frac{1}{\alpha_t} (\mathbf{x}_t - \beta_t \epsilon_\theta(\mathbf{x}_t, t)). \quad (1.15)$$

This corresponds to the **greedy search** in VAE. For random sampling, we can add an additional noise term:

$$\mathbf{x}_{t-1} = \frac{1}{\alpha_t} (\mathbf{x}_t - \beta_t \epsilon_\theta(\mathbf{x}_t, t)) + \sigma_t \mathbf{z} \quad \mathbf{z} \sim \mathcal{N}(0, \mathbf{I}). \quad (1.16)$$

In practice, we can keep $\beta_t = \sigma_t$, i.e., the variances are the same for both the forward and reverse diffusions.

It should be noted that the sampling process in DDPM is different from the one performed with Langevin equation. In DDPM, each sampling process starts from a new random noise, and is iterated for T steps to generate a new sample. On the other hand, Langevin sampling starts from a single random point and iterates definitively. In principle, all samples can be generated through infinite iterations. Therefore, although DDPM shares similarity with the Langevin sampling, they are fundamentally different models.

1.7. Hyperparameter settings

In DDPM, the upper limit T is usually set to 1000, which is much larger than what people usually expect. Why such a large value is necessary? In the original publication, the following function is used to set α_t :

$$\alpha_t = \sqrt{1 - \frac{0.02t}{T}}, \quad (1.17)$$

which is a monotonically decreasing function. Why such a choice is also necessary?

The above two questions are inter-dependent, and is related to the nature of the data. When the Euclidean distance is used for constructing the loss function, experience of image reconstructions using VAE shows that, only when the input and output images are sufficiently close to each other, we can generate new images of higher quality. Thus, using large time can

This basically says, we must first sample a noise from the normal distribution $\mathcal{N}(0, \mathbf{I})$, from which we can generate the noisy sample $\mathbf{x}_t = \bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \epsilon$. Then we plug this \mathbf{x}_t into the denoising network $\epsilon_\theta(\mathbf{x}_t)$ to recover ϵ_t . This expression is equivalent to the results⁵ derived using KL divergence between the destruction and regeneration process: $\text{argmin}_\theta \mathcal{D}_{KL}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p(\mathbf{x}_{t-1} | \mathbf{x}_t))$ from evidence lower bound, which upon further reparameterisation, gives the following expression for the KL divergence: $\frac{(1-\alpha_t)^2}{2\sigma_q^2(t)} \frac{\|\epsilon_0 - \epsilon_\theta(\mathbf{x}, t)\|_2^2}{(1-\bar{\alpha}_t)\alpha_t}$.

make sure the input and output in each time step are sufficiently similar to each other, reducing the blurring effects from using the Euclidean distance in the loss function.

A similar reason applies to the choice of a monotonically decreasing function for α_t . When t is small, x_t is close to the real image, so we need to decrease the distance between x_{t-1} and x_t , making it more suitable to apply the Euclidean distance in Eq. (1.6). This gives large α_t . When t is very large, x_t becomes very close to a true Gaussian noise, which can be modelled with the Euclidean distance, thus we can increase the distance between x_{t-1} and x_t with smaller α_t . Can we use a large α_t throughout the forward process then? In principle the answer is yes, but we need to increase T . Recall that when we derived Eq. (1.5), we wanted $\bar{\alpha}_T \approx 0$, which we can evaluate

$$\log \bar{\alpha}_T = \frac{1}{2} \sum_{t=1}^T \log \left(1 - \frac{0.02t}{T} \right) < \frac{1}{2} \sum_{t=1}^T \log \left(-\frac{0.02t}{T} \right) = -0.005(T+1) \quad (1.18)$$

With $T = 1000$, $\bar{\alpha}_T \approx 10^{-5}$, which is nearly zero as we wished. As a result, if we want to use a large α_t throughout, then we will have to increase T .

Finally, note that in the loss function [Eq. (1.14), $\varepsilon_\theta(\bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon, t)$], the time variable t has been explicitly included. This is because, at different t , we are dealing with data of different noise levels. Therefore, we are required to have T different reconstruction models, meaning T is also a model hyperparameter.

2. DDPM=Autoregressive VAE

In Section 1, we had established an analogy for the diffusion generative model DDPM, in which we described it as the process of demolition followed by reconstruction. Based on this analogy, we have completed a particular form of the mathematical derivations for the theoretical foundation of DDPM. We also pointed out that, fundamentally, DDPM is not the traditional diffusion model, but more like a VAE. In fact, in the seminal paper for DDPM³, it was derived based on the VAE framework. Therefore, in this section, we will reintroduce the DDPM diffusion model based on the framework of VAE.

Original blog post see <https://www.spaces.ac.cn/archives/9152>

2.1. Solving the problem step-wise

In the standard formulation of the VAE, both the encoding and generative processes are completed in a single step:

$$\text{ENCODER: } x \rightarrow z; \quad \text{GENERATION: } z \rightarrow x. \quad (2.1)$$

In the language of VAE, z is usually referred to as the *latent variable*.

In this case, we only need to deal with three data distributions: the **encoding distribution** $p(z|x)$, the **generative distribution** $q(x|z)$, and the prior distribution $q(z)$. The advantage behind such a formulation is its simplicity. There also exists a simple projective relationship between the sample x and the latent variable z , such that we can train and obtain both the encoder and the decoder (generator) simultaneously. It also allows us to manipulate the data in the latent space. However, VAE also has its notable limitations. Since all three distributions are usually modelled as the normal distributions, it has resulted in the limited expressive capability of the VAE model.

To overcome this limitation, DDPM breaks these processes into multiple steps:

$$\begin{aligned} \text{ENCODER: } & x = x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \cdots \rightarrow x_{T-1} \rightarrow x_T = z \\ \text{DECODER: } & z = x_T \rightarrow x_{T-1} \rightarrow x_{T-2} \rightarrow \cdots \rightarrow x_1 \rightarrow x_0 = x. \end{aligned} \quad (2.2)$$

In this way, every transition probability $p(x_t|x_{t-1})$ or $q(x_{t-1}|x_t)$ is only responsible for making a small change to the sample by the model. One may ask, if both p and q are normal distributions, why this alternative model will work better than the single-stepped one [Eq. (2.1)]? This is because, *Gaussian distributions work better in approximating small incremental changes*, which is similar to approximating a curve with linear relationship that will only work within a small region of Δx . As such, theoretically, the formulation proposed in Eq. (2.2) is able to improve the performance of a single-shot VAE.

2.2. The joint KL divergence

To arrive at a step-wise formulation for VAE, we have $p(x_t|x_{t-1})$ for every encoding step and $q(x_{t-1}|x_t)$ for every generation step. This allows us to write down the **joint probability**

In the following derivations, the convention behind VAE will be adopted, in which we will use p for the transition probability of the forward (encoding) process and q for the backward (decoding) process, which is opposite to the

distributions:

$$\begin{aligned} p(x_0, x_1, x_2, \dots, x_T) &= p(x_T|x_{T-1}) \cdots p(x_2|x_1)p(x_1|x_0)\tilde{p}(x_0) \\ q(x_0, x_1, x_2, \dots, x_T) &= q(x_0|x_1) \cdots q(x_{T-2}|x_{T-1})q(x_{T-1}|x_T)q(x_T). \end{aligned} \quad (2.3)$$

Here, x_0 represents the input samples, thus $\tilde{p}(x_0)$ represents the sample distribution, whereas $q(x_T)$ is the prior distribution, as x_T is the final encoded sample. The rest probability distributions of the type $p(x_t|x_{t-1})$ and $q(x_{t-1}|x_t)$ represent each small step in the encoding and generating processes, respectively.

The simplest way to understand VAE, is to minimise the KL divergence between the joint probability distributions, which is equally applicable to DDPM:

$$D_{KL}(p||q) = \int p(x_T|x_{T-1}) \cdots p(x_1|x_0)\tilde{p}(x_0) \log \frac{p(x_T|x_{T-1}) \cdots p(x_1|x_0)\tilde{p}(x_0)}{q(x_0|x_1) \cdots q(x_{T-1}|x_T)q(x_T)} dx_0 \cdots dx_T. \quad (2.4)$$

This is, therefore, the optimisation target for DDPM. What is left is to consolidate the forms of $p(x_t|x_{t-1})$ and $q(x_{t-1}|x_t)$ to simplify Eq. (2.4) further, making it possible to be computed numerically.

2.3. Divide-and-conquer

First of all, DDPM is only interested in the generating part of the process, therefore, in the encoder process, every step is modelled as a simple normal distribution:

$$p(x_t|x_{t-1}) = \mathcal{N}(x_t; \alpha_t x_{t-1}, \beta_t^2 \mathbf{I}).$$

This is an alternative way of writing Eq. (1.2)

This is much simpler than the VAE, the medium and variance of which are both learnt by a neural network. Here, the medium is basically x_{t-1} that is multiplied by a *scalar* α_t . Hence, DDPM has completely abandoned the capability to encode the data, and only established a generative model. The transition probability for the generating process $q(x_{t-1}|x_t)$ has now been formulated as a learnable model $\mathcal{N}(x_{t-1}; \mu(x_t), \sigma_t^2 \mathbf{I})$. Here, α_t , β_t and σ_t are all not learnables but parameterised. Hence, $\mu(x_t)$ is the only trainable object in the entire diffusion model.

Since p is not parameterised, the integration in Eq. (2.4) that involves purely the p distributions can be factored out as a constant. This makes Eq. (2.4) become equivalent to

$$\begin{aligned} & - \int p(x_T|x_{T-1}) \cdots p(x_1|x_0)\tilde{p}(x_0) \underbrace{\log[q(x_0|x_1) \cdots q(x_{T-1}|x_T)q(x_T)]}_{\text{make this into a summation}} dx_0 \cdots dx_T \\ &= - \int p(x_T|x_{T-1}) \cdots p(x_1|x_0)\tilde{p}(x_0) \left[\log q(x_T) + \sum_{t=1}^T \log q(x_{t-1}|x_t) \right] dx_0 \cdots dx_T. \end{aligned} \quad (2.5)$$

Since $q(x_T)$ is usually modelled as a normal distribution, it is also parameter-free. As such, we only need to compute *every term* of the form

$$\begin{aligned} & - \int \underbrace{p(x_T|x_{T-1}) \cdots p(x_1|x_0)\tilde{p}(x_0)}_{T+1 \text{ terms}} \log q(x_{t-1}|x_t) dx_0 \cdots dx_T \\ &= - \int \underbrace{p(x_t|x_{t-1}) \cdots p(x_1|x_0)\tilde{p}(x_0)}_{t+1 \text{ terms}} \log q(x_{t-1}|x_t) dx_0 \cdots dx_t \\ &= - \int p(x_t|x_{t-1})p(x_{t-1}|x_0)\tilde{p}(x_0) \log q(x_{t-1}|x_t) dx_0 dx_{t-1} x_t. \end{aligned} \quad (2.6)$$

This integral is maximally dependent up to x_t , so all the integrations from x_{t+1} to x_T gives unity, allowing us to remove some terms involving p .

Similarly, the integral is also independent from x_1, \dots, x_{t-2} .

2.4. Reconstructing the diffusion model

What follows from here is basically the same as discussed in Section 1.1.4:

1. Remove all the constants that are irrelevant to the optimisation problem, this left us with the term $-\log q(x_{t-1}|x_t)$ which gives the term $\frac{1}{2\sigma_t^2} \|x_{t-1} - \mu(x_t)\|_2^2$ in the target function for optimisation.

2. The transition probability $p(x_{t-1}|x_0)$ can be equivalently written as $x_{t-1} = \bar{\alpha}_{t-1}x_0 + \bar{\beta}_{t-1}\bar{\epsilon}_{t-1}$, similarly $p(x_t|x_{t-1})$ implies $x_t = \alpha_t x_{t-1} + \beta_t \epsilon_t$, in which $\bar{\epsilon}_{t-1}$ and $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$.
3. From $x_{t-1} = \frac{1}{\alpha_t}(x_t - \beta_t \epsilon_t)$, it inspires us to parameterise $\mu(x_t)$ as $\mu(x_t) = \frac{1}{\alpha_t}(x_t - \beta_t \epsilon_\theta(x_t, t))$.

Equipped with these transformations, the target function for optimisation can be rewritten in the following form:

This is equivalent to Eq. (1.10).

$$\frac{\beta_t^2}{\alpha_t^2 \sigma_t^2} \mathbb{E}_{\bar{\epsilon}_{t-1}, \epsilon_t \sim \mathcal{N}(0, \mathbf{I}), x_0 \sim \tilde{p}(x_0)} \left[\left\| \epsilon_t - \epsilon_\theta(\bar{\alpha}_t x_0 + \alpha_t \bar{\beta}_{t-1} \epsilon_{t-1} + \beta_t \epsilon_t, t) \right\|_2^2 \right]. \quad (2.7)$$

This is followed by the trick of changing the variables applied in Section 1.1.5, which gives

This is equivalent to Eq. (1.14).

$$\frac{\beta_t^4}{\alpha_t^2 \sigma_t^4} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), x_0 \sim \tilde{p}(x_0)} \left[\left\| \epsilon_t - \frac{\bar{\beta}_t}{\beta_t} \epsilon_\theta(\bar{\alpha}_t x_0 + \bar{\beta}_t \epsilon, t) \right\|_2^2 \right] \quad (2.8)$$

This leads to the loss function for DDPM. **As tested in the original paper, the model performance is even better if we drop the coefficient $\frac{\beta_t^4}{\alpha_t^2 \sigma_t^4}$ in the numerical implementation.** The above derivation starts from the loss function for VAE and sequentially simplifies the multivariate integral [Eq. (2.5)]. Although the derivation is longer, it is logically tractable. In comparison to the original DDPM paper³, where a conditional probability $q(x_{t-1}|x_t, x_0)$ was introduced which then split the terms that were subsequently cancelled out. This original derivation came somehow ‘out of the blue’, which may not be easily understood by many readers.

2.5. Hyperparameter settings

Now we move on to discuss the choice of the hyperparameters α_t , β_t and σ_t . For $p(x_t|x_{t-1})$, it is conventional to set $\alpha_t^2 + \beta_t^2 = 1$, which already allows us to reduce the number of hyperparameters by half, and helps to simplify the expressions. This has already been discussed in Section 1. From the superposition of normal distributions, we have the following formula one-shot sampling:

$$p(x_t|x_0) = \int p(x_t|x_{t-1}) \cdots p(x_1|x_0) dx_t dx_{t-1} \cdots dx_1 = \mathcal{N}(x_t; \bar{\alpha}_t x_0, \bar{\beta}_t^2 \mathbf{I}), \quad (2.9)$$

where $\bar{\alpha}_t = \alpha_1 \alpha_2 \cdots \alpha_t$, and $\bar{\beta}_t = \sqrt{1 - \bar{\alpha}_t^2}$.

So what inspires us to set $\alpha_t^2 + \beta_t^2 = 1$? The choice of $\mathcal{N}(x_t; \bar{\alpha}_t x_0, \bar{\beta}_t^2 \mathbf{I})$ means we have $x_t = \alpha_t x_{t-1} + \beta_t \epsilon_t$ with $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$, then it would require $\alpha_t^2 + \beta_t^2 = 1$ based on the normalisation condition for the linear combinations of normal distributions.

As mentioned before, we usually represent $q(x_T)$ as a standard normal distribution $\mathcal{N}(x_T; 0, \mathbf{I})$, and the training goal is to minimise the KL divergence of the two joint probability distributions, i.e. $p = q$. This means that the marginal distributions should also equal each other:

$$\begin{aligned} q(x_T) &= \int p(x_T|x_{T-1}) p(x_{T-1}|x_{T-2}) \cdots p(x_1|x_0) \tilde{p}(x_0) dx_0 dx_1 \cdots dx_{T-1} \\ &= \int p(x_T|x_0) \tilde{p}(x_0) dx_0. \end{aligned} \quad (2.10)$$

In general, the sample distribution is arbitrary, such that the only way to make Eq. (2.10) hold is to equate $p(x_T|x_0)$ to become $q(x_T)$ [Note that the integral in Eq. (2.10) is independent of x_T !]. **It means to evolve $\tilde{p}(x_0)$ into a Gaussian distribution that is independent from x_0 .** Recall that $q(x_T) \sim \mathcal{N}(x_T; \bar{\alpha}_t x_0, \bar{\beta}_t^2 \mathbf{I})$, this would require $\bar{\alpha}_t \approx 0$ when $t \rightarrow T$. More discussions on the choice of α_t can be referred back to Section 1.1.4.

For σ_t , in theory, different sample distributions $\tilde{p}(x_0)$ will lead to a different optimised σ_t . However, we do not really want to train the variance as an extra parameter. In this case, we can choose two special examples.

σ_t comes from the definition of $q(x_{t-1}|x_t)$, which we formulated it as a learnable distribution $\mathcal{N}(x_{t-1}; \mu(x_t), \sigma_t^2 \mathbf{I})$, this is needed for the reverse sampling process.

1. Assume there is only one training sample, x_* . In this case, we have $\tilde{p}(x_0) = \delta(x - x_*)$, the Dirac distribution. This will give an optimum value of $\sigma_t = \beta_{t-1}\beta_t/\bar{\beta}_t$.
2. Assume $\tilde{p}(x_0) \sim \mathcal{N}(0, I)$, then the optimum $\sigma_t = \beta_t$.

Both choices lead to comparable performances in practice. The proofs of these results are rather complicated, which will be reserved for the discussions in the next section.

3. DDPM=Bayes Plus Denoising

In the previous two sections, we have provided two different mathematical derivations for DDPM. The first one is based on an analogous comparison to the building of a house, which is easy to understand, but difficult to be extended theoretically. In the second approach, we draw our conclusions from the framework of autoregressive variational autoencoder. Although theoretically more robust, it is rather formal and lack of inspiration for further improvements. In this section, we will give another derivation of the DDPM, which applies the Bayes' rule to simplify the computations. The derivation contains more reasoning, which makes it more inspirational. Such derivation also has a deeper relationship with the DDIM, which is to be discussed in the next section.

Original blog post see <https://www.spaces.ac.cn/archives/9164>

Diffusion Denoising Implicit Model

3.1. Recap on the model setup

Recall that in the DDPM models, the forward process is modelled as a process of progressive transformation:

$$x = x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \cdots \rightarrow x_{T-1} \rightarrow x_T = z \quad (3.1)$$

in which the process gradually adds noise to the sample x and transform it into a pure noise z , whereas the reverse process gradually denoises z to recover the sample data x . This reverse process is the generative model that we aim to build.

The forward process is simple to achieve, in which every step is given by

$$x_t = \alpha_t x_{t-1} + \beta_t \varepsilon_t \quad \varepsilon_t \sim \mathcal{N}(0, I), \quad (3.2)$$

which can be equivalently written as $p(x_t|x_{t-1}) = \mathcal{N}(x_t; \alpha_t x_{t-1}, \beta_t^2 I)$. Under the constraint of $\alpha_t^2 + \beta_t^2 = 1$, we have the following:

$$\begin{aligned} x_t &= \alpha_t x_{t-1} + \beta_t \varepsilon_t \\ &= \alpha_t (\alpha_{t-1} x_{t-2} + \beta_{t-1} \varepsilon_{t-1}) + \beta_t \varepsilon_t \\ &= \cdots \\ &= (\alpha_t \cdots \alpha_1) x_0 + \underbrace{(\alpha_t \cdots \alpha_2) \beta_1 \varepsilon_1 + (\alpha_t \cdots \alpha_3) \beta_2 \varepsilon_2 + \cdots + \alpha_t \beta_{t-1} \varepsilon_{t-1} + \beta_t \varepsilon_t}_{\sim \mathcal{N}(0, (1 - \alpha_t^2 \cdots \alpha_1^2) I)} \end{aligned} \quad (3.3)$$

From which we have $p(x_t|x_0) = \mathcal{N}(x_t; \bar{\alpha}_t x_0, \bar{\beta}_t^2 I)$ with $\bar{\alpha}_t = \alpha_1 \cdots \alpha_t$ and $\bar{\beta}_t = \sqrt{1 - \bar{\alpha}_t^2}$. What DDPM is trying to achieve, is to figure out how to obtain $p(x_{t-1}|x_t)$ for the reverse generation process using the information above. With this, we could start from any arbitrary $z = x_T$, and through sampling $x_{T-1}, x_{T-2}, \dots, x_1$ stepwise, we could eventually generate a new sample $x = x_0$.

3.2. Using the Bayes' theorem

According to the Bayes' theorem, we have, in this case

$$p(x_{t-1}|x_t) = \frac{p(x_t|x_{t-1})p(x_{t-1})}{p(x_t)}. \quad (3.4)$$

However, we do not know what $p(x_{t-1})$ and $p(x)$ look like, so we cannot move further from this expression. However, we know that these distributions are tractable once they are conditioned on the input x_0 . In this case, we have

$$p(x_{t-1}|x_t, x_0) = \frac{p(x_t|x_{t-1})p(x_{t-1}|x_0)}{p(x_t|x_0)}. \quad (3.5)$$

As such, in what follows, we will focus on deriving expressions for $p(x_{t-1}|x_t)$

Since we know the expressions for $p(x_t|x_{t-1})$, $p(x_{t-1}|x_0)$ and $p(x_t|x_0)$ (the latter two from one-shot sampling), Eq. (3.5) can be simplified into:

$$p(x_{t-1}|x_t, x_0) = \mathcal{N}\left(x_{t-1}; \frac{\alpha_t \bar{\beta}_{t-1}^2}{\bar{\beta}_t^2} x_t + \frac{\bar{\alpha}_{t-1} \beta_t^2}{\bar{\beta}_t^2} x_0, \frac{\bar{\beta}_{t-1}^2 \beta_t^2}{\bar{\beta}_t^2} \mathbf{I}\right). \quad (3.6)$$

To prove Eq. (3.6), note that the exponent in Eq. (3.5), ignoring the factor $\frac{1}{2}$, is

$$\frac{\|x_t - \alpha_t x_{t-1}\|_2^2}{\beta_t^2} + \frac{\|x_{t-1} - \bar{\alpha}_{t-1} x_0\|_2^2}{\bar{\beta}_{t-1}^2} - \frac{\|x_t - \bar{\alpha}_t x_0\|_2^2}{\bar{\beta}_t^2}. \quad (3.7)$$

This expression is quadratic in x_{t-1} , so the final distribution must also be a normal distribution. This means that we only need to determine the corresponding mean and variance. It can be shown that the coefficient for the $\|x_{t-1}\|_2^2$ term is

$$\frac{\alpha_t^2}{\beta_t^2} + \frac{1}{\bar{\beta}_{t-1}^2} = \frac{\alpha_t^2 \bar{\beta}_{t-1}^2 + \beta_t^2}{\beta_t^2 \bar{\beta}_{t-1}^2} = \frac{\alpha_t^2 (1 - \bar{\alpha}_{t-1}^2) + (1 - \alpha_t^2)}{\beta_t^2 \bar{\beta}_{t-1}^2} = \frac{1 - \bar{\alpha}_t^2}{\beta_t^2 \bar{\beta}_{t-1}^2} = \frac{\bar{\beta}_t^2}{\beta_t^2 \bar{\beta}_{t-1}^2} \quad (3.8)$$

This means that the final exponent must be of the form $\frac{\bar{\beta}_t^2}{\beta_t^2 \bar{\beta}_{t-1}^2} \|x_{t-1} - \tilde{\mu}(x_t, x_0)\|_2^2$, such that the covariance matrix is $\frac{\beta_t^2 \bar{\beta}_{t-1}^2}{\bar{\beta}_t^2} \mathbf{I}$. On the other hand, taking out the linear term in x_{t-1} , the coefficient is $-2 \left(\frac{\alpha_t}{\beta_t} x_t + \frac{\bar{\alpha}_{t-1}}{\bar{\beta}_{t-1}} x_0 \right)$, which, upon dividing by $-\frac{2\bar{\beta}_t^2}{\beta_t^2 \bar{\beta}_{t-1}^2}$, we get the following:

$$\tilde{\mu}(x_t, x_0) = \frac{\alpha_t \bar{\beta}_{t-1}^2}{\beta_t} x_t + \frac{\bar{\alpha}_{t-1} \beta_t^2}{\bar{\beta}_t} x_0. \quad (3.9)$$

By now, we have acquired all the information we need for the transition probability $p(x_{t-1}|x_t, x_0)$.

3.3. The denoising process

Although we now have an expression for $p(x_{t-1}|x_t, x_0)$, this is not the final answer that we want. This is because ultimately we would like to predict x_{t-1} from x_t that is independent from x_0 . If, on the other hand, we can predict x_0 from x_t , then we would be able to estimate x_0 in $p(x_{t-1}|x_t, x_0)$ and make it only dependent upon x_t . This can be achieved by letting $x_0 = \bar{\mu}(x_t)$ as a learnable, with a loss function of $\|x_0 - \bar{\mu}(x_t)\|_2^2$. After the model has been trained, we can write

$$p(x_{t-1}|x_t, x_0) \approx p(x_{t-1}|x_t, x_0 = \bar{\mu}(x_t)) = \mathcal{N}\left(x_{t-1}; \frac{\alpha_t \bar{\beta}_{t-1}^2}{\bar{\beta}_t^2} x_t + \frac{\bar{\alpha}_{t-1} \beta_t^2}{\bar{\beta}_t^2} \bar{\mu}(x_t), \frac{\bar{\beta}_{t-1}^2 \beta_t^2}{\bar{\beta}_t^2} \mathbf{I}\right). \quad (3.10)$$

In the expression $\|x_0 - \bar{\mu}(x_t)\|_2^2$, x_0 represents the original data, x_t is the noised data. Hence, we are training a denoising model! This is the meaning of the first D in DDPM.

More specifically, $p(x_t|x_0) = \mathcal{N}(x_t; \bar{\alpha}_t x_0, \bar{\beta}_t^2 \mathbf{I})$, which means $x_t = \bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$, or equivalently, $x_0 = \frac{1}{\bar{\alpha}_t} (x_t - \bar{\beta}_t \varepsilon)$. This inspires us to parameterise $\bar{\mu}(x_t)$ in the form of

$$\bar{\mu}(x_t) = \frac{1}{\bar{\alpha}_t} (x_t - \bar{\beta}_t \varepsilon_\theta(x_t, t)). \quad (3.11)$$

In this case, the loss function for training $\bar{\mu}(x_t)$ becomes

$$\|x_0 - \bar{\mu}(x_t)\|_2^2 = \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \|\varepsilon - \bar{\beta}_t \varepsilon_\theta(\bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon, t)\|_2^2. \quad (3.12)$$

Upon ignoring the coefficient $\bar{\beta}_t^2/\bar{\alpha}_t^2$, we thus recovered the loss function in the original DDPM paper. It can be seen that the derivation here from x_t to x_{t-1} is more straight forward, which avoids the manipulations of integrals as done in the previous section. Now we can substitute Eq. (3.11) back into Eq. (3.10) and simplify, which lead to

$$p(x_{t-1}|x_t, x_0) \approx p(x_{t-1}|x_t, x_0 = \bar{\mu}(x_t)) = \mathcal{N}\left(x_{t-1}; \frac{1}{\bar{\alpha}_t} \left(x_t - \frac{\bar{\beta}_t^2}{\bar{\beta}_t} \varepsilon_\theta(x_t, t)\right), \frac{\bar{\beta}_{t-1}^2 \beta_t^2}{\bar{\beta}_t^2} \mathbf{I}\right). \quad (3.13)$$

Once trained, this will be the probability distribution that will be used for the reverse

sampling process.

To derive this result, note the following:

$$\frac{\bar{\alpha}_{t-1}\beta_t^2}{\bar{\beta}_t^2} \frac{1}{\bar{\alpha}_t} (x_t - \bar{\beta}_t \varepsilon_\theta(x_t, t)) = \frac{\beta_t^2}{\bar{\beta}_t^2 \alpha_t} x_t - \frac{\beta_t^2}{\bar{\beta}_t \alpha_t} \bar{\beta}_t \varepsilon_\theta(x_t, t).$$

Using the fact that $\bar{\alpha}_t = \bar{\alpha}_{t-1} \alpha_t$. Then for the coefficient of x_t , we have

$$\frac{\alpha_t \bar{\beta}_{t-1}^2}{\bar{\beta}_t^2} + \frac{\beta_t^2}{\bar{\beta}_t^2 \alpha_t} = \frac{\alpha_t^2 \bar{\beta}_{t-1}^2 + \beta_t^2}{\bar{\beta}_t^2 \alpha_t} = \frac{\alpha_t^2 (1 - \bar{\alpha}_{t-1}^2) + \beta_t^2}{\bar{\beta}_t^2 \alpha_t} = \frac{(\alpha_t^2 + \beta_t^2) - \alpha_t \bar{\alpha}_{t-1}^2}{\bar{\beta}_t^2 \alpha_t} = \frac{1 - \bar{\alpha}_t^2}{\alpha_t \bar{\beta}_t^2} = \frac{\bar{\beta}_t^2}{\alpha_t \bar{\beta}_t^2} = \frac{1}{\alpha_t}.$$

3.4. Correcting the estimates

An interesting question arises here, which is related to the introduction of $\bar{\mu}(x_t)$ that is used for estimating x_0 when approximating $p(x_{t-1}|x_t, x_0)$ with $p(x_{t-1}|x_t)$. If such a model can be used for estimating x_0 in on-shot, then why would we still need to transform x_t to x_0 slowly in a stepwise manner?

The problem is that, in reality, estimating x_0 from $\bar{\mu}(x_t)$ will not give very accurate result, at least in the first two steps of the reverse process. It only plays a forecasting roles to enable us taking a small step $p(x_{t-1}|x_t)$ backward. This is similar to the idea of ‘estimate-and-correct’ in numerical calculations, in which we start from a rough estimate and iterate forward until an accurate answer is reached. Similar idea can be found in Hinton’s paper ‘*Lookahead optimiser: k steps forward, 1 step back*’.⁶

3.5. Remaining questions

In the previous section, we pointed out that the challenge behind the direct application of Eq. (3.4) is the existence of the terms $p(x_{t-1})$ and $p(x_t)$ that are numerically intractable. This is because, according to the following definition using the marginalisation:

$$p(x_t) = \int p(x_t|x_0) \tilde{p}(x_0) dx_0, \quad (3.14)$$

in which we know the distribution $p(x_t|x_0)$ but $\tilde{p}(x_0)$ is not known a priori, so we cannot compute $p(x_t)$ directly, except in two special cases, which will be discussed here. It also provides the proofs for the choices of the standard deviation σ_t that was introduced in the previous section.

In the first case, we have only one sample in the entire dataset. Without lost of generality, we assume this sample is $\mathbf{0}$, then $\tilde{p}(x_0) = \delta(x_0)$. In this case, we have $p(x_t|\mathbf{0}) = p(x_t)$. Substitute this into Eq. (3.4) shows that this corresponds to the conditional distributions $p(x_{t-1}|x_t, x_0)$ under the special condition of $x_0 = 0$. As such, we now have

This can be obtained by setting $\mu(x_t) \equiv 0$ in Eq. (3.10).

$$p(x_{t-1}|x_t) = p(x_{t-1}|x_t, x_0 = 0) = \mathcal{N}\left(x_{t-1}; \frac{\alpha_t \bar{\beta}_{t-1}^2}{\bar{\beta}_t^2} x_t, \frac{\bar{\beta}_{t-1}^2 \beta_t^2}{\bar{\beta}_t^2} \mathbf{I}\right). \quad (3.15)$$

This shows that, in this case, we have $\sigma_t = \frac{\bar{\beta}_{t-1}^2 \beta_t^2}{\bar{\beta}_t^2}$, which is one particular choice of the standard deviation for the reverse sampling process.

In the second case, the sample follows a normal distribution $\tilde{p}(x_0) = \mathcal{N}(x_0; 0, \mathbf{I})$. From the one-shot sampling in the forward diffusion $p(x_t|x_0) = \mathcal{N}(x_t; \bar{\alpha}_t x_0, \bar{\beta}_t^2 \mathbf{I})$, we have $x_t = \bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon$, with $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$. Since $\tilde{p}(x_0) \sim \mathcal{N}(0, \mathbf{I})$, $p(x_t)$ should also follow a normal distribution. Substituting the expression of a normal distribution into Eq. (3.4) and removing the factor of $-\frac{1}{2}$, the result is

$$\frac{\|x_t - \alpha_{t-1} x_{t-1}\|^2}{\beta_t^2} + \|x_{t-1}\|^2 - \|x_t\|^2. \quad (3.16)$$

Similar to how we derived the expression for $p(x_{t-1}|x_t, x_0)$, the corresponding distribution can be written as

$$p(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \alpha_t x_t, \beta_t^2 \mathbf{I}), \quad (3.17)$$

which shows the variant $\sigma_t = \beta_t^2$, giving another choice if σ_t for the sampling process.

4. DDIM=High Level View of DDPM

Till now, we have shown three different interpretations of DDPM, each from a different perspective. So, is there a higher-level view of DDPM, from which we could acquire some new understanding on the diffusion model? The answer is yes. This leads to the DDIM that was introduced in the paper entitled “Denoising Diffusion Implicit Models”⁷, which will be discussed in this section.

Original blog post see <https://www.spaces.ac.cn/archives/9181>

4.1. The thought process

In the previous sections, we have mentioned that the derivation of DDPM is closely related to DDIM. More specifically, the derivations presented in the previous sections can be briefly summarised as following:

$$p(\mathbf{x}_t|\mathbf{x}_{t-1}) \xrightarrow{\text{derive}} p(\mathbf{x}_t|\mathbf{x}_0) \xrightarrow{\text{derive}} p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \xrightarrow{\text{approximate}} p(\mathbf{x}_{t-1}|\mathbf{x}_t). \quad (4.1)$$

The whole process is progressed in a stepwise manner. However, it is not difficult to find specific characteristics in the above formulation:

1. The loss function for model training is only dependent on $p(\mathbf{x}_t|\mathbf{x}_0)$;
2. The sampling process only relies on $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$.

It means that, although the entire derivation for the DDPM results is started from $p(\mathbf{x}_t|\mathbf{x}_{t-1})$ and progress forwards in a stepwise manner, the final results have actually nothing to do with $p(\mathbf{x}_t|\mathbf{x}_{t-1})$. This makes us to ask the following question:

High-Level Viewpoint 1

Since the final results for the diffusion generative model is independent from $p(\mathbf{x}_t|\mathbf{x}_{t-1})$, can this term be completely removed from the entire derivation?

This is the key idea that led to the birth of DDIM.

4.2. Method of undetermined coefficients

Some readers may find such an idea contradicts the Baye’s rule that was applied in the derivation presented in the previous chapter:

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{t-1}|\mathbf{x}_0)}{p(\mathbf{x}_t|\mathbf{x}_0)}. \quad (4.2)$$

In this case, if we do not have an ansatz for $p(\mathbf{x}_t|\mathbf{x}_{t-1})$, then how will we be able to obtain $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$? This question is actually raised from a very narrow mindset. In fact, without specifying $p(\mathbf{x}_t|\mathbf{x}_{t-1})$, the solution space for $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is even larger. In certain aspects, it even provides a much simpler derivation, because now $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ will only need to satisfy the following condition for the marginal distribution:

$$\int p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)p(\mathbf{x}_t|\mathbf{x}_0)d\mathbf{x}_t = p(\mathbf{x}_{t-1}|\mathbf{x}_0). \quad (4.3)$$

We can solve this equation using the method of undetermined coefficients. In the previous section, we have shown that the solution for $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is a normal distribution, so here, we can generalise this even further and let

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \kappa_t \mathbf{x}_t + \lambda_t \mathbf{x}_0, \sigma_t^2 \mathbf{I}), \quad (4.4)$$

in which κ_t , λ_t and σ_t are all coefficients to be determined. In order not to retrain the model, we keep $p(\mathbf{x}_{t-1}|\mathbf{x}_0)$ and $p(\mathbf{x}_t|\mathbf{x}_0)$ as before. As such, we have in which $\varepsilon, \varepsilon_1, \varepsilon_2 \sim \mathcal{N}(0, \mathbf{I})$, and based on the superposition of normal distributions, we have $\kappa_t \bar{\beta}_t \varepsilon_1 + \sigma_t \varepsilon_2 \sim \sqrt{\kappa_t^2 \bar{\beta}_t^2 + \sigma_t^2} \varepsilon$. By comparing the two different sampling approaches for computing \mathbf{x}_{t-1} , we found that, in order for Eq. (4.3) to hold, we must equate

$$\bar{\alpha}_{t-1} = \kappa_t \bar{\alpha}_t + \lambda_t, \quad \bar{\beta}_{t-1} = \sqrt{\kappa_t^2 \bar{\beta}_t^2 + \sigma_t^2}. \quad (4.5)$$

Notation	Representative Distribution	Sampling
$p(\mathbf{x}_{t-1} \mathbf{x}_0)$	$\mathcal{N}(\mathbf{x}_{t-1}; \bar{\alpha}_{t-1}\mathbf{x}_0, \bar{\beta}_{t-1}^2\mathbf{I})$	$\mathbf{x}_{t-1} = \bar{\alpha}_{t-1}\mathbf{x}_0 + \bar{\beta}_{t-1}\boldsymbol{\varepsilon}$
$p(\mathbf{x}_t \mathbf{x}_0)$	$\mathcal{N}(\mathbf{x}_t; \bar{\alpha}_t\mathbf{x}_0, \bar{\beta}_t^2\mathbf{I})$	$\mathbf{x}_t = \bar{\alpha}_t\mathbf{x}_0 + \bar{\beta}_t\boldsymbol{\varepsilon}_1$
$p(\mathbf{x}_{t-1} \mathbf{x}_t, \mathbf{x}_0)$	$\mathcal{N}(\mathbf{x}_{t-1}; \kappa_t\mathbf{x}_t + \lambda_t\mathbf{x}_0, \sigma_t^2\mathbf{I})$	$\mathbf{x}_{t-1} = \kappa_t\mathbf{x}_t + \lambda_t\mathbf{x}_0 + \sigma_t\boldsymbol{\varepsilon}_2$
$\int p(\mathbf{x}_{t-1} \mathbf{x}_t, \mathbf{x}_0) \cdot p(\mathbf{x}_t \mathbf{x}_0) d\mathbf{x}_t$		$\mathbf{x}_{t-1} = \kappa_t\mathbf{x}_t + \lambda_t\mathbf{x}_0 + \sigma_t\boldsymbol{\varepsilon}_2$ $= \kappa_t(\bar{\alpha}_t\mathbf{x}_0 + \bar{\beta}_t\boldsymbol{\varepsilon}_1) + \lambda_t\mathbf{x}_0 + \sigma_t\boldsymbol{\varepsilon}_2$ $= (\kappa_t\bar{\alpha}_t + \lambda_t)\mathbf{x}_0 + (\kappa_t\bar{\beta}_t\boldsymbol{\varepsilon}_1 + \sigma_t\boldsymbol{\varepsilon}_2)$

In which we have three unknowns but only two equations. This is why the solution space for $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is even larger if we do not specify $p(\mathbf{x}_t|\mathbf{x}_{t-1})$. By making σ_t as an independent variable, we have the following solutions:

$$\kappa_t = \frac{\sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2}}{\bar{\beta}_{t-1}}, \quad \lambda_t = \bar{\alpha}_{t-1} - \frac{\bar{\alpha}_t\sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2}}{\bar{\beta}_t^2} \quad (4.6)$$

or equivalently,

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_{t-1}; \frac{\sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2}}{\bar{\beta}_{t-1}}\mathbf{x}_t + \left(\bar{\alpha}_{t-1} - \frac{\bar{\alpha}_t\sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2}}{\bar{\beta}_t^2}\right)\mathbf{x}_0, \sigma_t^2\mathbf{I}\right). \quad (4.7)$$

For convenience, we choose $\bar{\alpha}_0 = 1$ and $\beta_0 = 0$. More specifically, we don't necessarily need to constrain $\alpha_t^2 + \beta_t^2 = 1$ here. But in order to align with the previous results, we still set $\alpha_t^2 + \beta_t^2 = 1$.

4.3. Same procedure as before

The above derivation shows that, without specifying $p(\mathbf{x}_t|\mathbf{x}_{t-1})$, but given $p(\mathbf{x}_t|\mathbf{x}_0)$ and $p(\mathbf{x}_{t-1}|\mathbf{x}_0)$, we arrive at a set of solutions for $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ that are dependent on the free parameter σ_t . Using the analogy of demolishing and rebuilding from Section 1, it means now we know how the building will finally be demolished into $[p(\mathbf{x}_t|\mathbf{x}_0), p(\mathbf{x}_{t-1}|\mathbf{x}_0)]$, but we do not know how each demolition step $[p(\mathbf{x}_t|\mathbf{x}_{t-1})]$ is achieved. From here, we want to learn how to rebuild the building through $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$. Of course, if we want to know how the demolition process works, we could use the Baye's theorem in the opposite way,

$$p(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = \frac{p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)p(\mathbf{x}_t|\mathbf{x}_0)}{p(\mathbf{x}_{t-1}|\mathbf{x}_0)}. \quad (4.8)$$

What follows is the same as shown in the previous chapter. Since we would like to use $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$ for the generation process, not $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ which is conditioned on \mathbf{x}_0 , we want to parameterise \mathbf{x}_0 with

$$\bar{\boldsymbol{\mu}}(\mathbf{x}_t) = \frac{1}{\bar{\alpha}_t} (\mathbf{x}_t - \bar{\beta}_t\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)). \quad (4.9)$$

Because we did not change $p(\mathbf{x}_t|\mathbf{x}_0)$, the target function for the training the model is still $\|\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}_\theta(\bar{\alpha}_t\mathbf{x}_0 + \bar{\beta}_t\boldsymbol{\varepsilon}, t)\|_2^2$. This means that the training process has not changed. We can use the model trained from DDPM and replace \mathbf{x}_0 in Eq. (4.7) with $\bar{\boldsymbol{\mu}}(\mathbf{x}_t)$, from which we get:

$$\begin{aligned} p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= (\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0 = \bar{\boldsymbol{\mu}}(\mathbf{x}_t)) \\ &= \mathcal{N}\left(\mathbf{x}_{t-1}; \frac{1}{\alpha_t} \left(\mathbf{x}_t - \left(\bar{\beta}_t - \alpha_t\sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2}\right)\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)\right), \sigma_t^2\mathbf{I}\right). \end{aligned} \quad (4.10)$$

This gives $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$ that is needed for the sampling procedure, where $\alpha_t = \bar{\alpha}_t/\bar{\alpha}_{t-1}$. Here, we did not modify the training procedure, as such, the model also did not change. However, the sampling process now has a new parameter σ_t , which brings some fresh results to DDPM.

4.4. A few examples

Refer back to the discussions on the choices of standard deviations σ_t in the previous two sections.

In principle, there is no restriction on the choice of the parameter σ_t . However, different choices of σ_t will lead to completely different sampling process. Here we will give a few examples.

First, let $\sigma_t = \bar{\beta}_{t-1}\beta_t/\bar{\beta}_t$, where $\beta_t = \sqrt{1 - \alpha_t^2}$. Substituting this back into Eq. (4.10) and we have

$$\begin{aligned} p(\mathbf{x}_{t-1}|\mathbf{x}_t) &\approx p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0 = \boldsymbol{\mu}(\mathbf{x}_t)) \\ &= \mathcal{N}\left(\mathbf{x}_{t-1}; \frac{1}{\alpha_t} \left(\mathbf{x}_t - \frac{\beta_t^2}{\bar{\beta}_t} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)\right), \frac{\bar{\beta}_{t-1}^2 \beta_t^2}{\bar{\beta}_t^2} \mathbf{I}\right), \end{aligned} \quad (4.11)$$

which is identical to the result presented in Eq. (3.13). In particular, in the DDIM paper, the results from setting different $\sigma_t = \eta \bar{\beta}_{t-1} \beta_t / \bar{\beta}_t$ with $\eta \in [0, 1]$ had been compared.

In the second example, we have $\sigma_t = \beta_t$, which is the other choice of σ_t that has been discussed previously. However, this does not lead to a more simplified version of Eq. (4.10). Nevertheless, the DDIM shows good performances in the numerical experiments using this standard choice of σ_t from DDPM.

In the most extreme case, we can let $\sigma_t = 0$. Then the transformation from \mathbf{x}_t to \mathbf{x}_0 becomes a well-defined one:

$$\mathbf{x}_{t-1} = \frac{1}{\alpha_t} \left(\mathbf{x}_t - (\bar{\beta}_t - \alpha_t \bar{\beta}_{t-1}) \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)\right). \quad (4.12)$$

This is the case that was the focus of the original DDIM paper. Strictly speaking, DDIM is referred to this particular case of $\sigma_t = 0$, where by the I (implicit) in DDIM means this is an **implicit probabilistic model**. This is because, now the generated \mathbf{x}_0 from a given $\mathbf{x}_T = \mathbf{z}$ is no longer random, which we shall see, this has some advantages both theoretically and practically.

4.5. Accelerated generations

It is worth reiterating that, in this section, we aim to eliminate $p(\mathbf{x}_t|\mathbf{x}_{t-1})$ in the derivation, as a result, all the expressions derived here become dependent on the hyperparameters $\bar{\alpha}_t$ and $\bar{\beta}_t$, whereas α_t and β_t can be derived from $\alpha_t = \bar{\alpha}_t/\bar{\alpha}_{t-1}$ and $\beta_t = \sqrt{1 - \alpha_t^2}$. From the loss function $\|\boldsymbol{\varepsilon} - \boldsymbol{\varepsilon}_\theta(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \boldsymbol{\varepsilon}, t)\|_2^2$, it can be seen that the training process is fixed once $\bar{\alpha}_t$ is specified.

Based on this, it is noticed in DDIM that

High-Level Viewpoint 2

The training outcomes of DDPM in fact include the training outcomes of the parameters for any of its sub-sequences.

More specifically, let $\tau = [\tau_1, \tau_2, \dots, \tau_{\dim(\tau)}]$ being an arbitrary sub-sequence of $[1, 2, \dots, \tau]$. If now we train a DDPM with parameters $\bar{\alpha}_{\tau_1}, \bar{\alpha}_{\tau_2}, \dots, \bar{\alpha}_{\tau_{\dim(\tau)}}$ with $\dim(\tau)$ steps, then the target functions are a subset of the target functions for the original DDPM model with parameters $\bar{\alpha}_1, \bar{\alpha}_2, \dots, \bar{\alpha}_\tau$ that is trained for τ steps. As such, when the model has been well trained, it actually includes all the parameters for models that consist of any of its arbitrary subsequences.

This means that, we can extract a smaller model of $\dim(\tau)$ steps, and then sample according to Eq. (4.10) as

$$p(\mathbf{x}_{\tau_i-1} | \mathbf{x}_{\tau_i}) \approx \mathcal{N}\left(\mathbf{x}_{\tau_i-1}; \frac{\bar{\alpha}_{\tau_i-1}}{\bar{\alpha}_{\tau_i}} \left(\mathbf{x}_{\tau_i} - \left(\bar{\beta}_{\tau_i} - \frac{\bar{\alpha}_{\tau_i}}{\bar{\alpha}_{\tau_i-1}} \sqrt{\bar{\beta}_{\tau_i-1}^2 - \bar{\sigma}_{\tau_i}^2} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_{\tau_i}, \tau_i)\right), \bar{\sigma}_{\tau_i}^2 \mathbf{I}\right)\right) \quad (4.13)$$

This is the process of accelerated sampling, which reduces the original τ diffusion steps to $\dim(\tau)$ steps. Notice that we cannot directly replace α_t in Eq. (4.10) by α_{τ_i} , as we mentioned, α_t is only a derived, auxiliary notation and equivalently, we should have $\alpha_{\tau_i} = \bar{\alpha}_{\tau_i}/\bar{\alpha}_{\tau_i-1}$. Similarly, $\bar{\sigma}_{\tau_i}$ is not σ_{τ_i} , which must now be determined by changing its notations into $\bar{\alpha}_t$ and $\bar{\beta}_t$, followed by replacing t and $t-1$ by τ_i and τ_i-1 , respectively. For instance,

$$\sigma_t = \frac{\bar{\beta}_{t-1} \beta_t}{\bar{\beta}_t} = \frac{\bar{\beta}_{t-1}}{\bar{\beta}_t} \sqrt{1 - \frac{\bar{\alpha}_t^2}{\bar{\alpha}_{t-1}^2}} \rightarrow \frac{\bar{\beta}_{\tau_i-1}}{\bar{\beta}_{\tau_i}} \sqrt{1 - \frac{\bar{\alpha}_{\tau_i}^2}{\bar{\alpha}_{\tau_i-1}^2}} = \bar{\sigma}_{\tau_i}. \quad (4.14)$$

Recall that in the previous section, the result shown in Eq. (4.11) was obtained by assuming $p(\mathbf{x}_t|\mathbf{x}_{t-1}) \sim \mathcal{N}\left(\mathbf{x}_{t-1}; \frac{\alpha_t \bar{\beta}_{t-1}^2}{\bar{\beta}_t^2} \mathbf{x}_t, \frac{\bar{\beta}_{t-1}^2 \beta_t^2}{\bar{\beta}_t^2} \mathbf{I}\right)$ and then apply the Bayes' rule with the known distributions of $p(\mathbf{x}_t|\mathbf{x}_0)$ and $p(\mathbf{x}_{t-1}|\mathbf{x}_0)$. The σ_t parameter can be traced back to the VAE formalism, in which we used the probability $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ for the decoder to calculate the KL-divergence. By taking the logarithm of this parameterised normal distribution, we get $-\log[q(\mathbf{x}_{t-1}|\mathbf{x}_t)] \sim \frac{1}{\sigma_t^2} \|\mathbf{x}_{t-1} - \bar{\boldsymbol{\mu}}(\mathbf{x}_t)\|_2^2$.

So why don't we just train a model of $\dim(\tau)$ steps, but must train a $\tau > \dim(\tau)$ model and then sample a sub-sequence? Possible explanation is that, on one hand, this enhances the model generalisability. On the other hand, it allows us to try other approaches of accelerated sampling without increasing the training costs.

4.6. Differential equation

We now go back to re-examine the case of $\sigma_t = 0$. In this case, Eq. (4.12) can be written equivalently as

$$\frac{x_t}{\bar{\alpha}_t} - \frac{x_{t-1}}{\bar{\alpha}_{t-1}} = \left(\frac{\bar{\beta}_t}{\bar{\alpha}_t} - \frac{\bar{\beta}_{t-1}}{\bar{\alpha}_{t-1}} \right) \varepsilon_\theta(x_t, t). \quad (4.15)$$

When T is sufficiently large, or when $\bar{\alpha}_t$ and $\bar{\alpha}_{t-1}$ are sufficiently small, the above equation becomes the finite-difference form of an ordinary differential equation. More specifically, by introducing an imaginary time variable s , we have

$$\frac{d}{ds} \left(\frac{x(s)}{\bar{\alpha}(s)} \right) = \varepsilon_\theta(x(s), t(s)) \frac{d}{ds} \left(\frac{\bar{\beta}(s)}{\bar{\alpha}(s)} \right). \quad (4.16)$$

Without loss of generality, let $s \in [0, 1]$, in which $s = 0$ corresponds to $t = 0$ and $s = 1$ corresponds to $t = T$. Note that in the original DDPM paper, the imaginary time was set to $\frac{\bar{\beta}(s)}{\bar{\alpha}(s)}$. This is a rather unsuitable choice, as it lies within $[0, \infty)$. This unbounded limit is rather not well suitable for numerical solutions.

What we need to do now is to solve for $x(0)$ given $x(1) \sim \mathcal{N}(0, I)$. The generation process of DDPM and DDIM corresponding to solving this ordinary differential equation (ODE, Eq. (4.15)) with Euler's method, which is known to be the slowest iterative method. Other approaches such as the Heun method and the R-K method may be the accelerated alternatives. This means that writing the generating process as ODE opens up new avenues to achieve accelerated samplings.

For example, using the default DDPM parameters with $T = 1000$ and $\alpha_t = \sqrt{1 - \frac{0.02t}{T}}$, repeating the process outlined in Section 1, we have

$$\log \bar{\alpha}_t = \sum_{i=k}^t \log \alpha_k = \frac{1}{2} \sum_{i=k}^t \log \left(1 - \frac{0.02k}{T} \right) < \frac{1}{2} \sum_{i=k}^t \log \left(-\frac{0.02k}{T} \right) = -\frac{0.005t(t+1)}{T}. \quad (4.17)$$

The above estimation is rather good given that every α_k is very close to 1. On the other hand, since the starting point for this chapter is $p(x_t|x_0)$, so we should start from $\bar{\alpha}_t$. According to the approximation given above, let us now take

$$\bar{\alpha}_t = \exp \left(-\frac{0.005t^2}{T} \right) = \exp \left(-\frac{5t^2}{T^2} \right). \quad (4.18)$$

If we let $s = \frac{t}{T}$, then $s \in [0, 1]$ (which is how we set it in the first place), and $\bar{\alpha}(s) = e^{-5s^2}$. Substituting this into Eq. (4.16) and simplify, we get:

$$\frac{dx(s)}{ds} = 10s \left(\frac{\varepsilon_\theta(x(s), sT)}{\sqrt{1 - e^{-10s^2}}} - x(s) \right). \quad (4.19)$$

We can also take $s = \frac{t^2}{T^2}$, now we still have $s \in [0, 1]$ but with $\bar{\alpha}(s) = e^{-5s}$. Substituting this into Eq. (4.16) and simplify, we get:

$$\frac{dx(s)}{ds} = 5 \left(\frac{\varepsilon_\theta(x(s), \sqrt{s}T)}{\sqrt{1 - e^{-10s}}} - x(s) \right). \quad (4.20)$$

5. The General Framework Based on Stochastic Differential Equation

When writing Section 1 for this series of blog post on diffusion model, some readers had already recommended the paper entitled "*Score-Based Generative Modelling through Stochastic Differential Equation*"⁸ published by Dr Yang Song. This paper has established a rather general theoretical framework for the diffusion model, which has linked up many results derived from DDPM, Stochastic Differential Equations (SDE) and Ordinary Differential Equations (ODE).

Original blog post see <https://www.spaces.ac.cn/archives/9209>

tion (ODE). Although this is a very good paper, it is not well suited for new comers in this field. This is because this paper has applied directly, and extensively, results from SDE, Fokker-Planck equations and score matching, which are rather challenging to be understood. However, with the foundations built from the previous few sections, we can now try to understand this particular paper. In the following chapters, we will try to recover the results derived in this paper, with minimum prerequisite knowledge.

5.1. Stochastic differentiations

In the DDPM formalism, the diffusion process has been divided into a fixed number of T steps. Using the analogy of house demolition and rebuilding that was established in Section 1, it means that both the demolition and rebuilding processes have been divided into T steps beforehand with T being a model hyperparameter. Such a choice is rather artificial. In reality, the entire process should be a transformative one that is continuous in time, which can be described by an SDE.

With this understanding, the forward process can be described by the following SDE:

$$d\mathbf{x} = f_t(\mathbf{x})dt + g_t d\mathbf{w}. \quad (5.1)$$

Some readers may be rather not familiar with this equation, this is of no problem. The equation can be regarded as the limiting case at $\Delta t \rightarrow 0$ for the following equation of finite differences:

$$\mathbf{x}_{t+\Delta t} - \mathbf{x}_t = f_t(\mathbf{x})\Delta t + g_t\sqrt{\Delta t}\boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I}). \quad (5.2)$$

In a more straightforward language, if the entire demolition process takes one day, then it is a transformative process from $t = 0$ to $t = 1$, in which every single step can be described by the above equation of finite differences. There is no specific restriction on Δt , except that the smaller Δt is, the better will Eq. (5.2) approximate the original SDE [Eq. (5.1)]. If we choose $\Delta t = 0.001$, then it corresponds to $T = 1000$ in the original DDPM formulation, and similarly $\Delta t = 0.01$ translates to $T = 100$. Therefore, under the viewpoint of a SDE in the continuous time domain, *different choices of T represent different levels of discretising the SDE*, but should all automatically lead to identical results. As such, *it is not necessary to pre-fix the T values, but changing it based on the requirement of the accuracy levels for different problems*. Therefore, the fundamental advantage of the SDE formalism of the diffusion model is *to separate theoretical analysis from numerical implementations*. We can conduct theoretical analysis based on the mathematical tools of SDE in the continuous time domain, while implementing it with a suitable discretisation strategy for numerical calculations.

Examining Eq. (5.2) in details, some readers may wonder why the first term in the RHS is $\mathcal{O}(\Delta t)$, whereas the second term is $\mathcal{O}(\sqrt{\Delta t})$. In another word, why the deterministic term is of a higher order than the random term? This is indeed a puzzling point in SDE. Simply speaking, $\boldsymbol{\varepsilon}$ always follows a normal distribution if this random term also follows $\mathcal{O}(\sqrt{\Delta t})$. Since the normal distribution has a mean of 0 and a standard deviation of \mathbf{I} , the neighbouring random effect could cancel out each other, and the long-term effect of the random process will only show up if we magnify it to $\mathcal{O}(\sqrt{\Delta t})$.

5.2. The reverse process

In the language of the probability theory, Eq. (5.2) is equivalent to the conditional probability

$$\begin{aligned} p(\mathbf{x}_{t+\Delta t}|\mathbf{x}_t) &= \mathcal{N}(\mathbf{x}_{t+\Delta t}; \mathbf{x}_t + f_t(\mathbf{x}_t)\Delta t, g_t^2\Delta t\mathbf{I}) \\ &\propto \exp\left(-\frac{\|\mathbf{x}_{t+\Delta t} - \mathbf{x}_t - f_t(\mathbf{x}_t)\Delta t\|_2^2}{2g_t^2\Delta t}\right), \end{aligned} \quad (5.3)$$

which, for simplicity, we have ignored the normalisation constant. Following DDPM, we want to know how to rebuild from the demolition process, i.e., we want to determine $p(\mathbf{x}_t|\mathbf{x}_{t+\Delta t})$. Again, from the Baye's rule,

$$\begin{aligned} p(\mathbf{x}_t|\mathbf{x}_{t+\Delta t}) &= \frac{p(\mathbf{x}_{t+\Delta t}|\mathbf{x}_t)p(\mathbf{x}_t)}{p(\mathbf{x}_{t+\Delta t})} = p(\mathbf{x}_{t+\Delta t}|\mathbf{x}_t)\exp[\log p(\mathbf{x}_t) - \log p(\mathbf{x}_{t+\Delta t})] \\ &\propto \exp\left(-\frac{\|\mathbf{x}_{t+\Delta t} - \mathbf{x}_t - f_t(\mathbf{x}_t)\Delta t\|_2^2}{2g_t^2\Delta t} + \log p(\mathbf{x}_t) - \log p(\mathbf{x}_{t+\Delta t})\right). \end{aligned} \quad (5.4)$$

This subsection discuss how the diffusion model, that was established in the discrete time domain, can be generalised into the continuous time domain.

In this case, we can ignore $\log p(\mathbf{x}_t) - \log p(\mathbf{x}_{t+\Delta t}) \rightarrow 0$, then the first term in Eq. (5.4) dominates. When $\Delta t \rightarrow 0$, the fraction can diverge to ∞ , making $p(\mathbf{x}_t|\mathbf{x}_{t+\Delta t})$ vanishes. This can be avoided if \mathbf{x}_t is sufficiently close to $\mathbf{x}_{t+\Delta t}$, such that $\|\mathbf{x}_{t+\Delta t} - \mathbf{x}_t\|_2^2 \sim \Delta t$ (i.e. they are of similar magnitude, so that $\|\mathbf{x}_{t+\Delta t} - \mathbf{x}_t\|_2^2/2g_t^2\Delta t \sim 1$, to avoid vanishing transition probability).

It can be seen that, when Δt is sufficiently small, $p(\mathbf{x}_t|\mathbf{x}_{t+\Delta t})$ will be non-zero only when $\mathbf{x}_{t+\Delta t}$ is very close to \mathbf{x}_t . Similar argument applies to $p(\mathbf{x}_{t+\Delta t}|\mathbf{x}_t)$.

In this case, we can make a Taylor expansion for $\log p(\mathbf{x}_{t+\Delta t})$ as an approximation when $\mathbf{x}_{t+\Delta t}$ is sufficiently close to \mathbf{x}_t :

$$\log p(\mathbf{x}_{t+\Delta t}) \approx \log p(\mathbf{x}_t) + (\mathbf{x}_{t+\Delta t} - \mathbf{x}_t) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \Delta t \frac{\partial}{\partial t} \log p(\mathbf{x}_t). \quad (5.5)$$

Note that the last $\frac{\partial}{\partial t}$ -term cannot be ignored, this is because $p(\mathbf{x}_t)$ is the probability of the random variable being equal to \mathbf{x}_t at time t , and $p(\mathbf{x}_{t+\Delta t})$ is the probability of the random variable being equal to $\mathbf{x}_{t+\Delta t}$ at time $t + \Delta t$. In this case, we must include the partial derivative with respect to the time variable. With this, we can now substitute the expression of $\log p(\mathbf{x}_{t+\Delta t})$ from Eq. (5.5) into Eq. (5.4) and complete the square, we then have,

$$\begin{aligned} p(\mathbf{x}_t|\mathbf{x}_{t+\Delta t}) &\propto \exp \left(-\frac{\|\mathbf{x}_{t+\Delta t} - \mathbf{x}_t - f_t(\mathbf{x}_t)\Delta t\|_2^2}{2g_t^2\Delta t} - (\mathbf{x}_{t+\Delta t} - \mathbf{x}_t) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) - \Delta t \frac{\partial}{\partial t} \log p(\mathbf{x}_t) \right) \\ &\propto \exp \left(-\frac{\|\mathbf{x}_{t+\Delta t} - \mathbf{x}_t - f_t(\mathbf{x}_t)\Delta t\|_2^2 - 2g_t^2(\mathbf{x}_{t+\Delta t} - \mathbf{x}_t) \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \Delta t - 2g_t^2\Delta t^2 \frac{\partial}{\partial t} \log p(\mathbf{x}_t)}{2g_t^2\Delta t} \right) \\ &\propto \exp \left(-\frac{\|\mathbf{x}_{t+\Delta t} - \mathbf{x}_t - [f_t(\mathbf{x}_t) - g_t^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)]\Delta t\|_2^2}{2g_t^2\Delta t} - \frac{1}{2} \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \Delta t + \frac{\partial}{\partial t} \log p(\mathbf{x}_t) \Delta t \right) \\ &\propto \exp \left(-\frac{\|\mathbf{x}_{t+\Delta t} - \mathbf{x}_t - [f_t(\mathbf{x}_t) - g_t^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)]\Delta t\|_2^2}{2g_t^2\Delta t} + \mathcal{O}(\Delta t) \right). \end{aligned} \quad (5.6)$$

When $\Delta t \rightarrow 0$, $\mathcal{O}(\Delta t)$ plays a negligible role, in this case,

$$\begin{aligned} p(\mathbf{x}_t|\mathbf{x}_{t+\Delta t}) &\approx \exp \left(-\frac{\|\mathbf{x}_{t+\Delta t} - \mathbf{x}_t - [f_t(\mathbf{x}_t) - g_t^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)]\Delta t\|_2^2}{2g_t^2\Delta t} \right) \\ &\approx \exp \left(-\frac{\|\mathbf{x}_t - \mathbf{x}_{t+\Delta t} - [f_{t+\Delta t}(\mathbf{x}_{t+\Delta t}) - g_{t+\Delta t}^2 \nabla_{\mathbf{x}_{t+\Delta t}} \log p(\mathbf{x}_{t+\Delta t})]\Delta t\|_2^2}{2g_{t+\Delta t}^2\Delta t} \right). \end{aligned} \quad (5.7)$$

This shows that $p(\mathbf{x}_t|\mathbf{x}_{t+\Delta t})$ is approximately a normal distribution with the following mean:

$$\boldsymbol{\mu} = \mathbf{x}_{t+\Delta t} - [f_{t+\Delta t}(\mathbf{x}_{t+\Delta t}) - g_{t+\Delta t}^2 \nabla_{\mathbf{x}_{t+\Delta t}} \log p(\mathbf{x}_{t+\Delta t})]\Delta t,$$

and standard deviation of $\boldsymbol{\sigma} = g_{t+\Delta t}^2 \Delta t \mathbf{I}$. Taking the limit of $\Delta t \rightarrow 0$, then the corresponding SDE for the reverse process is

$$d\mathbf{x} = [f_t(\mathbf{x}_t) - g_t^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})]dt + g_t d\mathbf{w}. \quad (5.8)$$

This first occurred in the paper entitled ‘‘Reverse-Time Diffusion Equation Model’’⁹. Here, the subscript t has been included in p to signify that this is the distribution at time t .

5.3. Score matching

Now we have obtained the SDE for the reverse diffusion process [Eq. (5.8)], if we further knows the quantity $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$, then we can rebuild the sample from the following discretised version of the SDE:

$$\mathbf{x}_t - \mathbf{x}_{t+\Delta t} = [f_{t+\Delta t}(\mathbf{x}_{t+\Delta t}) - g_{t+\Delta t}^2 \nabla_{\mathbf{x}_{t+\Delta t}} \log p(\mathbf{x}_{t+\Delta t})]\Delta t - g_{t+\Delta t} \sqrt{\Delta t} \boldsymbol{\varepsilon}, \quad (5.9)$$

in which $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I})$. This then completes the derivation of the diffusion generative model. In this case, how would we compute $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$? Note that the probability distribution $p_t(\mathbf{x})$ at t is identical to $p(\mathbf{x}_t)$ in the previous sections, it represents the marginal distribution at time t . This means that we can write the following:

$$p(\mathbf{x}_t|\mathbf{x}_0) = \lim_{\Delta t \rightarrow 0} \int \cdots \int p(\mathbf{x}_t|\mathbf{x}_{t-\Delta t})p(\mathbf{x}_{t-\Delta t}|\mathbf{x}_{t-2\Delta t}) \cdots p(\mathbf{x}_{\Delta t}|\mathbf{x}_0) d\mathbf{x}_{t-\Delta t} d\mathbf{x}_{t-2\Delta t} d\mathbf{x}_{\Delta t}, \quad (5.10)$$

which may be computed directly. This is because in real applications, we will design a model in which $p(\mathbf{x}_t|\mathbf{x}_0)$ has an analytical solution. This is possible when $f_t(\mathbf{x})$ is a linear function

in \mathbf{x} . With this, we have

$$p(\mathbf{x}_t) = \int p(\mathbf{x}_t|\mathbf{x}_0)\tilde{p}(\mathbf{x}_0)d\mathbf{x}_0 = \mathbb{E}_{\mathbf{x}_0}[p(\mathbf{x}_t|\mathbf{x}_0)], \quad (5.11)$$

and

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) = \frac{\mathbb{E}_{\mathbf{x}_0}[\nabla_{\mathbf{x}_t} p(\mathbf{x}_t|\mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0}[p(\mathbf{x}_t|\mathbf{x}_0)]} = \frac{\mathbb{E}_{\mathbf{x}_0}[p(\mathbf{x}_t|\mathbf{x}_0)\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)]}{\mathbb{E}_{\mathbf{x}_0}[p(\mathbf{x}_t|\mathbf{x}_0)]}. \quad (5.12)$$

It can be seen that the last expression in Eq. (5.12) contains a weighted average of $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)$. As we assume $p(\mathbf{x}_t|\mathbf{x}_0)$ has an analytical solution, in principle, Eq. (5.12) will also become computable. However, practically, it involves an average across all the training data \mathbf{x}_0 , which is both heavy in computational costs and lack of good generalisability. Hence we would like to train a neural network $s_\theta(\mathbf{x}_t, t)$, which can be used for computing $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)$.

Some readers may be familiar with the following expression:

$$\mathbb{E}[\mathbf{x}] = \underset{\boldsymbol{\mu}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}} [\|\boldsymbol{\mu} - \mathbf{x}\|_2^2], \quad (5.13)$$

meaning that, in order to make $\boldsymbol{\mu}$ to equal to be the mean of \mathbf{x} , we only need to minimise the mean of $\|\boldsymbol{\mu} - \mathbf{x}\|_2^2$. Similarly, in order to equate $s_\theta(\mathbf{x}_t, t)$ to the weighted average of $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)$ (i.e. $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$), we only need to minimise the weighted average of $\|s_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)\|_2^2$, i.e.

$$\frac{\mathbb{E}_{\mathbf{x}_0} [p(\mathbf{x}_t|\mathbf{x}_0)\|s_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)\|_2^2]}{\mathbb{E}_{\mathbf{x}_0} [p(\mathbf{x}_t|\mathbf{x}_0)]}. \quad (5.14)$$

Again, the purpose here is to seek a way to drop \mathbf{x}_0 in $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)$ for the sampling process such that it is not conditioned on \mathbf{x}_0 .

The denominator $\mathbb{E}_{\mathbf{x}_0} [p(\mathbf{x}_t|\mathbf{x}_0)]$ will only contribute to a weight to the loss function, which can be dropped in the following discussions without affecting the final solution. Finally, we integrate the loss function over \mathbf{x}_t , meaning we must minimise the loss function over every \mathbf{x}_t , from which we reached the final form of the loss function:

$$\begin{aligned} & \int \mathbb{E}_{\mathbf{x}_0} [p(\mathbf{x}_t|\mathbf{x}_0)\|s_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)\|_2^2] d\mathbf{x}_t \\ &= \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_t \sim p(\mathbf{x}_t|\mathbf{x}_0)\tilde{p}(\mathbf{x}_0)} [\|s_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{x}_0)\|_2^2]. \end{aligned} \quad (5.15)$$

This expression provides a physically meaningful explanation to the loss function, which is not necessarily the form that is numerically implemented in most cases.

Eq. (5.15) represents the loss function for the “(conditional) score-matching” generative model. The analytical solutions to the denoising loss function that was derived in Section 2 can be regarded as a special case for this result. The concept of score matching could be first traced back to the 2005 paper entitled “*Estimation of Non-Normalised Statistical Models by Score Matching*”¹⁰. As for the conditional score matching, the author believes that the earliest paper is “*A Connection between Score Matching and Denoising Autoencoder*”¹¹ that was published in 2011.

Although the results presented in this section is the same as the score matching, in our derivation, we did not use the concept of scores. Instead, the results come naturally from manipulating the loss function. The author thinks such an approach is more inspirational to the readers to establish its connection to the generative models. Hopefully, such an approach can also reduce the difficulties for understanding the score matching model.

5.4. The final results

Hence, we have now established the general procedure of a diffusion generative model:

1. Define the forward (demolition) process based on the SDE [Eq. (5.1)].
2. Find the expression for $p(\mathbf{x}_t|\mathbf{x}_0)$.
3. Train the score function $s_\theta(\mathbf{x}_t, t)$ using the loss function Eq. (5.15) (score-matching).
4. Substitute $s_\theta(\mathbf{x}_t, t)$ into Eq. (5.8), replacing the term $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$ to complete the reverse process of re-building.

Most readers may be frightened by the term such as the stochastic differential equations. However, fundamentally, SDE was only used as a bridge for our subsequent discussions. Once the SDE has been discretised into Eq. (5.2) and Eq. (5.3), we can basically forget about SDE, making it much easier to be understood conceptually.

It is also not difficult to see that, it is easy to define an SDE [Eq. (5.1)], but it is rather challenging to solve for $p(\mathbf{x}_t|\mathbf{x}_0)$ from it. In the remaining part of the original paper⁸, it focused on deriving the results for two practical examples. However, if solving for $p(\mathbf{x}_t|\mathbf{x}_0)$ from a pre-defined SDE is too difficult, it will be more trivial to first define $p(\mathbf{x}_t|\mathbf{x}_0)$, just like the case of DDIM, and then deduce the corresponding SDE.

For example, let us first define

$$p(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \bar{\alpha}_t \mathbf{x}_0, \bar{\beta}_t^2 \mathbf{I}). \quad (5.16)$$

Without loss of generality, let us also assume the starting point is $t = 0$ and the end point being $t = 1$, then the corresponding boundary conditions for $\bar{\alpha}_t$ and $\bar{\beta}_t$ are:

$$\bar{\alpha}_0 = 1, \quad \bar{\alpha}_1 = 0, \quad \bar{\beta}_0 = 0, \quad \bar{\beta}_1 = 1. \quad (5.17)$$

In reality, we only need to supply these boundary conditions approximately. For instance, in the previous section, DDPM is equivalent to the choice of $\bar{\alpha}_t = e^{-5t^2}$, such that when $t = 1$, $e^{-5} \approx 0$.

Once we have $p(\mathbf{x}_t|\mathbf{x}_0)$, we can deduce Eq. (5.1) in the reverse order, which, in principle, is to seek a solution for $p(\mathbf{x}_{t+\Delta t}|\mathbf{x}_t)$, which satisfy the condition of

$$p(\mathbf{x}_{t+\Delta t}|\mathbf{x}_0) = \int p(\mathbf{x}_{t+\Delta t}|\mathbf{x}_t) p(\mathbf{x}_t|\mathbf{x}_0) d\mathbf{x}_t. \quad (5.18)$$

Considering only the linear solution, then we have

$$d\mathbf{x} = f_t \mathbf{x} dt + g_t d\mathbf{w}. \quad (5.19)$$

Similar to Section 4, we have the following:

Notation	Representative Distribution	Sampling
$p(\mathbf{x}_{t+\Delta t} \mathbf{x}_0)$	$\mathcal{N}(\mathbf{x}_{t+\Delta t}; \bar{\alpha}_{t+\Delta t} \mathbf{x}_0, \bar{\beta}_{t+\Delta t}^2 \mathbf{I})$	$\mathbf{x}_{t+\Delta t} = \bar{\alpha}_{t+\Delta t} \mathbf{x}_0 + \bar{\beta}_{t+\Delta t} \boldsymbol{\varepsilon}$
$p(\mathbf{x}_t \mathbf{x}_0)$	$\mathcal{N}(\mathbf{x}_t; \bar{\alpha}_t \mathbf{x}_0, \bar{\beta}_t^2 \mathbf{I})$	$\mathbf{x}_t = \bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \boldsymbol{\varepsilon}_1$
$p(\mathbf{x}_{t+\Delta t} \mathbf{x}_t)$	$\mathcal{N}(\mathbf{x}_{t+\Delta t}; \mathbf{x}_t + f_t(\mathbf{x}_t)\Delta t, g_t^2 \Delta t \mathbf{I})$	$\mathbf{x}_{t+\Delta t} = (1 + f_t \Delta t) \mathbf{x}_t + g_t \sqrt{\Delta t} \boldsymbol{\varepsilon}_2$
$\int p(\mathbf{x}_{t+\Delta t} \mathbf{x}_t) \cdot p(\mathbf{x}_t \mathbf{x}_0) d\mathbf{x}_t$		$\mathbf{x}_{t+\Delta t}$ $= (1 + f_t \Delta t) \mathbf{x}_t + g_t \sqrt{\Delta t} \boldsymbol{\varepsilon}_2$ $= (1 + f_t \Delta t) (\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \boldsymbol{\varepsilon}_1) + g_t \sqrt{\Delta t} \boldsymbol{\varepsilon}_2$ $= (1 + f_t \Delta t) \bar{\alpha}_t \mathbf{x}_0 + ((1 + f_t \Delta t) \bar{\beta}_t \boldsymbol{\varepsilon}_1 + g_t \sqrt{\Delta t} \boldsymbol{\varepsilon}_2)$

From here, we get the

$$\begin{aligned} \bar{\alpha}_{t+\Delta t} &= (1 + f_t \Delta t) \bar{\alpha}_t \\ \bar{\beta}_{t+\Delta t}^2 &= (1 + f_t \Delta t)^2 \bar{\beta}_t^2 + g_t^2 \Delta t. \end{aligned} \quad (5.20)$$

By letting $\Delta t \rightarrow 0$, then we can solve for f_t and g_t as

$$\begin{aligned} f_t &= \frac{d}{dt}(\ln \bar{\alpha}_t) = \frac{1}{\bar{\alpha}_t} \frac{d\bar{\alpha}_t}{dt} \\ g_t &= \bar{\alpha}_t^2 \frac{d}{dt} \left(\frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \right) = 2\bar{\alpha}_t \bar{\beta}_t \frac{d}{dt} \left(\frac{\bar{\beta}_t}{\bar{\alpha}_t} \right). \end{aligned} \quad (5.21)$$

The first equation may be derived as following. From Eq. (5.20) we have

$$\bar{\alpha}_t + d\bar{\alpha}_t = (1 + f_t dt) \bar{\alpha}_t$$

$$d\bar{\alpha}_t = f_t dt \bar{\alpha}_t$$

$$\frac{d\bar{\alpha}_t}{\bar{\alpha}_t} = f_t dt$$

If we take $\bar{\alpha}_t \equiv 1$, then the results corresponds to the variance-exploding SDE (VE-SDE). On the other hand, if $\bar{\alpha}_t^2 + \bar{\beta}_t^2 = 1$, this corresponds to the variance-preserving SDE (VP-SDE).

As for the loss function, now we have

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_0) = -\frac{\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0}{\bar{\beta}_t^2} = -\frac{\boldsymbol{\varepsilon}}{\bar{\beta}_t}, \quad (5.22)$$

in which the second equal sign comes from the fact that $\mathbf{x}_t = \bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \boldsymbol{\varepsilon}$. To align with the previous results, let $s_\theta = -\frac{1}{\bar{\beta}_t} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)$, then Eq. (5.15) becomes

$$\frac{1}{\bar{\beta}_t^2} \mathbb{E}_{\mathbf{x}_0 \sim \bar{p}(\mathbf{x}_0), \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I})} [\|\boldsymbol{\varepsilon}_\theta(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \boldsymbol{\varepsilon}, t) - \boldsymbol{\varepsilon}\|_2^2]. \quad (5.23)$$

This is the same loss function as the DDPM after ignoring the coefficients. By replacing $\nabla_{\mathbf{x}_{t+\Delta t}} \log p(\mathbf{x}_{t+\Delta t})$ with $-\frac{1}{\bar{\beta}_{t+\Delta t}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_{t+\Delta t}, t + \Delta t)$, the result has the same form (to the first-order approximation) as for the reverse sampling process in DDPM.

6. The General Framework Based on Ordinary Differential Equation

In the previous chapter, we have discussed the paper entitled “*Score-Based Generative Modelling through Stochastic Differential Equation*” by Dr Yang Song⁸. However, as the title suggested, we had only discussed the part in the original paper that is related to the SDE, but have missed out the part that is related to the discussion around the probability flow ordinary differential equations (ODE), which will be discussed in details in this section.

In fact, the ODE part only constitutes a small portion of the original paper, but we need to discuss it with a separate section. This is because, after thinking about the derivations for a long time, there is no way that it can bypass the Fokker-Planck equation in the whole process. In this case, we must spend some lengths in this section to discuss the Fokker-Planck equation first, before moving on to the main topic of ODE.

6.1. Think again

Let us first summarise the results from the previous section. First of all, we have defined a forward (demolition) process with an SDE:

$$d\mathbf{x} = f_t(\mathbf{x})dt + g_t d\mathbf{w}. \quad (6.1)$$

Then, we have derived the SDE for the corresponding reverse (building) process:

$$d\mathbf{x} = [f_t(\mathbf{x}) - g_t^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g_t d\mathbf{w}. \quad (6.2)$$

Finally, we use a neural network $s_\theta(\mathbf{x}, t)$ to approximate $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ with the following loss function (score-matching):

$$\mathbb{E}_{\mathbf{x}_0, \mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_0) \bar{p}(\mathbf{x}_0)} [\|s_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_0)\|_2^2]. \quad (6.3)$$

At this point, we have established the general framework for training and inferencing the diffusion model, which is the most generalised framework for DDPM. However, as introduced Section 4, is there a similar higher level view on SDE as the generalisation for DDPM? The answer is yes, which is the probability flow ODE to be discussed here.

6.2. The Dirac function

What have we achieved in DDIM? Simply speaking, it was found that the training of DDPM is targeted at $p(\mathbf{x}_t | \mathbf{x}_0)$, that is independent from the transition probability $p(\mathbf{x}_t | \mathbf{x}_{t-1})$. Therefore, DDIM starts from $p(\mathbf{x}_t | \mathbf{x}_0)$ to seek for more general solutions for $p(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ and $p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0)$. The idea of probability flow ODE is similar, which asks the following question:

In the SDE framework, for a fixed $p(\mathbf{x}_t)$, how many different choices of $p(\mathbf{x}_{t+\Delta t} | \mathbf{x}_t)$ are possible, or, can we find different SDEs for the forward process?

To answer this question, similar to the previous section, let us first rewrite the SDE [Eq. (6.1)] in its finite-difference form:

$$\mathbf{x}_{t+\Delta t} = \mathbf{x}_t + f_t(\mathbf{x})\Delta t + g_t \sqrt{\Delta t} \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I}). \quad (6.4)$$

Original blog post see <https://www.spaces.ac.cn/archives/9228>

Eq. (6.4) describes the relationship amongst the random variables $\mathbf{x}_{t+\Delta t}$, \mathbf{x}_t and ε . Although we can easily find the expectation values for both sides, this is not what we want, as we are after the distribution $p(\mathbf{y}_t)$ and the relationship that it should satisfy. To achieve this, we can use the Dirac function:

$$p(\mathbf{x}) = \int \delta(\mathbf{x} - \mathbf{y}) p(\mathbf{y}) d\mathbf{y} = \mathbb{E}_{\mathbf{y}}[\delta(\mathbf{x} - \mathbf{y})]. \quad (6.5)$$

To strictly define the Dirac function, we need to use the tools from functional analysis. But generally, we can obtain the correct result if we just treat it as a normal function. From Eq. (6.5), we also know that the following holds for any $f(\mathbf{x})$:

$$p(\mathbf{x}) f(\mathbf{x}) = \int \delta(\mathbf{x} - \mathbf{y}) p(\mathbf{y}) f(\mathbf{y}) d\mathbf{y} = \mathbb{E}_{\mathbf{y}}[f(\mathbf{y}) \delta(\mathbf{x} - \mathbf{y})]. \quad (6.6)$$

Taking the derivative of both sides, we have

$$\nabla_{\mathbf{x}}[p(\mathbf{x}) f(\mathbf{x})] = \mathbb{E}_{\mathbf{y}}[\nabla_{\mathbf{x}}(f(\mathbf{y}) \delta(\mathbf{x} - \mathbf{y}))] = \mathbb{E}_{\mathbf{y}}[f(\mathbf{y}) \cdot \nabla_{\mathbf{x}} \delta(\mathbf{x} - \mathbf{y})]. \quad (6.7)$$

This is an important relationship that will be used latter. It can be seen that, fundamentally, it should that the derivative of a Dirac function can be transferred to the function that it multiplies to via integration.

6.3. The Fokker-Planck equation

With the preliminary knowledges discussed above, we can now write out the following based on Eq. (6.4):

$$\begin{aligned} & \delta(\mathbf{x} - \mathbf{x}_{t+\Delta t}) \\ &= \delta(\mathbf{x} - \mathbf{x}_t - f_t(\mathbf{x})\Delta t - g_t\sqrt{\Delta t}\varepsilon) \\ &= \delta(\mathbf{x} - \mathbf{x}_t) - (f_t(\mathbf{x})\Delta t + g_t\sqrt{\Delta t}\varepsilon) \nabla_{\mathbf{x}} \delta(\mathbf{x} - \mathbf{x}_t) + \frac{1}{2} (g_t\sqrt{\Delta t}\varepsilon \nabla_{\mathbf{x}})^2 \delta(\mathbf{x} - \mathbf{x}_t), \end{aligned} \quad (6.8)$$

in which we have Taylor expanded $\Delta(\mathbf{x})$ as a normal function, and kept the terms that are not exceeding $\mathcal{O}(\Delta t)$. Now we take the expectation values of both sides:

$$\begin{aligned} & p_{t+\Delta t}(\mathbf{x}) \\ &= \mathbb{E}_{\mathbf{x}_{t+\Delta t}}[\delta(\mathbf{x} - \mathbf{x}_{t+\Delta t})] \\ &\approx \mathbb{E}_{\mathbf{x}_t, \varepsilon} \left[\delta(\mathbf{x} - \mathbf{x}_t) - (f_t(\mathbf{x})\Delta t + g_t\sqrt{\Delta t}\varepsilon) \nabla_{\mathbf{x}} \delta(\mathbf{x} - \mathbf{x}_t) + \frac{1}{2} (g_t\sqrt{\Delta t}\varepsilon \nabla_{\mathbf{x}})^2 \delta(\mathbf{x} - \mathbf{x}_t) \right] \\ &= \mathbb{E}_{\mathbf{x}_t} \left[\delta(\mathbf{x} - \mathbf{x}_t) - f_t(\mathbf{x}_t)\Delta t \nabla_{\mathbf{x}} \delta(\mathbf{x} - \mathbf{x}_t) + \frac{1}{2} g_t^2 \Delta t \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{x}} \delta(\mathbf{x} - \mathbf{x}_t) \right] \\ &= p_t(\mathbf{x}) - \nabla_{\mathbf{x}}[f_t(\mathbf{x})\Delta t p_t(\mathbf{x})] + \frac{1}{2} g_t^2 \Delta t \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{x}} p_t(\mathbf{x}). \end{aligned} \quad (6.9)$$

Now use the fact that $\varepsilon \sim \mathcal{N}(0, I)$, such that $\mathbb{E}_{\varepsilon}[\varepsilon] = 0$

Now if we divide both side by Δt and take the limit of $\Delta t \rightarrow 0$, we then get the Fokker-Planck equation corresponding to Eq. (6.1):

$$\frac{\partial}{\partial t} p_t(\mathbf{x}) = -\nabla_{\mathbf{x}} f_t(\mathbf{x}) p_t(\mathbf{x}) + \frac{1}{2} g_t^2 \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{x}} p_t(\mathbf{x}), \quad (6.10)$$

which is the partial differential equation that describes the time-dependent evolution of the marginal probability distribution.

6.4. Equivalent transformation

There is no need to worry about the appearance of another partial differential equation, as we are not going to seek for its solution but only borrow it to perform an equivalent transformation. In Eq. (6.10), we can add in another function σ_t^2 to the second term and then subtract it out again. Providing that $\sigma_t^2 \leq g_t^2$, Eq. (6.10) is fully equivalent to

$$\frac{\partial}{\partial t} p_t(\mathbf{x}) = -\nabla_{\mathbf{x}} \left[f_t(\mathbf{x}) p_t(\mathbf{x}) - \frac{1}{2} (g_t^2 - \sigma_t^2) \nabla_{\mathbf{x}} p_t(\mathbf{x}) \right] + \frac{1}{2} \sigma_t^2 \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{x}} p_t(\mathbf{x})$$

$$= -\nabla_x \left[\left(f_t(x) - \frac{1}{2}(g_t^2 - \sigma_t^2) \nabla_x \log p_t(x) \right) p_t(x) \right] + \frac{1}{2} \sigma_t^2 \nabla_x \cdot \nabla_x p_t(x). \quad (6.11)$$

In the second line, we have used the fact that $\nabla_x \log p(x) \cdot p(x) = \left[\frac{1}{p(x)} \nabla_x p(x) \right] p(x) = \nabla_x p(x)$.

In terms of the functional form, Eq. (6.11) is equivalent to changing the term of $f_t(x)$ in Eq. (6.10) to $f_t(x) - \frac{1}{2}(g_t^2 - \sigma_t^2) \nabla_x \log p_t(x)$, and g_t^2 into σ_t^2 . Since Eq. (6.10) is to be matched with Eq. (6.1), this means that Eq. (6.11) is equivalent to:

$$dx = \left(f_t(x) - \frac{1}{2}(g_t^2 - \sigma_t^2) \nabla_x \log p_t(x) \right) dt + \sigma_t dw. \quad (6.12)$$

However, do not forget that Eq. (6.10) is fully equivalent to Eq. (6.11), this means that the corresponding marginal distributions $p_t(x)$ that come from Eq. (6.1) and Eq. (6.12) are also fully equivalent to each other. **This means that we have two forward processes with different variances, which will produce identical marginal distributions!**

This result is also equivalent to an updated version of DDIM, which we shall prove it later, that Eq. (6.12) is fully equivalent to DDIM when $f_t(x)$ is a linear function. Furthermore, according to the SDE results derived in the previous section, the corresponding reverse SDE for Eq. (6.12) is

$$dx = \left(f_t(x) - \frac{1}{2}(g_t^2 + \sigma_t^2) \nabla_x \log p_t(x) \right) dt + \sigma_t dw. \quad (6.13) \quad \text{Not quite sure how to derive this.}$$

6.5. Neural ODE

The main advantage of Eq. (6.12) is that it allows us to change the variance during the sampling process. Here, we consider the extreme case of $\sigma_t \equiv 0$, in which case the SDE falls back to an ODE:

$$dx = \left(f_t(x) - \frac{1}{2} g_t^2 \nabla_x \log p_t(x) \right) dt. \quad (6.14)$$

This is the so-called **Probability Flow ODE**. Practically, since $\nabla_x \log p_t(x)$ needs to be approximated by a neural network $s_\theta(x, t)$, so the above equation also corresponds to a neural ODE.

Why do we need to investigate the special case when $\sigma_t \equiv 0$? This is because now the propagation of our samples are noiseless, meaning we now have a well-defined transformation from x_0 to x_T . In this case, we can obtain the inverse transformation from x_T to x_0 by directly solving the ODE in the reverse order, which is also a definitive process. This is consistent with the **flow models**, which transforms noises into samples through an invertible transformation. As such, by using probability flow ODE allows us to link the diffusion model with the flow model. For example, in the original paper⁸, it has mentioned that, by applying the probability flow ODE, we can accurately estimate the likelihood, as well as acquiring the latent variables. These are all the advantages of a flow model. By harnessing the invertibility of the flow model, we can also perform image editing in the latent space.

On the other hand, when the transformation from x_T to x_0 is governed by an ODE, we can apply different higher-order numerical methods for solving the ODE to accelerate the transformation from x_T to x_0 . Although there are some methods to accelerate the solving of SDE, they are far less well investigated compared to ODE. Overall, compared to SDE, ODE is more straightforward for both theoretical analysis and practical solutions.

6.6. Revisiting DDIM

In the last part of Section 4, we demonstrated that the continuous version of the DDIM corresponds to the ODE:

$$\frac{d}{ds} \left(\frac{x(s)}{\bar{\alpha}(s)} \right) = \varepsilon_\theta(x(s), t(s)) \frac{d}{ds} \left(\frac{\bar{\beta}(s)}{\bar{\alpha}(s)} \right). \quad (6.15)$$

We will show below, that this result is a special case when we take the function $f_t(\mathbf{x})$ in Eq. (6.14) as a linear function in \mathbf{x} . As the end of Section 5, we had the following relationships:

$$\begin{cases} f_t = \frac{1}{\bar{\alpha}_t} \frac{\partial \bar{\alpha}_t}{\partial t} \\ g^2(t) = 2\bar{\alpha}_t \bar{\beta}_t \frac{\partial}{\partial t} \left(\frac{\bar{\beta}_t}{\bar{\alpha}_t} \right) \\ s_\theta(\mathbf{x}, t) = -\frac{\varepsilon_\theta(\mathbf{x}, t)}{\bar{\beta}_t}. \end{cases} \quad (6.16)$$

Substituting these relationships into Eq. (6.14) and replacing $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ with $s_\theta(\mathbf{x}, t)$. With further simplification, we obtain:

$$\frac{1}{\bar{\alpha}_t} \frac{d\mathbf{x}}{dt} - \frac{\mathbf{x}}{\bar{\alpha}_t^2} \frac{d\bar{\alpha}_t}{dt} = \varepsilon_\theta(\mathbf{x}, t) \frac{d}{dt} \left(\frac{\bar{\beta}_t}{\bar{\alpha}_t} \right), \quad (6.17)$$

in which the L.H.S. can be further simplified into $\frac{d}{dt} \left(\frac{\mathbf{x}}{\bar{\alpha}_t} \right)$, such that Eq. (6.17) becomes fully equivalent to Eq. (6.15).

As such, in this section, we have shown how to obtain a more general form of the differential equation for the forward process, from which we derived the probability flow ODE and demonstrated DDIM as one of its special case.

7. Optimum Estimation of the Variance (Part I)

For the diffusion models, how to choose the variance as one of its parameters is a critical issue, as the quality of the generated samples strongly depend on the different choices of the variances that are used for the model.

Original blog post see <https://www.spaces.ac.cn/archives/9245>

Recall in Section 1, it was given that, with the trained network $\varepsilon_\theta(\mathbf{x}_t, t)$, the greedy search in the backward direction is $\mathbf{x}_{t-1} = \frac{1}{\alpha_t}(\mathbf{x}_t - \beta_t \varepsilon_\theta(\mathbf{x}_t, t))$. This can be turned into a random search if we add a noise term with variance σ_t , such that $\mathbf{x}_{t-1} = \frac{1}{\alpha_t}(\mathbf{x}_t - \beta_t \varepsilon_\theta(\mathbf{x}_t, t)) + \sigma_t \mathbf{z}$ with $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and for practicality, we select $\sigma_t = \beta_t$.

In Section 2, we assumed two different sample distributions from which we deduced two different choices of σ_t .

Recall that when $\tilde{\mathbf{p}}(\mathbf{x}) = \delta(\mathbf{x}_0)$, i.e. when there exists only a single sample, then the optimum choice of $\sigma_t = \beta_{t-1} \beta_t^2 / \beta_t^2$. These are determined from $p(\mathbf{x}_{t-1} | \mathbf{x}_0)$ which can be computed from $p(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ when $\tilde{\mathbf{p}}(\mathbf{x}_0)$ is known.

In Section 4, when we discussed DDIM, we make the variance as a hyperparameter, and we even had allowed the variance to be zero, although such choice would deteriorate the quality of the generated samples.

Recall that in DDIM, the variance was first introduced in the generalised transition probability for the reversed process (conditioned on \mathbf{x}_0):

$$p(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \kappa_t \mathbf{x}_t + \lambda_t \mathbf{x}_0, \sigma_t^2 \mathbf{I}),$$

which we showed that both κ_t and λ_t can be written as a function of σ_t . This choice of the reverse transition probability has significantly expanded the solution space.

In Section 5, when we discussed the general framework of SDE, we showed that the variances for the forward and reverse SDE should be the same, but in principle, this only holds when $\Delta t \rightarrow 0$.

Recall that in Section 5, we derived $p(\mathbf{x}_t | \mathbf{x}_{t+\Delta t})$ for the reverse process using Baye's rule, Taylor expanded $\log p(\mathbf{x}_{t+\Delta t})$, and combined the Gaussian functions. When Δt is non-zero, $\sigma_t = g_{t+\Delta t}^2 \Delta t \mathbf{I}$, and when $\Delta t \rightarrow 0$, the stochastic term falls back to $g_t d\mathbf{w}$, which is the same as the forward process.

In the paper entitled “Improved Denoising Diffusion Probabilistic Models”¹², the variance was set to be a learnable parameter, but this increased the model training difficulties.

Therefore, how should we choose the variance after all? The two papers entitled “*Analytic-DPM: An Analytic Estimate of the Optimal Reverse Variance in Diffusion Probabilistic Models*”¹³ and “*Estimating the Optimal Covariance with Imperfect Mean in Diffusion Probabilistic Models*”¹⁴ have provided a perfect solution to this problem, which will be discussed here.

7.1. Uncertainty in the reconstruction process

The two publications mentioned above came from the same team. The first article, which will be dubbed as **Analytic-DPM**, derived an analytical solution for the unconditioned variance based on DDIM. In the second paper, which will be dubbed as the **Extended-Analytic-DPM**, had weakened the hypothesis in the first paper, and proposed an optimisation procedure for the conditional variance.

In Section 4, when we presented the DDIM model, we showed that, for a given transition probability $p(x_t|x_0) = \mathcal{N}(x_t; \bar{\alpha}_t x_0, \bar{\beta}_t^2 \mathbf{I})$, the general solution for the corresponding conditional transition probability for the reverse process $p(x_{t-1}|x_t, x_0)$ is

$$p(x_{t-1}|x_t, x_0) = \mathcal{N}\left(x_{t-1}; \frac{\sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2}}{\bar{\beta}_t} x_t + \gamma_t x_0, \sigma_t^2 \mathbf{I}\right), \quad (7.1)$$

in which $\gamma_t = \bar{\alpha}_{t-1} - \bar{\alpha}_t \sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2} / \bar{\beta}_t$, and σ_t is an adjustable variance. In what follows, we applied $\bar{\mu}(x_t)$ to estimate x_0 , and approximated

$$p(x_{t-1}|x_t) \approx p(x_{t-1}|x_t, x_0 = \bar{\mu}(x_t)). \quad (7.2)$$

However, from the viewpoint of Bayesian statistics, such a treatment is not very appropriate, as the prediction of x_0 from x_t cannot be very accurate but with some degrees of uncertainties. Strictly speaking, we should describe it with a probability distribution rather than a function. In fact, we have the following:

$$p(x_{t-1}|x_t) = \int p(x_{t-1}|x_t, x_0) p(x_0|x_t) dx_0. \quad (7.3)$$

There is no way we can acquire an exact solution for $p(x_0|x_t)$, but we only need a rough estimation. In this case, we can use a normal distribution $\mathcal{N}(x_0; \bar{\mu}(x_t), \bar{\sigma}_t^2 \mathbf{I})$ to regress $p(x_0|x_t)$ (details later). With this approximated distribution, we have,

$$\begin{aligned} x_{t-1} &= \frac{\sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2}}{\bar{\beta}_t} x_t + \gamma_t x_0 + \sigma_t \varepsilon_1 \\ &= \frac{\sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2}}{\bar{\beta}_t} x_t + \gamma_t (\bar{\mu}(x_t) + \bar{\sigma}_t \varepsilon_2) + \sigma_t \varepsilon_1 \\ &= \left(\frac{\sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2}}{\bar{\beta}_t} x_t + \gamma_t \bar{\mu}(x_t) \right) + \underbrace{(\sigma_t \varepsilon_1 + \gamma_t \bar{\sigma}_t \varepsilon_2)}_{\sim \sqrt{\sigma_t^2 + \gamma_t^2 \bar{\sigma}_t^2} \varepsilon}, \end{aligned} \quad (7.4)$$

Marginalisation over all possible transitions from x_t to x_0 .

Directly drawing sample from Eq. (7.1)

Draw x_0 from the approximated normal distribution

in which $\varepsilon_1, \varepsilon_2$ and $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$. By comparing Eq. (7.1), Eq. (7.2) and Eq. (7.4), it can be seen that $p(x_{t-1}|x_t)$ is much closer to a normal distribution with the mean of $\frac{\sqrt{\bar{\beta}_{t-1}^2 - \sigma_t^2}}{\bar{\beta}_t} x_t + \gamma_t \bar{\mu}(x_t)$ and standard deviation of $(\sigma_t^2 + \gamma_t^2 \bar{\sigma}_t^2) \mathbf{I}$. In particular, the mean is identical to that in Eq. (7.1), except now the standard deviation contains an extra term of $\gamma_t^2 \bar{\sigma}_t^2$, which will not be zero even if $\sigma_t = 0$. This extra term is the correction term to the optimum variance in the **Analytic-DPM**.

7.2. Optimisation for the mean

Now we move on to discuss why we can use a normal distribution $\mathcal{N}(x_0; \bar{\mu}(x_t), \bar{\sigma}_t^2 \mathbf{I})$ to regress $p(x_0|x_t)$. This problem is equivalent of determining the mean and the variance for the probability distribution $p(x_0|x_t)$. For the mean $\bar{\mu}(x_t)$, it depends on x_t , so it must be

regressed with a model. To do this, we will need a loss function, which we can start with

$$\mathbb{E}_{\mathbf{x}}[\mathbf{x}] = \underset{\boldsymbol{\mu}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}}[\|\mathbf{x} - \boldsymbol{\mu}\|_2^2], \quad (7.5)$$

from which we have

$$\begin{aligned} \bar{\boldsymbol{\mu}}(\mathbf{x}_t) &= \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)}[\mathbf{x}_0] \\ &= \underset{\boldsymbol{\mu}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)}[\|\mathbf{x}_0 - \boldsymbol{\mu}\|_2^2] \\ &= \underset{\boldsymbol{\mu}(t)}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t), \mathbf{x}_t \sim p(\mathbf{x}_t)}[\|\mathbf{x}_0 - \boldsymbol{\mu}(\mathbf{x}_t)\|_2^2] \\ &= \underset{\boldsymbol{\mu}(t)}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0), \mathbf{x}_t \sim p(\mathbf{x}_t|\mathbf{x}_0)}[\|\mathbf{x}_0 - \boldsymbol{\mu}(\mathbf{x}_t)\|_2^2]. \end{aligned} \quad (7.6)$$

This expression gives the formal definition of for the expectation value of \mathbf{x}_0 , if we can draw sample from the distribution of $p(\mathbf{x}_0|\mathbf{x}_t)$.

Applying Eq. (7.5).

Now we also need to sample \mathbf{x}_t as $\boldsymbol{\mu}(\mathbf{x}_t)$ is a function of it

This is the target function for training $\boldsymbol{\mu}(\mathbf{x}_t)$. If we introduce the parameterisation procedure as before, where we let

$$\boldsymbol{\mu}(\mathbf{x}_t) = \frac{1}{\bar{\alpha}_t} (\mathbf{x}_t - \bar{\beta}_t \varepsilon_{\theta}(\mathbf{x}_t, t)), \quad (7.7)$$

then, we recover the loss function for training the DDPM model, which is $\|\varepsilon - \varepsilon_{\theta}(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \varepsilon, t)\|_2^2$. So the procedure from optimising the mean leads to the identical results as before.

7.3. Estimation of the variance: part I

Similar to the above discussions, the covariance matrix is given by

$$\begin{aligned} \boldsymbol{\Sigma}(\mathbf{x}_t) &= \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)}[(\mathbf{x}_0 - \bar{\boldsymbol{\mu}}(\mathbf{x}_t))(\mathbf{x}_0 - \bar{\boldsymbol{\mu}}(\mathbf{x}_t))^{\top}] \\ &= \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)}\left[\left((\mathbf{x}_0 - \boldsymbol{\mu}) - (\bar{\boldsymbol{\mu}}(\mathbf{x}_t) - \boldsymbol{\mu})\right)\left((\mathbf{x}_0 - \boldsymbol{\mu}) - (\bar{\boldsymbol{\mu}}(\mathbf{x}_t) - \boldsymbol{\mu})\right)^{\top}\right] \\ &= \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)}[(\mathbf{x}_0 - \boldsymbol{\mu}_0)(\mathbf{x}_0 - \boldsymbol{\mu}_0)^{\top}] - (\bar{\boldsymbol{\mu}}(\mathbf{x}_t) - \boldsymbol{\mu}_0)(\bar{\boldsymbol{\mu}}(\mathbf{x}_t) - \boldsymbol{\mu}_0)^{\top}, \end{aligned} \quad (7.8)$$

in which $\boldsymbol{\mu}_0$ is an arbitrary vector quantity. The last equation holds due to the translational invariance of the covariance matrix.

The above derivation gives us the full covariance matrix, but this is not what we want, since we would like to use $\mathcal{N}(\mathbf{x}_0; \bar{\boldsymbol{\mu}}(\mathbf{x}_t), \bar{\sigma}_t^2 \mathbf{I})$ to regress $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$, in which, by design, the covariance matrix is $\bar{\sigma}_t^2 \mathbf{I}$. The latter has two characters:

1. It is independent from \mathbf{x}_t , as such, in order to eliminate its dependency on \mathbf{x}_t , we can take the average over all \mathbf{x}_t , such that $\boldsymbol{\Sigma}_t = \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)}[\boldsymbol{\Sigma}(\mathbf{x}_t)]$;
2. It is a multiple of the identity matrix, meaning we only need to consider the diagonal elements and take the mean of the trace, *i.e.* $\bar{\sigma}_t^2 = \operatorname{Tr}(\boldsymbol{\Sigma}_t)/d$, with $d = \dim(\mathbf{x})$.

In this case, we have

$$\begin{aligned} \bar{\sigma}_t^2 &= \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} \left[\frac{\|\mathbf{x}_0 - \boldsymbol{\mu}_0\|_2^2}{d} \right] - \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} \left[\frac{\|\bar{\boldsymbol{\mu}}(\mathbf{x}_t) - \boldsymbol{\mu}_0\|_2^2}{d} \right] \\ &= \frac{1}{d} \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t|\mathbf{x}_0)} [\|\mathbf{x}_0 - \boldsymbol{\mu}_0\|_2^2] - \frac{1}{d} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} [\|\bar{\boldsymbol{\mu}}(\mathbf{x}_t) - \boldsymbol{\mu}_0\|_2^2] \\ &= \frac{1}{d} \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0)} [\|\mathbf{x}_0 - \boldsymbol{\mu}_0\|_2^2] - \frac{1}{d} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} [\|\bar{\boldsymbol{\mu}}(\mathbf{x}_t) - \boldsymbol{\mu}_0\|_2^2] \end{aligned} \quad (7.9)$$

This is one possible analytical form of $\bar{\sigma}_t^2$ that was derived by the author. When the model of $\bar{\boldsymbol{\mu}}(\mathbf{x}_t)$ has been trained, it can be approximately evaluated by first sampling a series of samples \mathbf{x}_0 and \mathbf{x}_t . Most importantly, if we select $\boldsymbol{\mu}_0 = \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0)}[\mathbf{x}_0]$, then we have the following:

$$\bar{\sigma}_t^2 = \operatorname{Var}[\mathbf{x}_0] - \frac{1}{d} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} [\|\bar{\boldsymbol{\mu}}(\mathbf{x}_t) - \boldsymbol{\mu}_0\|_2^2]. \quad (7.10)$$

Here, $\text{Var}[\mathbf{x}_0]$ is the variance of the training data at the pixel level. If the pixel level falls into the range of $[a, b]$ for every pixel, then its variance will not be larger than $\left(\frac{b-a}{2}\right)^2$, such that we have the following inequality:

$$\bar{\sigma}_t^2 \leq \text{Var}[\mathbf{x}_0] \leq \left(\frac{b-a}{2}\right)^2. \quad (7.11)$$

7.4. Estimation of the variance: part II

The above derivation is a relatively straightforward one that was came up by the author (J. Su). In the original Analytic-DPM paper, a slightly different, but less straightforward solution was given. Substituting Eq. (7.7) into the definition of the covariance matrix, we have

$$\begin{aligned} \Sigma(\mathbf{x}_t) &= \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [(\mathbf{x}_0 - \bar{\boldsymbol{\mu}}(\mathbf{x}_t))(\mathbf{x}_0 - \bar{\boldsymbol{\mu}}(\mathbf{x}_t))^\top] \\ &= \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} \left[\left(\left(\mathbf{x}_0 - \frac{\mathbf{x}_t}{\bar{\alpha}_t} \right) + \frac{\bar{\beta}_t}{\bar{\alpha}_t} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right) \left(\left(\mathbf{x}_0 - \frac{\mathbf{x}_t}{\bar{\alpha}_t} \right) + \frac{\bar{\beta}_t}{\bar{\alpha}_t} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right)^\top \right] \\ &= \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} \left[\left(\mathbf{x}_0 - \frac{\mathbf{x}_t}{\bar{\alpha}_t} \right) \left(\mathbf{x}_0 - \frac{\mathbf{x}_t}{\bar{\alpha}_t} \right)^\top \right] - \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)^\top \\ &= \frac{1}{\bar{\alpha}_t^2} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [(\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0)(\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0)^\top] - \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)^\top. \end{aligned} \quad (7.12)$$

Now we take the average over \mathbf{x}_t on both sides, we get

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [(\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0)(\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0)^\top] \\ &= \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t|\mathbf{x}_0)} [(\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0)(\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0)^\top]. \end{aligned} \quad (7.13)$$

Since $p(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \bar{\alpha}_t \mathbf{x}_0, \bar{\beta}_t^2 \mathbf{I})$, so $\bar{\alpha}_t \mathbf{x}_0$ is the mean for the distribution $p(\mathbf{x}_t|\mathbf{x}_0)$, therefore, $\mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t|\mathbf{x}_0)} [(\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0)(\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0)^\top]$ is the covariance matrix for $p(\mathbf{x}_t|\mathbf{x}_0)$ which is obviously $\bar{\beta}_t^2 \mathbf{I}$. Therefore,

$$\mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t|\mathbf{x}_0)} [(\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0)(\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0)^\top] = \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0)} [\bar{\beta}_t^2 \mathbf{I}] = \bar{\beta}_t^2 \mathbf{I}. \quad (7.14)$$

Furthermore, we have

$$\Sigma_t = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}_t)} [\Sigma(\mathbf{x}_t)] = \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} (\mathbf{I} - \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}_t)} [\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)^\top]). \quad (7.15)$$

Taking the trace and then dividing d on both sides, we obtain the following upper limit:

$$\bar{\sigma}_t^2 = \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \left(1 - \frac{1}{d} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}_t)} [\|\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)\|_2^2] \right) \leq \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2}, \quad (7.16)$$

which is the original result from Analytic-DPM.

7.5. Further thoughts

Till this point, we have completely finished introducing Analytic-DPM. The entire derivation is a bit tricky but nonetheless, not too difficult and relatively clear in its logic. This is much friendly to be understood compared to the original paper, which completed the derivations using 7 pages involving 13 Lemmas in the Appendix.

Surely, we admire the observations made by the original authors, who first obtained analytical solutions for the variance. However, the entire original derivation for Analytic-DPM is rather cumbersome and lack of inspirations. For example, the notation $\nabla_{\mathbf{x}_t} p(\mathbf{x}_t)$ had been used throughout the original paper, which causes three problems: Firstly, the whole derivation was not straightforward, making it difficult to understand how the authors came up with the ideas in the first place. Secondly, it makes the understanding of score-matching a prerequisite. Finally, one must change $\nabla_{\mathbf{x}_t} p(\mathbf{x}_t)$ back to either $\bar{\boldsymbol{\mu}}(\mathbf{x}_t)$ or $\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)$ in practical applications, which is an unnecessary detour.

The starting point for our derivations here, is to estimate the parameters for the normal distributions, which its moment estimation is the same as estimating the maximum likelihood. As such, we only need to estimate the corresponding mean and covariance. There is no

real need to match the result with those in the score-matching scheme, since the baseline model for Analytic-DPM is DDIM, which the latter was not based on the score-matching model anyway. Hence it is not beneficial to derive the Analytic-DPM results using the score-matching formalism both theoretically and practically. Using $\bar{\mu}(x_t)$ or $\varepsilon_\theta(x_t, t)$ directly in the derivation, leads to more obvious mathematical form, and much straightforward for numerical implementation.

8. Optimum Estimation of the Variance (Part II)

In the previous section, we have introduced the framework of Analytic-DPM, which leads to an optimum estimation of the variance in the diffusion model that has already been trained. More importantly, the optimum variance was derived as an analytical solution. Numerical experimentations showed that this newly estimated optimum variance could indeed lead to improved qualities of the samples that are generated from the diffusion models.

In this section, we will continue to discuss the upgraded version of Analytic-DPM. This idea came from the paper entitled “*Estimating the Optimum Covariance with Imperfect Mean in Diffusion Probabilistic Models*”¹⁴ published by the same team, and it was referred to as the **Extended-Analytic-DPM**, a term that will also be used here.

8.1. Revisiting the results

We had, in the previous section, started from DDIM, and derived the optimum variance for the generation process being

$$\sigma_t^2 + \gamma_t^2 \bar{\sigma}_t^2, \quad (8.1)$$

in which $\bar{\sigma}_t$ is the variance for $p(x_0|x_t)$. It has the following estimate (from the second approach):

$$\bar{\sigma}_t^2 = \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \left(1 - \frac{1}{d} \mathbb{E}_{x \sim p(x_t)} [\|\varepsilon_\theta(x_t, t)\|_2^2] \right). \quad (8.2)$$

In retrospect, it was not difficult to come up with such a result. Let us assume that

$$\bar{\mu}(x_t) = \frac{1}{\bar{\alpha}_t} (x_t - \bar{\beta}_t \varepsilon_\theta(x_t, t)) \quad (8.3)$$

has already provided us with an accurate estimate of the mean vector for the distribution $p(x_0|x_t)$, then according to the definition of the covariance matrix:

$$\begin{aligned} \Sigma(x_t) &= \mathbb{E}_{x_0 \sim p(x_0|x_t)} [(x_0 - \bar{\mu}(x_t))(x_0 - \bar{\mu}(x_t))^\top] \\ &= \frac{1}{\bar{\alpha}_t^2} \mathbb{E}_{x_0 \sim p(x_0|x_t)} [(x_t - \bar{\alpha}_t x_0)(x_t - \bar{\alpha}_t x_0)^\top] - \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \varepsilon_\theta(x_t, t) \varepsilon_\theta(x_t, t)^\top. \end{aligned} \quad (8.4)$$

To remove the explicit dependency on x_t , we then take the average over $x_t \sim p(x_t)$, which gives

$$\Sigma_t = \mathbb{E}_{x \sim p(x_t)} [\Sigma(x_t)] = \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} (\mathbf{I} - \mathbb{E}_{x \sim p(x_t)} [\varepsilon_\theta(x_t, t) \varepsilon_\theta(x_t, t)^\top]). \quad (8.5)$$

If we then tool the trace of the diagonal element and average it over the matrix dimension, i.e. $\bar{\sigma}_t^2 = \text{Tr}(\Sigma_t)/d$, then we recover the scalar variance, as given in Eq. (8.2).

8.2. Different approaches to improve the results

Before we formally introduce the Extended-Analytic-DPM, we can briefly reflect, are there any other rooms to improve Analytic-DPM? Actually there are quite a few possible improvements.

For example, in Analytic-DPM, we assumed that the normal distribution that we used to approximate $p(x_0|x_t)$ has a covariance matrix of $\sigma_t^2 \mathbf{I}$, in which all the diagonal matrix elements are identical. The most trivial improvement would be to make the diagonal elements different from each other, i.e. $\text{diag}(\bar{\sigma}_t^2)$. Here, we restrict the vector multiplications are all being carried out with the Hamamard product, for example, $x^2 = x \otimes x$. In this case, we only need to consider the diagonal part of Σ_t , and from Eq. (8.5), the corresponding estimate for the

Original blog post see <https://www.spaces.ac.cn/archives/9246>

variance is

$$\bar{\sigma}_t^2 = \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} (\mathbf{1}_d - \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}_t)} [\varepsilon_\theta^2(\mathbf{x}_t, t)]), \quad (8.6)$$

in which $\mathbf{1}_d$ is a d -dimensional unit vector. A further improvement is to retain the dependency of $\bar{\sigma}_t^2$ on \mathbf{x}_t , *i.e.* we consider $\bar{\sigma}_t^2(\mathbf{x}_t)$, making it similar to $\bar{\mu}(\mathbf{x}_t)$, which now $\bar{\sigma}_t^2(\mathbf{x}_t)$ also becomes a learnable model on the input \mathbf{x}_t .

On the other hand, can we consider using the full matrix of Σ_t ? This is workable only in principle but not very practical. This is because Σ_t is a $(d \times d)$ matrix. For high resolution images, as d is the total number of pixels, it makes the storage and computation of such a large matrix very costly.

However, as $\bar{\mu}(\mathbf{x}_t)$ is learnt from the model, it may not be able to accurately represent $\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}_0|\mathbf{x}_t)}[\mathbf{x}_0]$. This is the meaning of the **imperfect mean** in the title for the Extended-Analytic-DPM paper. Therefore, it would be more meaningful and practical to further improve the model under the ideal of imperfect mean.

8.3. Maximum likelihood

Suppose the model for the mean $\bar{\mu}(\mathbf{x}_t)$ has been trained, then, in the distribution $p(\mathbf{x}_0|\mathbf{x}_t) \sim \mathcal{N}(\mathbf{x}_0; \bar{\mu}(\mathbf{x}_t), \bar{\sigma}_t^2 \mathbf{I})$, the only remaining parameter will be $\bar{\sigma}_t^2$. The corresponding negative log-likelihood is

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [-\log \mathcal{N}(\mathbf{x}_0; \bar{\mu}(\mathbf{x}_t), \bar{\sigma}_t^2 \mathbf{I})] \\ &= \frac{1}{2\bar{\sigma}_t^2} \mathbb{E}_{\mathbf{x}_t, \mathbf{x}_0 \sim p(\mathbf{x}_t|\mathbf{x}_0) \tilde{p}(\mathbf{x}_0)} [\|\mathbf{x}_0 - \bar{\mu}(\mathbf{x}_t)\|_2^2] + \frac{d}{2} \log \bar{\sigma}_t^2 + \frac{d}{2} \log 2\pi, \end{aligned} \quad (8.7)$$

The key idea here is, since we know \mathbf{x}_0 , and thus $\tilde{p}(\mathbf{x}_0)$ very well, we should seek an estimate of $\bar{\sigma}_t$ directly from sampling \mathbf{x}_0 rather than sample from \mathbf{x}_t and map it back to \mathbf{x}_0 .

for which the minimum is reached when

$$\bar{\sigma}_t^2 = \frac{1}{d} \mathbb{E}_{\mathbf{x}_t, \mathbf{x}_0 \sim p(\mathbf{x}_t|\mathbf{x}_0) \tilde{p}(\mathbf{x}_0)} [\|\mathbf{x}_0 - \bar{\mu}(\mathbf{x}_t)\|_2^2]. \quad (8.8)$$

This also seems to be the formal definition for the variance.

Here, $\bar{\mu}(\mathbf{x}_t)$ may not be sufficiently accurate, so the second equality in Eq. (8.4) does not really hold, we can only work from the first equality in Eq. (8.4), so we substitute Eq. (8.3) into Eq. (8.8) and we reached:

$$\bar{\sigma}_t^2 = \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2 d} \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0), \varepsilon \sim \mathcal{N}(0, \mathbf{I})} [\|\varepsilon - \varepsilon_\theta(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \varepsilon, t)\|_2^2]. \quad (8.9)$$

This is the result when the covariance matrix is a scaled identity matrix, *i.e.*, $\bar{\sigma}_t^2 \mathbf{I}$. For a more general case where the covariance matrix is a diagonal one, *i.e.*, $\mathcal{N}(\mathbf{x}_0; \bar{\mu}(\mathbf{x}_t), \text{diag}(\bar{\sigma}_t^2))$, the corresponding result is:

$$\bar{\sigma}_t^2 = \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0), \varepsilon \sim \mathcal{N}(0, \mathbf{I})} [\|\varepsilon - \varepsilon_\theta(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \varepsilon, t)\|_2^2]. \quad (8.10)$$

8.4. Conditional variance

If we want to obtain the covariance $\text{diag}(\bar{\sigma}_t^2(\mathbf{x}_t))$ that is conditioned on \mathbf{x}_t , then we need to eliminate the step that takes the average over \mathbf{x}_t , from which we get

$$\bar{\sigma}_t^2(\mathbf{x}_t) = \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [(\mathbf{x}_0 - \bar{\mu}(\mathbf{x}_t))^2] = \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [(\varepsilon_t - \varepsilon_\theta(\mathbf{x}_t, t))^2], \quad (8.11)$$

in which $\varepsilon_t = (\mathbf{x}_t - \bar{\alpha}_t \mathbf{x}_0) / \bar{\beta}_t$. Same as in the previous section, taking the use of

$$\mathbb{E}_{\mathbf{x}}[\mathbf{x}] = \underset{\mu}{\text{argmin}} \mathbb{E}_{\mathbf{x}} [\|\mathbf{x} - \mu\|_2^2], \quad (8.12)$$

we get the following

$$\bar{\sigma}_t^2(\mathbf{x}_t) = \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [(\varepsilon_t - \varepsilon_\theta(\mathbf{x}_t, t))^2]$$

$$\begin{aligned}
&= \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \operatorname{argmin}_g \mathbb{E}_{x_0 \sim p(x_0|x_t)} \left[\|(\varepsilon_t - \varepsilon_\theta(x_t, t))^2 - g\|_2^2 \right] \\
&= \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \operatorname{argmin}_{g(x_t)} \mathbb{E}_{x_t \sim p(x_t)} \mathbb{E}_{x_0 \sim p(x_0|x_t)} \left[\|(\varepsilon_t - \varepsilon_\theta(x_t, t))^2 - g(x_t)\|_2^2 \right] \\
&= \frac{\bar{\beta}_t^2}{\bar{\alpha}_t^2} \operatorname{argmin}_{g(x_t)} \mathbb{E}_{x_t, x_0 \sim p(x_t|x_0) \tilde{p}(x_0)} \left[\|(\varepsilon_t - \varepsilon_\theta(x_t, t))^2 - g(x_t)\|_2^2 \right]. \quad (8.13)
\end{aligned}$$

This is the “NPR-DPM” scheme for learning the conditional variance in the Extended-Analytic-DPM. In addition, the paper also introduced another variance estimation scheme, dubbed “SN-DPM”, which was based on the assumption of a “perfect mean” rather than an imperfect one. However, the result shows that SN-DPM gives better results than NPR-DPM. What this means is that even though the paper was originally designed to solve the imperfect mean problem, the better results of SN-DPM suggests that the perfect mean is more close to practical realisms, such that the problem of imperfect mean is somehow meaningless.

8.5. Two stages

It was mentioned in the paper “*Improved Denoising Diffusion Probabilistic Models*”¹² that, by treating the covariance as a learnable will increase the difficulty in model training. In this case, why we need to train another model for the variance in the Extended-Analytic-DPM?

In the DDPM framework, there exist two possible choices for the variance, namely $\sigma_t = (\bar{\beta}_{t-1}/\bar{\beta}_t)\beta_t$ and $\sigma_t = \beta_t$, which can already provide reasonably good generative outcomes. In such case, further fine tuning the variance will not affect the generative outcomes significantly, at least for a generative process with a large number of steps. The model training, in this case, should still be focusing on learning the mean $\bar{\mu}(x_t)$. On the other hand, if both the mean and variance are learnt simultaneously, then, as the training progresses, the constantly changing variance can interfere with the learning of the mean, which violates the goal to focus on training the mean $\bar{\mu}(x_t)$.

The Extended-Analytic-DPM has devised a genius solution to overcome the problem through a two-staged approach, in which the original variance was used to train the model for the mean $\bar{\mu}(x_t)$, then the trained model was fixed, and some of its parameters were used to train another model for the variance. The main advantages of this approach are:

1. Reducing the number of parameters and the corresponding training costs.
2. Enabling the re-usage of the pre-trained model for the mean, and
3. Achieving a more stable training process.

9. Conditional Generation

In the previous few sections, we have mainly been focusing on discussing the theoretical frameworks for the diffusion models. In this section, we shall focus on a more practical topic, that is the conditional generation of new samples with the diffusion models.

As a generative model, diffusion model shares a similar developmental track as VAE, GAN and the flow model, for which the conditional generative model was soon proposed after the birth of the unconditional generative model. Generally speaking, the unconditional generation can facilitate the exploration of the upper limit in the qualities of the samples that are generated by the model, whereas the conditional generation is more for practical applications, as now we can control the samples that are being generated based on our specific needs. Since the birth of DDPM, many different conditional diffusion generative models have been developed. The capability to conditionally control the diffusion generative process is the driving force that makes the diffusion models become so popular these days, such as the text-to-image models DALL.E2 and Imagen.

In this section, we will briefly study and summarise the theoretical foundations for the conditional diffusion models.

9.1. Theoretical analysis

Practically, there are two pathways to realise conditional generation, which are usually dubbed as either **Classifier-Guidance** or **Classifier-Free**.

For most people, the cost for training a diffusion model at the SOTA-level is too large. On the other hand, the cost associated with training a classifier is more or less acceptable. As such, most people wants to use a pre-trained diffusion model and add in an additional classifier to

Original blog post see <https://www.spaces.ac.cn/archives/9257>

realise the conditional sample generation, This is the *ad-hoc Classifier-Guidance* approach. Nevertheless, for Tech Giants such as Google and OpenAI, that are not lack of both data and computational power, they are more keen to include the conditional signals during the training of the diffusion model, in order to achieve better sample qualities, this is the so-called **Classifier-Free** approach.

The **Classifier-Guidance** approach was first proposed in the article “*Diffusion Models Beats GANs on Image Synthesis*”¹⁵, which was aimed to achieve sample generations according to different classifications. Subsequently, in the paper “More Control for Free! Image Synthesis with Semantic Diffusion Guidance”¹⁶, the concept of the classifier has been further generalised, making it possible to generate samples that are guided by either images or texts. Generally speaking, the training costs for the Classifier-Guidance models are lower. However, the cost for inferencing is higher and the sample quality is lower compared the Classifier-Free approach.

As for the **Classifier-Free** approach, it was first proposed in the article “*Classifier-Free Diffusion Guidance*”¹⁷. Models such as DALL.E2 and Imagen are all based on this approach. Strictly speaking, there is not much technical tricks involved in the Classifier-Free approach. It is the most straightforward generalisation of diffusion model to its conditional variant. It appears later than the Classifier-Guidance approach most likely because of its high training cost. Nevertheless, when there exists sufficient data and computational power for training, the samples generated from the Classifier-Free approach have demonstrated amazing capability in controlling the details in the process of sample generation.

9.2. Conditional input

Simply put, there is not much technicalities in the Classifier-Free approach, except for its high training costs. As such, we will focus the subsequent discussions on the Classifier-Guidance approach, and briefly outline the Classifier-Free approach at the end.

Based on the analysis from the previous sections, it should be clear to the readers that the key step in the generation process is to build the transition probability $p(x_{t-1}|x_t)$. As such, for sample generation based on some condition y , all we need to do is introduce y into the generation process via the conditional probability $p(x_{t-1}|x_t, y)$ which replaces $p(x_{t-1}|x_t)$. In order to reuse the $p(x_{t-1}|x_t)$ from the trained unconditional generative model, we apply the Baye’s rule

$$p(x_{t-1}|y) = \frac{p(x_{t-1})p(y|x_{t-1})}{p(y)}, \quad (9.1)$$

and putting x_t back into every term, we get:

$$p(x_{t-1}|x_t, y) = \frac{p(x_{t-1}|x_t)p(y|x_{t-1}, x_t)}{p(y|x_t)}. \quad (9.2)$$

Note that x_t is generated from x_{t-1} by adding the noise in the forward process. In this case, adding x_t will not bring any benefit to the classifier, as such, $p(y|x_{t-1}, x_t) = p(y|x_{t-1})$, this leads to the following:

$$p(x_{t-1}|x_t, y) = \frac{p(x_{t-1}|x_t)p(y|x_{t-1})}{p(y|x_t)} = p(x_{t-1}|x_t) \exp[\log p(y|x_{t-1}) - \log p(y|x_t)]. \quad (9.3)$$

9.3. Approximated distribution

Readers who are familiar with the derivations in Section 5 on SDE will not be surprised by the following derivations. When T is sufficiently large, the standard deviations for the distribution $p(x_t|x_{t-1})$ will be sufficiently small. This is saying that the transition probability $p(x_t|x_{t-1})$ will only be significantly larger than zero when x_t is very close to x_{t-1} . This is also true for the conditional probabilities, *i.e.* either $p(x_{t-1}|x_t, y)$ or $p(x_t|x_{t-1}, y)$ will only be significantly larger than zero when x_t is very close to x_{t-1} . As such, we only need to consider the change in the probability distributions within this small range of x . With this, we can apply the Taylor expansion as

$$\log p(y|x_{t-1}) - \log p(y|x_t) \approx (x_{t-1} - x_t) \nabla_{x_t} \log p(y|x_t). \quad (9.4)$$

Strictly speaking, there should also exist an extra t -dependent term, but it is irrelevant to x_{t-1} , as it is a constant that will not change the probability of x_{t-1} . Therefore, this term

is ignored in Eq. (9.4). Assuming that we already have $p(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}(\mathbf{x}_t), \sigma_t^2 \mathbf{I}) \propto \exp[-\|\mathbf{x}_{t-1} - \boldsymbol{\mu}(\mathbf{x}_t)\|^2/2\sigma_t^2]$, then we have the following approximation:

$$\begin{aligned} p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}) &\propto \exp[-\|\mathbf{x}_{t-1} - \boldsymbol{\mu}(\mathbf{x}_t)\|^2/2\sigma_t^2 + (\mathbf{x}_{t-1} - \mathbf{x}_t)\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)] \\ &\propto \exp[-\|\mathbf{x}_{t-1} - \boldsymbol{\mu}(\mathbf{x}_t) - \sigma_t^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)\|^2/2\sigma_t^2 \log p(\mathbf{y}|\mathbf{x}_t)]. \end{aligned} \quad (9.5)$$

This result shows that the distribution $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})$ is approximately a normal distribution of $\mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}(\mathbf{x}_t) + \sigma_t^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t), \sigma_t^2 \mathbf{I})$. In this case, we only need to modify the sampling process to:

$$\mathbf{x}_{t-1} = \boldsymbol{\mu}(\mathbf{x}_t) + \sigma_t^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) + \sigma_t \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I}). \quad (9.6)$$

This is the core result in the Classifier-Guidance Approach. Note that \mathbf{x}_t is the sample after noises are being added, which is the input for $p(\mathbf{y}|\mathbf{x}_t)$. It means that we need a model that can predict the noised sample. If we only have a model to predict the noise-free sample $p_0(\mathbf{y}|\mathbf{x}_0)$, then what we way consider is

$$p(\mathbf{y}|\mathbf{x}_0) = p_0(\mathbf{y}|\boldsymbol{\mu}(\mathbf{x}_t)), \quad (9.7)$$

i.e., we use $\boldsymbol{\mu}(\cdot)$ to denoise \mathbf{x}_t before plugging it into $p(\mathbf{y}|\mathbf{x}_t)$. This avoids the cost of training a classifier with noisy samples.

9.4. Scaling the classifier gradient

In the original paper entitled “*Diffusion Model Beats GAN on Image Synthesis*”, it was found that if a scaling parameter γ was introduced into the gradient of the classifier, one can better attenuate the generated samples. In this case, we have:

$$\mathbf{x}_{t-1} = \boldsymbol{\mu}(\mathbf{x}_t) + \sigma_t^2 \gamma \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) + \sigma_t \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I}). \quad (9.8)$$

When $\gamma > 1$, the generation process will utilise more information from the classifier. This will increase the correlation between the output and the input condition \mathbf{y} , at the price of reducing the diversity of the output samples. On the other hand, when $\gamma < 1$, the generative process will reduce the correlation between the generated output and the conditional input, thus leading to more diversified outcomes.

How should we understand this parameter? In the original paper, it was proposed that γ could be understood by adjusting the form of the probability distribution through the use of a power operation:

$$\tilde{p}(\mathbf{y}|\mathbf{x}_t) = \frac{p^\gamma(\mathbf{y}|\mathbf{x}_t)}{\mathcal{Z}(\mathbf{x}_t)}, \quad \mathcal{Z}(\mathbf{x}_t) = \sum_{\mathbf{y}} p^\gamma(\mathbf{y}|\mathbf{x}_t). \quad (9.9)$$

As γ increases, $\tilde{p}(\mathbf{y}|\mathbf{x}_t)$ will slowly merge into an one-hot distribution. If this distribution was used to substitute $p(\mathbf{y}|\mathbf{x}_0)$ in the classifier as the guidance, then the generative process will tend to select samples that have higher confidence level in being correctly classified. Although this explanation can facilitate our understanding to some extent, it was not completely right. This is because

$$\nabla_{\mathbf{x}_t} \log \tilde{p}(\mathbf{y}|\mathbf{x}_t) = \gamma \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \mathcal{Z}(\mathbf{x}_t) \neq \gamma \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t). \quad (9.10)$$

The authors of the original paper had wrongly took $\mathcal{Z}(\mathbf{x}_t)$ as a constant, which gave $\nabla_{\mathbf{x}_t} \mathcal{Z}(\mathbf{x}_t) = 0$. However, when $\gamma \neq 1$, it is obvious that $\mathcal{Z}(\mathbf{x}_t)$ will also be dependent on \mathbf{x}_t . Unfortunately, there does not seem to be a suitable solution to overcome this problem, we can only roughly assume that the property of the gradient function at $\gamma = 1$ can be generalised to the case when $\gamma \neq 1$.

9.5. Similarity control

As a matter of fact, the best way to understand the case when $\gamma \neq 1$ is to abandon the Bayes’ rule [Eq. (9.2) and Eq. (9.3)] when trying to understand the term $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})$. Instead, we could straightforwardly define

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}) = \frac{p(\mathbf{x}_{t-1}|\mathbf{x}_t) e^{\gamma \cdot \text{sim}(\mathbf{x}_{t-1}, \mathbf{y})}}{\mathcal{Z}(\mathbf{x}_t, \mathbf{y})}, \quad \mathcal{Z}(\mathbf{x}_t|\mathbf{y}) = \sum_{\mathbf{x}_{t-1}} p(\mathbf{x}_{t-1}|\mathbf{x}_t) e^{\gamma \cdot \text{sim}(\mathbf{x}_{t-1}, \mathbf{y})}, \quad (9.11)$$

in which $\text{sim}(\mathbf{x}_{t-1}, \mathbf{y})$ is a similarity measure between the generated sample \mathbf{x}_{t-1} and the conditional input \mathbf{y} . From this perspective, γ comes naturally from the definition of $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})$, which directly controls the correlation between the output and the condition. As γ increases, the model will generate \mathbf{x}_{t-1} that has stronger correlation with \mathbf{y} .

In order to arrive at an approximation that will further facilitate samplings, we could expand the similarity term approximately the point \mathbf{x}_{t-1} that is close to \mathbf{x}_t as:

$$e^{\gamma \cdot \text{sim}(\mathbf{x}_{t-1}, \mathbf{y})} \approx e^{\gamma \cdot \text{sim}(\mathbf{x}_t, \mathbf{y}) + \gamma(\mathbf{x}_{t-1} - \mathbf{x}_t) \nabla_{\mathbf{x}_t} \text{sim}(\mathbf{x}_t, \mathbf{y})}. \quad (9.12)$$

Assuming that this approximation is sufficiently accurate, we can then eliminate the term that are irrelevant to \mathbf{x}_{t-1} , from which we get

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}) \propto p(\mathbf{x}_{t-1}|\mathbf{x}_t) e^{\gamma(\mathbf{x}_{t-1} - \mathbf{x}_t) \nabla_{\mathbf{x}_t} \text{sim}(\mathbf{x}_t, \mathbf{y})}. \quad (9.13)$$

Similarly, substituting in $p(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}(\mathbf{x}_t), \sigma_t^2 \mathbf{I})$ and completing the square, we get

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t) \propto \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}(\mathbf{x}_t) + \gamma \sigma_t^2 \nabla_{\mathbf{x}_t} \text{sim}(\mathbf{x}_t, \mathbf{y}), \sigma_t^2 \mathbf{I}), \quad (9.14)$$

which now we no longer need to worry about the probabilistic interpretation of $p(\mathbf{y}|\mathbf{x}_t)$, but only a suitable definition of the similarity measure $\text{sim}(\mathbf{x}_t, \mathbf{y})$. Here, \mathbf{y} becomes no longer restricted to a class, it could be any type of inputs (image or text). The general treatment is to encode them into a vector representation using their own encoder and then apply the cosine similarity:

$$\text{sim}(\mathbf{x}_t, \mathbf{y}) = \frac{E_1(\mathbf{x}_t) \cdot E_2(\mathbf{y})}{\|E_1(\mathbf{x}_t)\| \cdot \|E_2(\mathbf{y})\|}. \quad (9.15)$$

It should be noted that, since the intermediate state \mathbf{x}_t contains Gaussian noise, it is generally better not to use the encoder for the noise-free sample but fine-tune it with noised samples. Above is the core idea in the article “*More Control for Free! Image Synthesis with Semantic Diffusion Guidance*”.

9.6. Continuous scenarios

We saw that in all the previous derivations, that the correction to the mean is either $\sigma_t^2 \gamma \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)$ or $\gamma \sigma_t^2 \nabla_{\mathbf{x}_t} \text{sim}(\mathbf{x}_t, \mathbf{y})$, which both contain the term σ_t , meaning that when $\sigma_t = 0$, the correcting term will also be zero, making the conditional input playing no effect on sample generation.

So in this case, can we still set σ_t to be zero? The answer is yes. For example, in DDIM introduced in Section 4, it has its corresponding generation process with $\sigma_t = 0$. How should we realise conditional generation? In this case, we need to apply the SDE formalism introduced in Section 6, in which we showed that for a forward SDE:

$$d\mathbf{x} = f_t(\mathbf{x})dt + g_t d\mathbf{w}, \quad (9.16)$$

with the corresponding reverse SDE

$$d\mathbf{x} = \left(f_t(\mathbf{x}) - \frac{1}{2}(g_t^2 + \sigma_t^2) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}) \right) dt + \sigma_t d\mathbf{w}. \quad (9.17)$$

Here, we are allowed to choose σ_t freely. Both DDPM and DDIM can be considered as its special cases. Especially when $\sigma_t = 0$, this corresponds to a generalised DDIM. It can be seen that, in the reverse SDE, the term that is input-dependent is $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x})$. If now we want to perform conditional generation, we should just need to change it to $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}|\mathbf{y})$. From Bayes' rule, we have

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}|\mathbf{y}) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) + \nabla_{\mathbf{x}} \log p_t(\mathbf{y}|\mathbf{x}), \quad (9.18)$$

noting that $\nabla_{\mathbf{x}} \log p_t(\mathbf{y}) = 0$. With the parameterisation trick that we had used before, we have $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = -\varepsilon_{\theta}(\mathbf{x}_t, t)/\bar{\beta}_t$, thus we have

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}|\mathbf{y}) = -\frac{\varepsilon_{\theta}(\mathbf{x}_t, t)}{\bar{\beta}_t} + \nabla_{\mathbf{x}} \log p_t(\mathbf{y}|\mathbf{x}) = -\frac{\varepsilon_{\theta}(\mathbf{x}_t, t) - \bar{\beta}_t \nabla_{\mathbf{x}} \log p_t(\mathbf{y}|\mathbf{x})}{\bar{\beta}_t}. \quad (9.19)$$

This means that, regardless of what the variance σ_t is, all we need to do is to replace $\varepsilon_\theta(x_t, t)$ with $\varepsilon_\theta(x_t, t) - \beta_t \nabla_x \log p_t(y|x)$, and we can realise conditional generation. As such, under the unified framework of SDE, we could arrive at the generation framework of Classifier-Guidance very easily.

9.7. Classifier-free approach

Finally, we will briefly introduce the Classifier-Free approach. It is actually very simple, we just need to define

$$p(x_{t-1}|x_t, y) = \mathcal{N}(x_{t-1}; \mu(x_t, y), \sigma_t^2 I). \quad (9.20)$$

Using the DDPM results from the previous sections, we can parameterise $\mu(x_t, y)$ generally as

$$\mu(x_t, y) = \frac{1}{\alpha_t} \left(x_t - \frac{\beta_t^2}{\beta_t^2} \varepsilon_\theta(x_t, y, t) \right), \quad (9.21)$$

with the corresponding loss function

$$\mathbb{E}_{x_0, y \sim \tilde{p}(x_0, y), \varepsilon \sim \mathcal{N}(0, I)} [\|\varepsilon\|_2^2 - \varepsilon_\theta(x_t, y, t)]. \quad (9.22)$$

The advantage here is that the additional input y is already included in the training. Theoretically, the more information we have, the easier it is to train. On the other hand, this also brings a major disadvantage, whereby the model needs to be retrained for every new set of conditional input signals.

More specifically, like the Classifier-Guidances, the Classifier-Free approach also introduced a γ -parameter to balance specificity and diversity. Practically, we rewrite Eq. (9.8) using the following modified mean:

$$\mu(x_t) + \sigma_t^2 \gamma \nabla_{x_t} \log p(y|x_t) = \gamma [\mu(x_t) + \sigma_t^2 \nabla_{x_t} \log p(y|x_t)] - (\gamma - 1)\mu(x_t), \quad (9.23)$$

which shows that the Classifier-Free approach is equivalent to the usage of a trainable model to directly regress the term $\mu(x_t) + \sigma_t^2 \nabla_{x_t} \log p(y|x_t)$. In comparison, we can also introduce $\omega = \gamma - 1$ in the Classifier-Guidance and use the following

$$\hat{\varepsilon}_\theta = (1 + \omega)\varepsilon_\theta(x_t, y, t) - \omega\varepsilon_\theta(x_t, t) \quad (9.24)$$

to replace $\varepsilon_\theta(x_t, y, t)$ in the generation. How should we obtain the unconditioned $\varepsilon_\theta(x_t, t)$ in this case? We can introduce a specific input ϕ , which the target image is all the image samples. Once this is added for model training, we can treat that $\varepsilon_\theta(x_t, t) = \varepsilon_\theta(x_t, \phi, t)$.

10. Unified Diffusion Model: The Theoretical Framework

Through the previous sections, we had provided a thorough introduction to the diffusion generative models. It can be seen that there were a lot of theoretical contents involved. Nevertheless, it is clear that in all previous derivations, we have primarily been focusing on diffusion models that dealt with samples that were continuous, and the forward process was built upon a noising process draws noises from a standard normal distribution. In this section, we hope to establish a framework of Unified Diffusion Model (UDM) that may overcome the above limitations, in which

Original blog post see <https://www.spaces.ac.cn/archives/9262>

1. We impose no restriction on the the type of the samples, it can be either discrete or continuous.
2. We impose no restriction on how to realise the forward process. It could be noising, blurring, masking or deletions.
3. We impose no restriction on the nature of the time variable t , it can be either discrete or continuous.
4. It will include all the known results, including DDPM, DDIM, SDE and ODE that had been discussed previously.

It is the aim of this section that we will develop a unified framework that satisfies the above

desires.

10.1. The forward process

From the discussions in the previous sections, it is clear that, in order to build a diffusion model, we will need three components, namely, the forward deconstruction process, the reverse reconstruction process and a loss function for training. In this subsection, we begin our analysis on the forward process.

In the earliest DDPM model, the forward process was established based on the transition probability $p(x_t|x_{t-1})$. In the subsequent works on DDIM, it was realised that both the training targets and the generative process do not have any direct connection with $p(x_t|x_{t-1})$, but the probability $p(x_t|x_0)$. It was also found that, to derive $p(x_t|x_0)$ from $p(x_t|x_{t-1})$ was also rather difficult. Therefore, **a more practical approach is to directly treat $p(x_t|x_0)$ as the forward process.**

The main purpose of establishing $p(x_t|x_0)$ is to build the training data for the diffusion model, so **the basic requirement for $p(x_t|x_0)$ is that sampling from such a distribution must be easy.** For this purpose, we could re-parameterise the forward process as

$$x_t = \mathcal{F}_t(x_0, \varepsilon), \quad (10.1)$$

in which \mathcal{F} is a function of t , x_0 and ε . In particular, ε is a random variable that is sampled from some standard distribution $q(\varepsilon)$, with normal distribution being the most common choice, although other distributions are also feasible. It could be anticipated that Eq. (10.1) is **a generalisation of the forward process that includes all possible transformations from x_0 to x_t , without any restriction on the data type.** Generally, the only restriction is imposed on t , such that the smaller the time variable t is, $\mathcal{F}_t(x_0, \varepsilon)$ will contain more information about x_0 , and it becomes easier to reconstruct x_0 from $\mathcal{F}_t(x_0, \varepsilon)$. Conversely, it becomes more difficult to reconstruct x_0 at large t . At a certain upper limit T , $\mathcal{F}_T(x_0, \varepsilon)$ will contain almost no information about x_0 , from which it becomes literally impossible to reconstruct x_0 .

10.2. The reverse process

In the reverse process, we aim to generate realistic samples through multiple steps of denoising. The key quantity involves is the reverse transition probability $p(x_{t-1}|x_t)$. In general, we have the following:

$$p(x_{t-1}|x_t) = \int p(x_{t-1}|x_t, x_0)p(x_0|x_t)dx_0. \quad (10.2)$$

If x_0 are discrete data, then we can write Eq. (10.2) equivalently by changing the integral into a discrete summation. Since we would like the sampling from $p(x_{t-1}|x_t)$ to be easy to achieve numerically, we will also require the sampling from both $p(x_{t-1}|x_t, x_0)$ and $p(x_0|x_t)$ to be easy. In this case, the sampling from the distribution $p(x_{t-1}|x_t)$ can be achieved via the following procedure:

$$\hat{x}_0 \sim p(x_0|x_t) \quad \& \quad x_{t-1} \sim p(x_{t-1}|x_t, x_0 = \hat{x}_0) \quad \rightarrow \quad p(x_{t-1}|x_t). \quad (10.3)$$

From this formalism, every sampling step $x_t \rightarrow x_{t-1}$ contains two sub-steps:

1. **Estimation:** make an estimate of x_0 from $p(x_0|x_t)$.
2. **Correction:** put the estimated x_0 into $p(x_{t-1}|x_t, x_0)$, which will be incorporated with $p(x_0|x_t)$ to obtain $p(x_{t-1}|x_t)$.

Hence, reverse sampling is an iterative process of estimation followed by correction. By iteratively incorporating the estimate of $x_t \rightarrow x_0$ into the reverse process, we obtain the corrected regression sequence of $x_T \rightarrow \dots \rightarrow x_t \rightarrow x_{t-1} \rightarrow \dots \rightarrow x_0$. This breaks the generative process, which is difficult to be achieved in one-shot, into many sequentially dependent small steps.

10.3. Training target

Of course, this proposed reverse sampling process is rather abstractive, as we do not have any knowledge on the distributions of either $p(x_{t-1}|x_t, x_0)$ or $p(x_0|x_t)$. In this subsection, we will discuss the choice of $p(x_0|x_t)$. Obviously, $p(x_0|x_t)$ is the probabilistic model that allows one

Reminder: The advantage of writing $p(x_{t-1}|x_t)$ in this integral form is that, the conditional probability $p(x_{t-1}|x_t, x_0)$ has a much well defined solution space, and we can approximate $p(x_0|x_t)$ with $q(x_t|x_0)$ from the forward process.

to predict x_0 from x_t , and it must be a distribution that is easy to be calculated and sampled from. When x_0 is continuous data, there are not many choices of suitable distributions, and the most common choice is a normal distribution with a trainable mean:

$$p(x_0|x_t) \approx q(x_t|x_0) = \mathcal{N}(x_0; \mathcal{G}_t(x_t), \bar{\sigma}_t^2 \mathbf{I}). \quad (10.4)$$

To reduce the training cost, we usually do not treat the standard deviation $\bar{\sigma}_t^2$ as a trainable, but using the approaches outlined in Section 7 to estimate it afterward. On the other hand, when x_0 are discrete data, we can use either auto-regressive or non-auto-regressive language models (Seq2Seq) to build our model. Both building the probabilistic models and samplings are relatively much easier for discrete samples. With a concrete form of $q(x_0|x_t)$, the training target is easy to establish. A natural choice is to use the cross-entropy

$$\mathbb{E}_{x_0 \sim \tilde{p}(x_0), x_t \sim p(x_t|x_0)} [-\log q(x_0|x_t)] = \mathbb{E}_{x_0 \sim \tilde{p}(x_0), \varepsilon \sim q(\varepsilon)} [-\log q(x_0|\mathcal{F}_t(x_0, \varepsilon))]. \quad (10.5)$$

As such, this solves the problem of estimating $p(x_0|x_t)$ and designing the corresponding training target. If $q(x_0|x_t)$ follows the normal distribution, then after eliminating the constants, the training target simplifies to

$$\mathbb{E}_{x_0 \sim \tilde{p}(x_0), \varepsilon \sim q(\varepsilon)} \left[\frac{1}{2\sigma_t^2} \|x_0 - \mathcal{G}_t(\mathcal{F}_t(x_0, \varepsilon))\|_2^2 \right]. \quad (10.6)$$

10.4. Conditional probability

Now we are left with the term $p(x_{t-1}|x_t, x_0)$ [in Eq. (10.2)] to be dealt with, which is the probability of estimating x_{t-1} from x_t and x_0 . There is also some room for designing this probability distribution. However, any given distributions for $p(x_{t-1}|x_t, x_0)$ must satisfy the following equality for marginal distribution:

$$\int p(x_{t-1}|x_t, x_0) p(x_t|x_0) dx_t = p(x_{t-1}|x_0). \quad (10.7)$$

Obviously, the simplest choice that satisfies the above constrain is to let

$$p(x_{t-1}|x_t, x_0) = p(x_{t-1}|x_0), \quad (10.8)$$

which means now we have made $p(x_{t-1}|x_t, x_0)$ become completely independent from x_t . In principle, there is no problem of making such a choice. However, the generated outcome based on such a simplification is usually not very good. This is because x_{t-1} now becomes completely dependent upon x_0 , which represents the original sample. In the reverse process, however, we can only sample x_0 approximately from the approximated distribution of $q(x_0|x_t)$. Since $q(x_0|x_t)$ is usually not very accurate, errors in the generative process will continuously accumulate throughout. At the same time, the distribution $p(x_{t-1}|x_0)$ will carry noises in the sampling process, which may significantly affect the estimated \hat{x}_0 , thus worsening the generative outcomes.

Fortunately, under most situations, we can derive a new result based on this simple choice of $p(x_{t-1}|x_t, x_0) = p(x_{t-1}|x_0)$. Based on Eq. (10.1), we have the following:

$$\begin{aligned} x_{t-1} \sim p(x_{t-1}|x_0) &\Leftrightarrow x_{t-1} = \mathcal{F}_{t-1}(x_0, \varepsilon) \\ x_t \sim p(x_t|x_0) &\Leftrightarrow x_t = \mathcal{F}_t(x_0, \varepsilon). \end{aligned} \quad (10.9)$$

Assuming that $\mathcal{F}_t(x_0, \varepsilon)$ is an invertible function of ε , then we have the solution $\varepsilon = \mathcal{F}_t^{-1}(x_0, x_t)$. With this, we can now replace ε in $\mathcal{F}_{t-1}(x_0, \varepsilon)$, which gives us

$$x_{t-1} = \mathcal{F}_{t-1}(x_0, \mathcal{F}_t^{-1}(x_0, x_t)), \quad (10.10)$$

which is equivalent to

$$p(x_{t-1}|x_t, x_0) = \delta(x_{t-1} - \mathcal{F}_{t-1}(x_0, \mathcal{F}_t^{-1}(x_0, x_t))). \quad (10.11)$$

This is a design of $p(x_{t-1}|x_t, x_0)$ that simultaneously depends on both x_0 and x_t , in which some dependency of x_{t-1} on x_0 has been spread to x_t , and we also eliminated the noise, such that we have a more stabilised reversed process to achieve better generative outcomes. Furthermore, if $q(\varepsilon)$ was a normal distribution, we may arrive at a more general result. Based

Reminder: The cross-entropy between two distributions $p(x)$ and $q(x)$ is defined as $-\mathbb{E}_{x \sim p(x)} [\log q(x)] = -\int p(x) \log q(x) dx$, in which $\mathbb{E}_{x \sim p(x)}$ represents the expectation value taken over the true data distribution. In this case, Eq. (10.5) is trying to minimise the discrepancy between the distribution $p(x_0|x_t)$ that is what we try to learn and the true distribution $q(x_t|x_0)$ from training data generated from forward process.

This result is analogous to the general expression for $p(x_{t-1}|x_t, x_0)$ in Eq. (4.7) from DDIM.

on the superposition of normal distributions, we get

$$\mathbf{x}_{t-1} = \mathcal{F}_{t-1}(\mathbf{x}_0, \varepsilon) \quad \Leftrightarrow \quad \mathbf{x}_{t-1} = \mathcal{F}_{t-1} \left(\mathbf{x}_0, \sqrt{1 - \tilde{\sigma}_t^2} \varepsilon_1 + \tilde{\sigma}_t \varepsilon_2 \right). \quad (10.12)$$

In this case, the solution of ε from \mathbf{x}_0 and \mathbf{x}_t can only be used to substitute one of ε_1 and ε_2 . The final sampling process follows:

$$\mathbf{x}_{t-1} = \mathcal{F}_{t-1} \left(\mathbf{x}_0, \sqrt{1 - \tilde{\sigma}_t^2} \mathcal{F}_t^{-1}(\mathbf{x}_0, \mathbf{x}_t) + \tilde{\sigma}_t \varepsilon \right). \quad (10.13)$$

10.5. The thought process

By now, we have completed the building of a unified theoretical framework for the diffusion generative model. In the next section, we will provide some concrete examples of how to derive the results of the existing diffusion models from this unified framework, as well as some new outcomes from the new framework. Here, we will provide some further analysis on the derivation process that has been presented in this section.

Some readers may feel confused with the derivations here. This is because it was established based on author's understanding on all the theoretical details of the diffusion model that have been discussed in the previous nine sections, from which a unified framework may be established. It is not technically challenging but logically not trivial to be understood. Foremost, our goal is to design a unified theoretical framework that satisfies the targets outlined at the beginning of this section. The key in our design is to **understand and control parts of the model that can be freely adjusted and those that must be fixed under certain constraints**.

More specifically, for those readers who are familiar with the diffusion model, they would appreciate that the core idea behind the diffusion model is to “*learn how to reconstruct from destruction*”. As such, we could choose any arbitrary destruction protocols in theory, but must be able to learn how to reconstruct from it. Nevertheless, the destruction process is not completely free of constraints. Generally, one must destruct the samples progressively in order for us to learn how to rebuild them, also progressively. This allows us to design the forward destruction process as presented in Eq. (10.1), where \mathcal{F} can be any destruction protocol and t represents the progress of the destruction process. There is no specific constraints on the type of the original data \mathbf{x}_0 . Finally, ε encodes the random factor that may exist in the forward process. With these, we have established the most generic form of the destruction process.

As for the backward reconstruction process, we first deconstructed $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$ with Eq. (10.2) derived from the probability theory, which can be regarded as a constraint. How did it come up? In retrospect, since the forward process follows $\mathbf{x}_0 \rightarrow \mathbf{x}_t$, then it is natural to match it with the reverse process with $\mathbf{x}_t \rightarrow \mathbf{x}_0$. Eq. (10.2) contains two parts, in which $p(\mathbf{x}_0|\mathbf{x}_t)$ is very clear, it represents the probability of predicting \mathbf{x}_0 from \mathbf{x}_t , which does not have much room for further simplification and can only be acquired from a learnable model. The freedom comes from $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$, which must be a distribution that is easy to be sampled from, while satisfying the constraint imposed by Eq. (10.7) from the probability theory, which had taken the author a rather long time to come up with a satisfactory solution.

11. Unified Diffusion Model: Applications

In the previous section, we have tried to establish a UDM, which allows the application of a more generic diffusion model and data type. In this section, we will be providing some concrete examples to showcase how we may achieve our model design targets.

Original blog post see <https://www.spaces.ac.cn/archives/9271>

11.1. Recap on the general framework

Firstly, UDM builds the forward process by choosing a suitable noise distribution $q(\varepsilon)$ and a transformation \mathcal{F} :

$$\mathbf{x}_t = \mathcal{F}_t(\mathbf{x}_0, \varepsilon), \quad \varepsilon \sim q(\varepsilon). \quad (11.1)$$

Then, through the following partition

$$\hat{\mathbf{x}}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t) \quad \& \quad \mathbf{x}_{t-1} \sim p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0 = \hat{\mathbf{x}}_0), \quad (11.2)$$

we will achieve the sampling $\mathbf{x}_{t-1} \sim p(\mathbf{x}_{t-1}|\mathbf{x}_t)$ in the reverse process. Here $p(\mathbf{x}_0|\mathbf{x}_t)$ is the probability of estimating \mathbf{x}_0 from \mathbf{x}_t , which we can generally approximate it with a

simple distribution such as $q(x_0|x_t)$. In this case, the corresponding training target will be $-\log q(x_0|x_t)$ or some simple variants of it. When x_0 is of the continuous type, we can choose the conditional normal distribution for $q(x_0|x_t)$. When x_0 are discrete data, we could approximate $q(x_0|x_t)$ with either an autoregressive or a non-autoregressive model.

As for the conditional distribution $p(x_{t-1}|x_t, x_0)$, the most basic choice is

$$p(x_{t-1}|x_t, x_0) = p(x_{t-1}|x_0) \Leftrightarrow x_{t-1} = \mathcal{F}_{t-1}(x_0, \varepsilon). \quad (11.3)$$

From this starting point, we can arrive at different results under different assumptions. When $\mathcal{F}_t(x_0, \varepsilon)$ is an invertible function of ε , we can get a solution of $\varepsilon = \mathcal{F}_t^{-1}(x_0, x_t)$, from which we can establish a sampling process with better certainties:

$$x_{t-1} = \mathcal{F}_{t-1}(x_0, \mathcal{F}_t^{-1}(x_0, x_t)). \quad (11.4)$$

Moving one step forward, if $q(\varepsilon)$ is a standard normal distribution, then we have

$$x_{t-1} = \mathcal{F}_{t-1} \left(x_0, \sqrt{1 - \tilde{\sigma}_t^2} \mathcal{F}_t^{-1}(x_0, x_t) + \tilde{\sigma}_t \varepsilon \right). \quad (11.5)$$

11.2. Hot diffusion

Here, we shall first prove that the concept of *hot diffusion* is a special case of UDM, which corresponds to the mainstream diffusion models such as DDPM and DDIM. This idea came from the paper that introduced *cold diffusion*, which will be discussed next.

In the mainstream diffusion models, we are primarily dealing with continuous data, in which the forward process corrupts the data by adding Gaussian noise:

$$x_t = \bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \mathbf{I}), \quad (11.6)$$

Note how such a choice corresponds to an invertible function \mathcal{F} of the noise ε .

and subsequently we choose $\mathcal{N}(x_0; \bar{\mu}(x_t), \bar{\sigma}_t \mathbf{I})$ for $q(x_0|x_t)$. Generally, we do not treat $\bar{\sigma}_t$ as a trainable parameter, such that when we ignore the constant term, we arrive at the following target function:

$$-\log q(x_0|x_t) = \frac{1}{2\bar{\sigma}_t^2} \|x_0 - \bar{\mu}(x_t)\|_2^2. \quad (11.7)$$

If we further introduce the parameterisation of $\bar{\mu}(x_t) = \frac{1}{\bar{\alpha}_t}(x_t - \bar{\beta}_t \varepsilon_\theta(x_t, t))$ and combine with $x_t = \bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon$, then we have

$$-\log q(x_0|x_t) = \frac{\bar{\beta}_t^2}{2\bar{\sigma}_t^2 \bar{\alpha}_t^2} \|\varepsilon - \varepsilon_\theta(\bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon, t)\|_2^2. \quad (11.8)$$

Numerical experiments shows that the performance is even better if we ignore the prefactor, such that the final training target is $\|\varepsilon - \varepsilon_\theta(\bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon, t)\|_2^2$. The choices of the standard derivations follows the discussions in Section 4 and Section 8. Finally, for the distribution of $p(x_{t-1}|x_t, x_0)$, we have

$$\begin{aligned} x_{t-1} &= \bar{\alpha}_{t-1} x_0 + \bar{\beta}_{t-1} \varepsilon \\ &\sim \bar{\alpha}_{t-1} x_0 + \sqrt{\bar{\beta}_{t-1}^2 - \bar{\sigma}_t^2} \varepsilon_1 + \sigma_t \varepsilon_2, \end{aligned} \quad (11.9)$$

in which $\varepsilon, \varepsilon_1$ and $\varepsilon_2 \sim \mathcal{N}(0, \mathbf{I})$. From $x_t = \bar{\alpha}_t x_0 + \bar{\beta}_t \varepsilon$, we have $\varepsilon = (x_t - \bar{\alpha}_t x_0)/\bar{\beta}_t$, which can be used for substituting ε_1 . This leads to the following general form of $p(x_{t-1}|x_t, x_0)$:

$$x_{t-1} = \bar{\alpha}_{t-1} x_0 + \sqrt{\bar{\beta}_{t-1}^2 - \bar{\sigma}_t^2} \frac{x_t - \bar{\alpha}_t x_0}{\bar{\beta}_t} + \bar{\sigma}_t \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \mathbf{I}). \quad (11.10)$$

On the other hand, $\hat{x}_0 \sim p(x_0|x_t)$ implies that

$$\hat{x}_0 = \bar{\mu}(x_t) + \bar{\sigma}_t \varepsilon = \frac{1}{\bar{\alpha}_t} (x_t - \bar{\beta}_t \varepsilon_\theta(x_t, t)) + \bar{\sigma}_t \varepsilon. \quad (11.11)$$

By combining Eq. (11.10) and Eq. (11.11), we get the generating process for all mainstream diffusion models, in which DDPM has chosen $\bar{\sigma}_t = 0$ and $\sigma_t = \bar{\beta}_{t-1} \beta_t / \bar{\beta}_t$, DDIM has chosen

$\bar{\sigma}_t = 0$, $\sigma_t = 0$, whereas in Analytic-DPj. the value of $\bar{\sigma}_t$ had been further optimised.

11.3. Cold diffusion

We now move on to show that the idea of *cold diffusion*, introduced in the paper entitled “*Cold Diffusion: Inverting Arbitrary Image Transforms without Noise*”¹⁸, is also a special case of the UDM. Cold diffusion also focuses on modelling continuous data, and it can be seen from the paper title that, this method aims at building a forward process with any arbitrary noiseless transformations. According to our best knowledge, this is the first work that tries to build a generic forward process, which has inspired us in formulating the UDM.

Cold diffusion builds the forward process with a definitive transformation $x_t = \mathcal{F}(x_0)$. To facilitate the subsequent discussions, let us generalise this forward process further by introducing the noise term as

$$x_t = \mathcal{F}(x_0) + \sigma \varepsilon, \quad \varepsilon \sim q(\varepsilon). \quad (11.12)$$

Here, \mathcal{F} represents any arbitrary way to destruct the raw data. For images, this could be either blurring, painting, snowification, and so on. For a definitive transformation, we only need to let $\sigma \rightarrow 0$. Subsequently, we choose $q(x_0|x_t)$ as a normal distribution with ℓ_1 -norm:

$$q(x_0|x_t) = \frac{1}{\mathcal{Z}(\tau)} \int e^{-\|x_0 - \mathcal{G}_t(x_t)\|_1/\tau} dx_0, \quad (11.13)$$

in which $\mathcal{Z}(\tau)$ is the normalisation constant. If we choose τ to be a fixed constant, then ignore the constant term, we have $-\log q(x_0|x_t) \propto \|x_0 - \mathcal{G}_t(x_t)\|_1$. Combining this with $x_t = \mathcal{F}(x_0)$, we have the following training target, which is to minimise

$$\|x_0 - \mathcal{G}_t(\mathcal{F}(x_0))\|_1. \quad (11.14)$$

For the reverse process, the variance in $q(x_0|x_t)$ is completely ignored, *i.e.* we let $\tau \rightarrow 0$. As such, we have $\hat{x}_0 = \mathcal{G}_t(x_t)$. If we further choose $p(x_{t-1}|x_t, x_0) = p(x_{t-1}|x_0)$, *i.e.* $x_{t-1} = \mathcal{F}_{t-1}(x_0) + \sigma \varepsilon$, then we substitute in \hat{x}_0 and take the limit of $\sigma \rightarrow 0$, we reached

$$\hat{x}_0 = \mathcal{G}_t(x_t), \quad x_{t-1} = \mathcal{F}_{t-1}(\hat{x}_0). \quad (11.15)$$

This result corresponds to the **naïve sampling** proposed in the paper. On the other hand, if we started from $x_t = \mathcal{F}(x_0) + \sigma \varepsilon$ and solved $\varepsilon = (x_t - \mathcal{F}(x_0))/\sigma$, before we substitute this into $x_{t-1} = \mathcal{F}_{t-1}(x_0) + \sigma \varepsilon$, we arrived at the result of **improved sampling**:

$$\hat{x}_0 = \mathcal{G}_t(x_t), \quad x_{t-1} = x_t + \mathcal{F}_{t-1}(\hat{x}_0) - \mathcal{F}_t(\hat{x}_0). \quad (11.16)$$

Overall, cold diffusion was successful in realising a generic forward diffusion process. However, it over emphasised the ideal of “without noise”, which has led to shortages that are almost impossible to overcome theoretically. For example, when using a blurring operation in the forward process, it blurs an $(w \times w \times 3)$ dimensional image data into a 3-dimensional vector with a definitive transformation. As the same time, another definitive reverse process is also applied to rebuild the $3w^2$ -dimensional image from a 3-dimensional vector. It is unavoidable that a lot of informations would have been lost in this process, which limited the resolutions of the reconstructed images. To solve this problem, we cannot eliminate the noise term in the diffusion process. This is because noises represent uncertainties, which promote a “many-to-one” mapping in the forward process. In another word, it allows information loss. As a matter of fact, this problem has been made aware of in the implementation of the cold diffusion, in which a $3w^3$ -dimensional random noise had been injected into the 3-dimensional vector. Numerical experiments showed that this has indeed let to improved image qualities.

11.4. The editing model

The two examples given above are both dealing with continuous data. However, we have mentioned that the purpose of UDM is to not restricting the data type. Therefore, in this subsection, we will discuss an example with discrete data. We aim to show that the text generation model for editing can also be viewed as a special case of UDM.

For simplicity, we consider a sentence with a fixed length ℓ . It is possible also to deal with variable length but it is slightly more difficult. With this, we define the forward process $x_t = \mathcal{F}(x_0, \varepsilon)$ as ‘random substitution’, *i.e.*, we randomly select t tokens in the sentence and substitute them with other tokens. When $t = \ell$, x_t is a combination of ℓ random tokens. In

this case, $q(x_0|x_t)$ is a model that predicts the original token sequence from a sequence of fully randomised tokens, which can be trained with regressive models and cross-entropy as the loss function. Note that now the forward process $\mathcal{F}(x_0, \varepsilon)$ is certainly not invertible, as there are multiple pathways to evolve x_0 to x_t . As such, we can only apply $p(x_{t-1}|x_t, x_0) = p(x_{t-1}|x_0)$, which leads to the following generative process:

1. Randomly choose ℓ tokens as the initial x_ℓ .
2. Make an estimate of \hat{x}_0 from $q(x_0|x_t)$.
3. Randomly choose $\ell - 1$ tokens to substitute other tokens to form x_{t-1} .
4. Repeat steps 2 and 3 until a final x_0 is obtained.

However, this algorithm will not lead to very good outcomes. This is because the estimate from step 2 may be significantly destroyed by the noises from step 3. In order to improve the sampling quality, we need to use a better sampling procedure, this requires $\mathcal{F}(x_0, \varepsilon)$ to be an invertible function of ε . This means that from a give pair of x_0 and x_t , we can figure out how the forward transformation was achieved. To achieve this, we redefined the forward process as “*randomly selecting t tokens in the sentence and substitute them with different ones*”. The main difference is now we forces the new tokens must substitute with ones that are different from the original ones. Otherwise, one might have sampled the same tokens as before. With this restriction, we could directly compare x_0 and x_t to figure out what the transformation was. In this case, we change step 3 into a transformation from \hat{x}_0 to x_t :

1. Randomly choose ℓ tokens as the initial x_ℓ .
2. Make an estimate of \hat{x}_0 from $q(x_0|x_t)$, under the condition that t tokens in \hat{x}_0 must be different from x_t , which may be easily achievable with a non-autoregressive model.
3. Randomly choose a token in x_t that is different from \hat{x}_0 , substitute it with the token of \hat{x}_0 in the corresponding positions to form x_{t-1} .
4. Repeat steps 2 and 3 until a final x_0 is obtained.

With this sampling procedure, in every sampling step, the part in \hat{x}_0 that is identical to x_t can be kept. Comparing x_{t-1} and x_t , only one token would have been changed, so we have a much more stable progressive generative process.

11.5. The masking model

Here, we provide another example to further help the readers in understanding the models. Similar to above, we consider a sentence with ℓ tokens, and define the forward process $x_t = \mathcal{F}(x_0, \varepsilon)$ as random masking, *i.e.*, we select t tokens and randomly replace them with [MASK], in which $t \leq \ell$. When $t = \ell$, x_t is a sentence with ℓ [MASK].

In this case, $q(x_0|x_t)$ is a model that predicts the original sequence from a sequence that contains [MASK]. This can be achieved with non-autoregressive models such as BERT based on MLM, with cross-entropy as the loss function. The basic generative process is:

1. Using a sequence containing ℓ [MASK] as the starting x_ℓ .
2. Sample \hat{x}_0 from $q(x_0|x_t)$.
3. Randomly selecting $t - 1$ tokens in \hat{x}_0 and replace it with [MASK], this gives x_{t-1} .
4. Repeat steps 2 and 3 until a final x_0 is obtained.

Noticing that $\mathcal{F}_t(x_0, \varepsilon)$ is an invertible function of ε , meaning that we can figure out how the sequence has been transformed from x_0 to x_t , *i.e.* which tokens have been changed into [MASK]. In this case, we have the following improved sampling procedure:

1. Using a sequence containing ℓ [MASK] as the starting x_ℓ .
2. Sample \hat{x}_0 from $q(x_0|x_t)$, where we only need to sample for the tokens that were originally the [MASK], and do not change those that were not [MASK].
3. Select $t - 1$ positions randomly from x_t in which the t positions were occupied by [MASK]. Changing those into the corresponding token at \hat{x}_0 and make this to be x_{t-1} .
4. Repeat steps 2 and 3 until a final x_0 is obtained.

In which step 2 and 3 can be combined into a single step:

Randomly selecting a position from t possible positions that are occupied by [MASK] in x_t , and according to the probability $q(x_0|x_t)$, sample a token and replace the [MASK] at that position, giving a new x_{t-1} .

This is basically the same as the Gibbs sampling in MLM model. From these two examples, we could appreciate that many models based on progressive changes can all be described by the UDM, or equivalently, using USM to guide the design of new models based on progressive changes.

11.6. The encoding model

In all the models described above, there exists no trainable parameters in the forward process, meaning these forward processes are pre-designed. This is not completely necessary, which we can generalise the forward diffusion in DDPM into

$$x_t = \bar{\alpha}_t \mathcal{F}(x_0) + \bar{\beta}_t \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, I), \quad (11.17)$$

in which $\mathcal{F}(x_0)$ is the encoding model for x_0 , which can contain trainable parameter. The corresponding training target is

$$-\log q(x_0|x_t) = -\log q(x_0|\bar{\alpha}_t \mathcal{F}(x_0) + \bar{\beta}_t \varepsilon). \quad (11.18)$$

The only difference is that now \mathcal{F} is also parameterised. The reverse sampling process is also similar, except now the last step $\hat{x}_0 \sim q(x_0|x_1)$ will return \hat{x}_0 directly. Most importantly, because now we have an encoding model $\mathcal{F}(x_0)$, the input data x_0 could be either continuous or discrete, which is similar to VAE that transforms the data distribution into the normal distribution of the latent variables.

12. Conquering the Diffusion ODE

In Section 5, we aimed at developing an understanding of the diffusion model from the perspective of SDE. Subsequently, in Section 6, we showed that, in the diffusion model that was established based on SDE, there exists an underlying ODE model. Coincidentally, in Section 4 where we discussed DDIM, we have demonstrated that the original DDPM model which performs stochastic samples, also implicitly contains a DDIM model, in which the sampling procedure is deterministic, and its continuous limit also corresponds to an ODE.

Thinking of the above arguments in detail, it is not difficult to discover that both DDPM→DDIM and SDE→ODE are procedures that convert a stochastic sampling model into a deterministic one. On the other hand, if our initial goal was directly based on ODE, then these procedures of establishing a deterministic sampling framework from a stochastic starting point are obviously too cumbersome. In this section, we will try to provide a direct derivation for the diffusion model based on ODE, and demonstrate its linkages with the Jacobian and the diffusion equations for heat conduction.

12.1. Differential equation

The core idea behind the generative models, such as GAN, is to find a deterministic transformation that is capable of transforming a random variable sampled from a simple probabilistic distribution, such as the standard normal distribution, into a realistic sample that follows a more data-specific distribution. The flow mode is another type of generative model, which follows the opposite design concept. It tries to find a deterministic transformation that transforms the unknown data distribution into a simple one, then it solves for the corresponding

Original blog post see <https://www.spaces.ac.cn/archives/9280>

inverse transformation to establish the generative model.

The traditional flow model achieves this invertible transformation through a careful design of the coupling layers. However, it was later realised that the same type of transformation can be achieved with the usage of differential equations, and it is theoretically more elegant. This leads to the new field of neural ODE which builds generative models by combining neural network with differential equations.

Consider the first-order ordinary differential equation on $\mathbf{x}_t \in \mathbb{R}^d$:

$$\frac{d\mathbf{x}_t}{dt} = f_t(\mathbf{x}_t). \quad (12.1)$$

Suppose $t \in [0, T]$, then with a given \mathbf{x}_0 , (under some conditions that are easy to realise) we could solve for \mathbf{x}_T deterministically. It implies that this ODE describes a definitive transformation from \mathbf{x}_0 to \mathbf{x}_T . More specifically, this transformation will be invertible, meaning that we could also solve the ODE in the reverse order to obtain the transformation from \mathbf{x}_T to \mathbf{x}_0 . Therefore, differential equations lead to an elegant theoretical solution for constructing an invertible transformation.

12.2. The Jacobian

Similar to the diffusion models discussed before, we treat \mathbf{x}_0 as a sample from the input data, and \mathbf{x}_T as a sample drawn from a simple probability distribution. Our goal here is to obtain a transformation from \mathbf{x}_0 to \mathbf{x}_T through the differential equation.

To do this, we first write out the finite-difference form for the differential equation [Eq. (12.1)] as:

$$\mathbf{x}_{t+\Delta t} - \mathbf{x}_t = f_t(\mathbf{x}_t)\Delta t. \quad (12.2)$$

Since this is a definitive transformation, we have the following expression for the probability distributions at \mathbf{x}_t and $\mathbf{x}_{t+\Delta t}$:

$$p_t(\mathbf{x}_t)d\mathbf{x}_t = p_{t+\Delta t}(\mathbf{x}_{t+\Delta t})d\mathbf{x}_{t+\Delta t} = p_{t+\Delta t}(\mathbf{x}_{t+\Delta t}) \left| \frac{\partial \mathbf{x}_{t+\Delta t}}{\partial \mathbf{x}_t} \right| d\mathbf{x}_t. \quad (12.3)$$

Here, $\frac{\partial \mathbf{x}_{t+\Delta t}}{\partial \mathbf{x}_t}$ is the Jacobian matrix for the transformation and $|\cdot|$ takes its determinant. By taking the derivatives on both sides of Eq. (12.2), we have

$$\frac{\partial \mathbf{x}_{t+\Delta t}}{\partial \mathbf{x}_t} = \mathbf{I} + \frac{\partial f_t(\mathbf{x}_t)}{\partial \mathbf{x}_t} \Delta t, \quad (12.4)$$

and based on the definition on the derivative of a determinant, we have

$$\left| \frac{\partial \mathbf{x}_{t+\Delta t}}{\partial \mathbf{x}_t} \right| \approx 1 + \text{Tr} \left[\frac{\partial f_t(\mathbf{x}_t)}{\partial \mathbf{x}_t} \right] \Delta t = 1 + \nabla_{\mathbf{x}_t} \cdot f_t(\mathbf{x}_t) \Delta t \approx \exp [\nabla_{\mathbf{x}_t} \cdot f_t(\mathbf{x}_t) \Delta t]. \quad (12.5)$$

By taking the logarithm to both sides of Eq. (12.3) with the use of Eq. (12.5), we obtain the following equation for the change in the probability distributions in the finite-difference form:

$$\log p_{t+\Delta t}(\mathbf{x}_{t+\Delta t}) - \log p_t(\mathbf{x}_t) \approx -\nabla_{\mathbf{x}_t} \cdot f_t(\mathbf{x}_t) \Delta t. \quad (12.6)$$

12.3. Taylor approximation

Now suppose $p_t(\mathbf{x}_t)$ defines a series of probability density functions that vary continuously with respect to the time variable t , in which $p_0(\mathbf{x}_0)$ is the true sample distribution and $p_T(\mathbf{x}_T)$ is a simple distribution. When both Δt and $\mathbf{x}_{t+\Delta t} - \mathbf{x}_t$ are sufficiently small, we can write the following first-order Taylor approximation for Eq. (12.6):

$$\log p_{t+\Delta t}(\mathbf{x}_{t+\Delta t}) - \log p_t(\mathbf{x}_t) \approx (\mathbf{x}_{t+\Delta t} - \mathbf{x}_t) \nabla_{\mathbf{x}_t} \cdot \log p_t(\mathbf{x}_t) + \Delta t \frac{\partial}{\partial t} \log p_t(\mathbf{x}_t). \quad (12.7)$$

Substituting the expression for $\mathbf{x}_{t+\Delta t} - \mathbf{x}_t$ from Eq. (12.2) and compare the result to Eq. (12.6), we get the equation which $f_t(\mathbf{x}_t)$ must satisfy:

$$-\nabla_{\mathbf{x}_t} \cdot f_t(\mathbf{x}_t) = f_t(\mathbf{x}_t) \nabla_{\mathbf{x}_t} \cdot \log p_t(\mathbf{x}_t) + \Delta t \frac{\partial}{\partial t} \log p_t(\mathbf{x}_t). \quad (12.8)$$

In other words, any arbitrary function $f_t(\mathbf{x}_t)$ that satisfies the above equation will be a suitable candidate for constructing the ODE of the form shown in Eq. (12.1), solving which will allow us to achieve the transformation between the data and a simple probability distribution. Rearranging Eq. (12.8), we get the following

$$\frac{\partial}{\partial t} \log p_t(\mathbf{x}_t) = -\nabla_{\mathbf{x}_t} (f_t(\mathbf{x}_t) p_t(\mathbf{x}_t)), \quad (12.9)$$

which is actually the special case of the Fokker-Planck equation introduced in Section 6 when $g_t = 0$.

12.4. The equation of heat conduction

Now consider a solution to Eq. (12.9) that is of the following form:

$$f_t(\mathbf{x}_t) = -D_t(\mathbf{x}_t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t), \quad (12.10)$$

in which $D_t(\mathbf{x}_t)$ could be either a matrix or a scalar, depending on the complexity of the problem. Why should we consider this particular form of the solution? The author's initial intention was to match with the results from DDIM. It was later found that, after some generalisation, it can be matched with the diffusion model to be discussed below. In retrospect, if $D_t(\mathbf{x}_t)$ is a non-negative scalar function and if we substitute Eq. (12.10) into Eq. (12.2), we found that it led to a form that is similar to the case of gradient descent, in which the transformation from \mathbf{x}_0 to \mathbf{x}_T is a process that progressively searches the regions of low probabilities, and conversely, the transformation from \mathbf{x}_T to \mathbf{x}_0 gradually leads to the searches over the regions of high probabilities. This follows our basic intuitions, which may be considered as the inspiration behind Eq. (12.10).

Now by substituting Eq. (12.10) into Eq. (12.9), we get

$$\frac{\partial}{\partial t} p_t(\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \cdot (D_t(\mathbf{x}_t) \nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t)). \quad (12.11)$$

This is the diffusion equation in differential equations. Here, we will only consider the most simplified case, in which $D_t(\mathbf{x}_t)$ is a scalar function that is independent on \mathbf{x}_t , such that the diffusion equation simplifies into

$$\frac{\partial}{\partial t} p_t(\mathbf{x}_t) = D_t \nabla_{\mathbf{x}_t}^2 p_t(\mathbf{x}_t). \quad (12.12)$$

This gives the equation for heat conduction, which forms the basis for our subsequent discussions.

12.5. Solving for the probability distribution

Using the Fourier transformation, we can convert the equation for heat conduction into an ODE, from which we can solve for $p_t(\mathbf{x}_t)$, and the result is

$$p_t(\mathbf{x}_t) = \int \frac{1}{(2\pi\sigma_t^2)^{\frac{d}{2}}} \exp\left(-\frac{\|\mathbf{x}_t - \mathbf{x}_0\|_2^2}{2\sigma_t^2}\right) p_0(\mathbf{x}_0) d\mathbf{x}_0 = \int \mathcal{N}(\mathbf{x}_t; \mathbf{x}_0, \sigma_t^2 \mathbf{I}) p_0(\mathbf{x}_0) d\mathbf{x}_0, \quad (12.13)$$

in which $\sigma_t^2 = 2 \int_0^t D_s ds$ or $D_t = \dot{\sigma}_t \sigma_t$ (with $\sigma_0 = 0$). From which we see that the solution is a Gaussian mixture model with $p_0(\mathbf{x}_0)$ as the initial distribution.

The Procedure for solving Eq. (12.12) Here, we shall briefly introduce the procedure for solving the equation of heat conduction, which can be skipped by the readers who are not interested or who are familiar with the equation.

It is not difficult to solve Eq. (12.12) with the use of Fourier transformation. Taking the Fourier transformation to the variable \mathbf{x}_t on both sides of the equation and use the fact that $\nabla_{\mathbf{x}_t} \rightarrow i\omega$, the result is:

$$\frac{\partial}{\partial t} \mathcal{F}_t(\omega) = -D_t \omega^2 \mathcal{F}_t(\omega). \quad (12.14)$$

This is an ordinary differential equation with the following solution:

$$\mathcal{F}_t(\omega) = \mathcal{F}_0(\omega) \exp\left(-\frac{\sigma_t^2 \omega^2}{2}\right), \quad (12.15)$$

in which $\sigma_t^2 = 2 \int_0^t D_s ds$, and $\mathcal{F}_0(\omega)$ is the Fourier transformation of $p_0(x_0)$. Now we take the inverse Fourier transformation on both size, for which $\mathcal{F}_t(\omega)$ transforms back to $p_t(x_t)$, $\mathcal{F}_0(\omega)$ transforms back to $p_0(x_0)$ and $\exp(-\sigma_t^2 \omega^2/2)$ corresponds to the normal distribution $\mathcal{N}(x_t; x_0, \sigma_t^2 \mathbf{I})$. The final form of the solution given in Eq. (1.12) can be obtained by following the convolution property of the Fourier transformation.

12.6. Completing the model design

Let us briefly summarise the results so far. From solving the equation for heat conduction, we have obtained

$$p_t(x_t) = \int \mathcal{N}(x_t; x_0, \sigma_t^2 \mathbf{I}) p_0(x_0) dx_0, \quad (12.16)$$

which corresponds to the solution for the following equation

$$\frac{dx_t}{dt} = -\dot{\sigma}_t \sigma_t \nabla_{x_t} \log p_t(x_t). \quad (12.17)$$

This gives a definitive transformation from $p_0(x_0)$ to $p_T(x_T)$. If $p_T(x_T)$ is a distribution that is easy to sample from, and we know how to compute the score function $\nabla_{x_t} \log p_t(x_t)$, then we can perform a random sampling from $x_T \sim p_T(x_T)$, followed by solving the differential equation in the reverse order, to synthesise new sample $x_0 \sim p_0(x_0)$.

The first equation is, when is $p_T(x_T)$ a distribution that is easy to sample from? Based on Eq. (12.16), we know that

$$x_T \sim p_T(x_T) \Leftrightarrow x_T = x_0 + \sigma_T \varepsilon, \quad x_0 \sim p_0(x_0), \quad \varepsilon \sim \mathcal{N}(0, \mathbf{I}). \quad (12.18)$$

When σ_T is sufficiently large, the influence of x_0 upon x_T will be negligible, at which point we can simply take

$$x_T \sim p_T(x_T) \Leftrightarrow x_T = \sigma_T \varepsilon \quad \varepsilon \sim \mathcal{N}(0, \mathbf{I}), \quad (12.19)$$

which makes $p_T(x_T)$ a distribution that is easy to sample from. As such, the choice of σ_t should follow the general rule of $\sigma = 0$, $\sigma_T \gg 1$, with σ_t being a continuous growing function over $[0, T]$.

The second question is, how should we compute $\nabla_{x_t} \log p_t(x_t)$? This is the same as the score-matching that was discussed in Section 5, in which we use a neural network $s_\theta(x_t, t)$ to surrogate it, with the following training target:

$$\mathbb{E}_{x_0, x_t \sim \mathcal{N}(x_t; x_0, \sigma^2 \mathbf{I}) p_0(x_0)} \left[\|s_\theta(x_t, t) - \nabla_{x_t} \log \mathcal{N}(x_t; x_0, \sigma^2 \mathbf{I})\|_2^2 \right]. \quad (12.20)$$

This is called the conditional score matching, which has been derived before and will not be repeated here.

In summary, we have provided a different approach to derive the ODE-based diffusion model. Starting from the ODE, by combining the Jacobian we obtained the first-order approximation to the probability distribution. By comparing this approximation to that obtained directly from the first-order Taylor expansion, we obtained the equation that the $f_t(x_t)$ term should satisfy. Subsequently, this equation can be transformed into the equation of diffusion or heat conduction, that is to be solved. Compared to the previous derivation from the SDE and FP equations, the procedure presented here should be much more straight forward.

13. From the Law of Gravitational Force to Diffusion Model

For many readers, the diffusion generative model is probably the first deep-learning model that they have encountered, which has utilised so many different mathematical methods. In this series of articles, we have demonstrated the strong connection between the diffusion models and many different fields in mathematics, including mathematical analysis, probability theory,

Original blog post see <https://www.spaces.ac.cn/archives/9305>

statistics, ordinary and stochastic partial differential equations. It is reasonable to say, even for researchers in the field of pure theory, such as mathematical physics, it will be very likely for them to find their skills useful in the development of diffusion generative models.

In this section, we shall introduce another diffusion model that has a strong link with mathematical physics, that is the ODE-based diffusion model that was inspired by the law of the gravitational force. It came from the paper entitled “**Poisson Flow Generative Model**”¹⁹ (PFGM for short), which has provided a new perspective for constructing an ODE-based diffusion model.

13.1. The gravitational force

We have all learned about the law of gravitational force in high school. Basically, it states that, the attractive force between two point masses, is proportional to the product of their masses and the inverse square of the distance between them. Here, we ignore the masses and the constant of the gravitational field, but focus on the relationship between the gravitational force with distance and direction. Assuming the source of this attractive force is at \mathbf{y} , then the force experienced by an object at \mathbf{x} is given by:

$$\mathbf{F}(\mathbf{x}) = -\frac{1}{4\pi} \frac{\mathbf{x} - \mathbf{y}}{\|\mathbf{x} - \mathbf{y}\|_2^3}. \quad (13.1)$$

We can further ignore the factor of $1/4\pi$ for the time being, as it will not affect our subsequent discussions. Strictly speaking, the above equation describes the gravitational field in the three-dimensional space. For the d -dimensional case, the corresponding gravitational field is

$$\mathbf{F}(\mathbf{x}) = -\frac{1}{S_d(1)} \frac{\mathbf{x} - \mathbf{y}}{\|\mathbf{x} - \mathbf{y}\|_2^d}, \quad (13.2)$$

in which $S_d(1)$ is the surface area of a d -dimensional hypersphere. This expression corresponds to the gradient of the Green’s function for the d -dimensional Poisson distribution, thus the word “Poisson” in the title of the paper mentioned above.

13.2. Following the lines of the gravitational field

If there exist multiple sources of attractions, we just need to add up the attractive forces from different sources, based on the additivity of the forces. Generally, we observe that, for most of the field lines, they start from infinitely far away, and end at the attractors. Now, an intuitive but smart idea arises. If every source of the gravitational force represents a real sample to be generated, then for an arbitrary far away point, it will evolve into a real sample if it moves by following the field line. This is the key genius thought behind the article “Poisson Flow generative Model”.

13.3. Equivalent centre-of-mass

Nevertheless, to convert a smart idea into a workable model, there are still many details that will need to be filled in. For instance, when we say “an arbitrarily far away point”, which corresponds to the prior distribution in the diffusion model, how far is really far enough, and how should we sample an arbitrary point? Apparently, if the sampling procedure is too complicated, then it becomes valueless.

Fortunately, there exists an equivalence property for the gravitational field: *The gravitational field experienced by a point that is infinitive far away from multiple sources of attraction, is equivalent to the field generated from a single centre-of-mass with the sum of all masses.*

This is saying that, when the distance is sufficiently large, we can treat the field as originating from a single point source of attraction. What is the property of the gravitational field from a single point mass? It is isotropic in space! This means that when the radius is sufficiently large, we can consider the gravitational field lines will pass the surface of the sphere centred at the point mass evenly. Therefore, for the procedure of initial sampling, it should be performed on the surface of a hypersphere with a sufficiently large radius, and we shall discuss later how large is sufficiently large.

13.4. Mode Collapsing

Has this completed the generative model? Not yet! The isotropic nature of the gravitational field has made the initial sampling easy. However, it can also cause the gravity sources cancelling each other, leading to the so-called phenomenon of “mode-collapsing”. More specifically, we consider the field distribution around a spherical shell, in which it can be observed

that, outside the shell, the gravitational field distributes isotropically. However, inside the shell, the field is empty! It shows that the fields cancel each other, which essentially forms a vacuum inside the shell.

the cancellation of the field implies that, because the gravitational forces form uniformly distributed sources on a spherical shell cancel out each other, it is equivalent to the case where these sources did not exist in the first place. On the other hand, the core idea of the generative model discussed here, is to let a point far away to move along the gravitational field until it reaches a particular source. If, the gravitational field cancels each other, it means that we may never reach certain sources of the gravitational field, or equivalently, we may never be able to synthesise some of the realistic samples. In this case, we lost the diversity in the generative outcomes. This is the phenomenon of mode collapse.

13.5. Adding one extra dimension

It seems like the problem of mode collapsing is unavoidable anyhow. This is because, when we are building the generative model, we usually assume that the real samples follow a continuous distribution. In this case, we can always choose a sphere, even though the samples do not distribute uniformly on its surface, we can always find a subset from this distribution that is uniformly distributed, such that the forces from this subset cancel out each other. This is equivalent to saying this subset of samples being non-existent, leading to the problem of mode collapse.

Does this mean we are stuck? No, as this brings to the second genius idea behind PFGM, that is, to add an one extra dimension.

From our previous analysis, the problem of mode collapse become unavoidable, is because the continuous sample distribution which we assumed in the first place has led to the situation in which the isotropy in the field distribution becomes unavoidable. In order to avoid the mode collapse, we need to try avoiding isotropy in the sample distribution. However, the sample distribution is our target distribution, which we cannot change. Nevertheless, we can add one extra dimension to the samples. If now we tackle the problem in the $(d+1)$ -dimensional space, we can no longer have isotropy on a place. Using a low dimensional analogous example, we know that a circle in the two-dimensional space is isotropic. However, in the three-dimensional space, only a sphere is isotropic, a circle that is isotropic in the two-dimensional space is no longer isotropic in the three-dimensional space.

In this regard, let us assume that the real samples that we would like to generate are $\mathbf{x} \in \mathbb{R}^d$, and we can introduce another new dimension t , such that the sample data now change into the form of $(\mathbf{x}, t) \in \mathbb{R}^{d+1}$. Correspondingly, the original sample follows a distribution of $\mathbf{x} \sim \tilde{p}(\mathbf{x})$. Now the distribution becomes $(\mathbf{x}, t) \sim \tilde{p}(\mathbf{x})\delta(t)$, in which $\delta(t)$ is the Dirac distribution. This is a transformation that places the real sample onto the $t = 0$ plane in the $(d+1)$ -dimensional space. With this, in the $(d+1)$ -dimensional space, all the realistic samples will always have $t = 0$, which cannot be isotropically distributed any more.

13.6. Everything clicked

From the first look of it, adding one extra dimension is only a small mathematical trick. However, as one thinks more about it, we could find more brilliant points behind this approach. Many problems that are difficult to be solved in the d -dimensional space have suddenly become easy when we moved to the $(d+1)$ -dimension.

According to Eq. (13.2) and the linear superposition of the gravitational fields, the corresponding field in the $(d+1)$ -dimensional space is

$$\begin{aligned} \mathbf{F}(\mathbf{x}, t) &= -\frac{1}{S_{d+1}(1)} \iint \frac{(\mathbf{x} - \mathbf{x}_0, t - t_0)}{(\|\mathbf{x} - \mathbf{x}_0\|^2 + (t - t_0)^2)^{(d+1)/2}} \delta(t_0) \tilde{p}(\mathbf{x}_0) d\mathbf{x}_0 dt_0 \\ &= -\frac{1}{S_{d+1}(1)} \iint \frac{(\mathbf{x} - \mathbf{x}_0, t)}{(\|\mathbf{x} - \mathbf{x}_0\|^2 + t^2)^{(d+1)/2}} \delta(t_0) \tilde{p}(\mathbf{x}_0) d\mathbf{x}_0 \\ &\triangleq (\mathbf{F}_x, \mathbf{F}_t), \end{aligned} \quad (13.3)$$

in which \mathbf{F}_x is the first d components of $\mathbf{F}(\mathbf{x}, t)$ and \mathbf{F}_t is its $(d+1)$ -th component. We will discuss later on how to learn $\mathbf{F}(\mathbf{x}, t)$. Assuming that we have already know the form of $\mathbf{F}(\mathbf{x}, t)$, then we need to move along the field line, meaning that the trajectories of our sample movement should always follow the direction of $\mathbf{F}(\mathbf{x}, t)$, *i.e.*,

$$(d\mathbf{x}, dt) = (\mathbf{F}_x, \mathbf{F}_t)d\tau \quad \Rightarrow \quad \frac{d\mathbf{x}}{dt} = \frac{\mathbf{F}_x}{\mathbf{F}_t}. \quad (13.4)$$

This is the differential equation that we will need for the generative model. Before, in the d -dimensional space, apart from the problem of mode collapse, when to stop the movement is also a problem. Intuitively, we just need to follow the field line and stop when we have collided with the sample. But how do we judge whether we have collided with a real sample? This is not a trivial question to answer. However, in the $(d+1)$ -dimensional space, all the real samples lie on the plane of $t = 0$, hence we should stop when $t = 0$ is reached.

As for the prior distribution, based on the discussions above, it should be a uniform distribution on a $(d+1)$ -dimensional hypersphere with a sufficiently large radius. However, since we use $t = 0$ as our stop signal, then we could fix a sufficiently large τ , in the order of 40 to 100, and then sample on the $t = T$ plane. In this case, the generative process becomes solving the equation of motion from $t = T$ to $t = 0$ as given by the differential equation. This gives well-defined starting and end points for the generative process.

If we perform sampling on the $t = T$ place, it will not be uniform. In fact, we have the following:

$$p_{prior}(\mathbf{x}) \propto \frac{1}{(\|\mathbf{x}\|^2 + T^2)^{(d+1)/2}}. \quad (13.5)$$

It can be seen that the probability density is only dependent on the modulus $\|\mathbf{x}\|$. Hence, for sampling, we need to first sample the modulus from a given distribution, then we perform uniform sampling for different directions and combine the two together. For sampling the modulus, let $r = \|\mathbf{x}\|$ and converts to hyperspherical harmonics, then we have $p_{prior}(r) \propto r^{d-1}(r^2 + T^2)^{(d+1)/2}$, from which we could sample from this distribution with the inverse cumulative distribution function.

13.7. Training the field

Now, we have the prior distribution and the differential equation, all we have left is to develop a protocol to train the function of a vector field $F(\mathbf{x}, t)$. Based on the differential equation (13.4), it can be seen that the movement of the particle is only dependent on the relative value of the vector field. As such, the results will not be affected if we shrink or expand the vector field. Based on Eq. (13.3), we can write the vector field as an expectation value:

$$\mathbf{F}(\mathbf{x}, t) = \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0)} \left[-\frac{(\mathbf{x} - \mathbf{x}_0, t)}{(\|\mathbf{x} - \mathbf{x}_0\|^2 + t^2)^{(d+1)/2}} \right]. \quad (13.6)$$

Using the following equality that we had applied many times:

$$\mathbb{E}_{\mathbf{x}} = \underset{\boldsymbol{\mu}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}} [\|\mathbf{x} - \boldsymbol{\mu}\|^2], \quad (13.7)$$

we can introduce a trainable function $s_{\theta}(\mathbf{x}, t)$ to learn $\mathbf{F}(\mathbf{x}, t)$ with the following target function:

$$\mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0)} \left[\left\| s_{\theta}(\mathbf{x}, t) + \frac{(\mathbf{x} - \mathbf{x}_0, t)}{(\|\mathbf{x} - \mathbf{x}_0\|^2 + t^2)^{(d+1)/2}} \right\|^2 \right]. \quad (13.8)$$

However, we still need to perform sampling on \mathbf{x} and t in order to evaluate the above target function, and there is no clear definition on how this should be done. This is a major character behind PFGM, it defined the reverse generative process, without the need of defining the forward process. This one-step sampling is practically equivalent to the forward process. As such, the original paper considered an approach to construct samples of \mathbf{x} and t through perturbing each real sample as:

$$\mathbf{x} = \mathbf{x}_0 + \|\varepsilon_{\mathbf{x}}\|(1 + \tau)^m \mathbf{u}, \quad t = \|\varepsilon_t\|(1 + \tau)^m, \quad (13.9)$$

in which $(\varepsilon_{\mathbf{x}}, \varepsilon_t) \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_{(d+1) \times (d+1)})$ and $m \sim \mathcal{U}[0, M]$. \mathbf{u} is the unit vector that distribute uniformly on the d -dimensional hypersphere. τ , σ and M are all constants.

Finally, the training target in the original PFGM paper is slightly different from Eq. (13.8), it is approximately equivalent to

$$\left\| s_{\theta}(\mathbf{x}, t) + \operatorname{Normalise} \left(\mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}(\mathbf{x}_0)} \left[\frac{(\mathbf{x} - \mathbf{x}_0, t)}{(\|\mathbf{x} - \mathbf{x}_0\|^2 + t^2)^{(d+1)/2}} \right] \right) \right\|^2 \quad (13.10)$$

In the real implementation, we can only sample a finite number of x_0 evaluate the expectation value inside the bracket. As such, this target function contains some bias. Although a biased target function is not necessarily worse than an unbiased one. Why such a target function was chosen in the original paper was unclear. Our best guess is because the biased estimation has introduced the normalisation operations in the vectors, which has made the training process numerically more stable. But this will also increase the batch size to achieve better accuracy, thus the training cost.

Overall, in this section, we have introduced an ODE-based diffusion model that is inspired by the law of the gravitational forces. It has overcome the reliance on the Gaussian distributions on the diffusion models discussed before, and is a brand new theoretical framework that builds ODE-based diffusion model from the field theory, which is worth to be studied in details.

14. The General Procedure of Constructing an ODE (Part I)

Continuing from our discussions in the previous section, in which we have introduced an ODE-based diffusion generative model that was inspired by the gravitational field with a clear geometric meaning. Some readers may ask, it seems like the gravitational field may not be the only choice, is it possible to use other type of forces to construct the diffusion model under the same physical picture? On the other hand, even though the model is physically very trivial, there is a lack of mathematical proof that we could indeed learn the real sample distribution from it.

Here, we aim to provide a mathematically more accurate answer on what type of force field is suitable for constructing an ODE-based diffusion generative model.

14.1. The basic conclusions from the previous sections

To answer this question, we need to use a conclusion that we have derived in Section 12, which finds the corresponding ODE for describing the change in the probability distributions from the starting ODE that describes the evolution of the samples.

Consider the first-order ODE on $x_t \in \mathbb{R}^d$ and $t \in [0, T]$:

$$\frac{dx_t}{dt} = f_t(x_t), \quad (14.1)$$

which describes an (invertible) transformation from x_0 to x_T . If x_0 was a random variable, then x_t will also be a random variable throughout the transformation. The corresponding changes in the probability distribution can be described by the following ODE:

$$\frac{d}{dt}p_t(x_t) = -\nabla_{x_t} \left(f_t(x_t)p_t(x_t) \right). \quad (14.2)$$

This equation could be derived by using the Jacobian and Taylor expansion as shown in Section 12, or as we have done in Section 6, we derived the full Fokker-Plank equation first, and then set $g_t = 0$. It should be pointed out that, Eq. (14.2) is the famous **continuity equation** in physics, which reflects many different conservation laws.

Returning to the diffusion model, our goal here, is to construct a transformation, that is capable of transforming a simple distribution into the target sample distribution. In principle, we could apply Eq. (14.2) to solve for a possible $f_t(x_t)$ with a given distribution $p_t(x_t)$, then using Eq. (14.1) to complete the generation process. However, noticing that we only have one equation [Eq. (14.2)] here but the target function $f_t(x_t)$ has d components, so this is an indeterminate equation. In principle, we can solve the equation by assigning an arbitrary complete set of $p_t(x_t)$ and solve for $f_t(x_t)$, not just the distributions at the boundaries corresponding to $t = 0$ and $t = T$.

As such, the construction of an ODE-based diffusion model is simply to solve an ODE with almost no constraints. Although this is conceptually true, it is particularly difficult to achieve in coding. Therefore, a more accurate way to rephrase the question is how can we solve for some practically more useful solutions from Eq. (14.2).

Original blog post see <https://www.spaces.ac.cn/archives/9370>

14.2. Simplifying the equation

Noticing that, Eq. (14.2) can be rewritten in the following form:

$$\underbrace{\left(\frac{\partial}{\partial t}, \nabla_{\mathbf{x}_t}\right)}_{\nabla_{(t, \mathbf{x}_t)}} \cdot \underbrace{\left(p_t(\mathbf{x}_t), f_t(\mathbf{x}_t)p_t(\mathbf{x}_t)\right)}_{\mathbf{u} \in \mathbb{R}^{d+1}} = 0. \quad (14.3)$$

As shown above, we can treat $\left(\frac{\partial}{\partial t}, \nabla_{\mathbf{x}_t}\right)$ as the gradient $\nabla_{(t, \mathbf{x}_t)}$ in the $(d+1)$ -dimensional space, and $\left(p_t(\mathbf{x}_t), f_t(\mathbf{x}_t)p_t(\mathbf{x}_t)\right)$ can be written as a $(d+1)$ -dimensional vector, such that Eq. (14.2) can be casted into a simple divergence equation:

$$\nabla_{(t, \mathbf{x}_t)} \mathbf{u}(t, \mathbf{x}_t) = 0. \quad (14.4)$$

Under this form, we have

$$\frac{d\mathbf{x}_t}{dt} = f_t(\mathbf{x}_t) = \frac{\mathbf{u}_{>1}(t, \mathbf{x}_t)}{\mathbf{u}_1(t, \mathbf{x}_t)}, \quad (14.5)$$

in which \mathbf{u}_1 and $\mathbf{u}_{>1}$ represent the first and last d components of $\mathbf{u}(t, \mathbf{x}_t)$. of course, we cannot forget the boundary conditions:

$$\begin{cases} \mathbf{u}_1(0, \mathbf{x}_0) = p_0(\mathbf{x}_0) & \text{(initial condition)} \\ \int \mathbf{u}_1(t, \mathbf{x}_t) dt = 1 & \text{(integral condition)} \end{cases} \quad (14.6)$$

in which $p_0(\mathbf{x}_0)$ is the target sample distribution. For the distribution at $t = T$, we only require it to be a distribution that is easy to sample from, with no other constraint to be imposed at this stage.

14.3. The Green's function

With the above transformation, we can treat $\mathbf{u}(t, \mathbf{x}_t)$ as a vector field in the $(d+1)$ -dimensional space, and Eq. (14.5) determines the trajectory of a point mass that is moving along the field direction. this then becomes consistent with the physical picture established in the last section based on the gravitational forces.

In order to solve for a general solution for $\mathbf{u}(t, \mathbf{x}_t)$, we can apply the mathematical tool of the Green's function. Firstly, let us try to solve the following problem

$$\begin{cases} \nabla_{(t, \mathbf{x}_t)} \cdot G(t, 0; \mathbf{x}_t, \mathbf{x}_0) = 0 \\ G_1(0, 0; \mathbf{x}_t, \mathbf{x}_0) = \delta(\mathbf{x}_t - \mathbf{x}_0) \\ \int G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0) d\mathbf{x}_t = 1. \end{cases} \quad (14.7)$$

It is not difficult to prove that, if the above equation holds, then we have

$$\mathbf{u}(t, \mathbf{x}_t) = \int G(t, 0; \mathbf{x}_t, \mathbf{x}_0) p_0(\mathbf{x}_0) d\mathbf{x}_0 = \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0)}[G(t, 0; \mathbf{x}_t, \mathbf{x}_0)], \quad (14.8)$$

which is the solution to Eq. (14.4) that satisfies the corresponding constraints. In this case, we have successfully expressed $\mathbf{u}(t, \mathbf{x}_t)$ as some expectation value over the training samples, which will help the model training. It is not difficult to see that the Green's function $G(t, 0; \mathbf{x}_t, \mathbf{x}_0)$ is indeed the conditional probability $p(\mathbf{x}_t | \mathbf{x}_0)$ in the diffusion model.

As a matter of face, the Green's function $G(t, 0; \mathbf{x}_t, \mathbf{x}_0)$ that is defined in Eq. (14.7) is not the green's function that is usually defined. In the traditional definition, the Green's function corresponds to the solution at the point source of the field, whereas here, we have put the source at the boundary points. Even so, the Green's function $G(t, 0; \mathbf{x}_t, \mathbf{x}_0)$ defined here still possesses the properties of a normal Green's function, which is equivalent to a forcefield that is sourced from a point and is continuously distributed in the space.

14.4. The gravitational force

Now we will seek for some concrete solutions based on the above established mathematical framework. As we have mentioned before, both Eq. (14.4) and Eq. (14.7) are undetermined

equations that give one equation to solve for $d + 1$ unknowns. Theoretically, there exist infinitely many solutions of different forms. Therefore, to actually solve these equations, we must introduce some additional assumptions in order to arrive at a more concrete solution. The first assumption is the isotropic assumption, which corresponds to the results given in Section 13.

14.4.1. An ansatz

We note that the isotropic property is referring to the state (t, \mathbf{x}_t) lying in the $(d + 1)$ -dimensional space being isotropic. It means that $G(t, 0; \mathbf{x}_t, \mathbf{x}_0)$ points to the origin at $(0, \mathbf{x}_0)$, and the modulus is only depending on $R = \sqrt{t^2 + \|\mathbf{x}_t - \mathbf{x}_0\|^2}$. As such, we can let

$$G(t, 0; \mathbf{x}_t, \mathbf{x}_0) = \varphi(R)(t, \mathbf{x}_t - \mathbf{x}_0). \quad (14.9)$$

From which we have

The first boundary condition as in Eq. (14.7).

$$\begin{aligned} 0 &= \nabla_{(t, \mathbf{x}_t)} \cdot G(t, 0; \mathbf{x}_t, \mathbf{x}_0) \\ &= [\nabla_{(t, \mathbf{x}_t)} \varphi(R)] \cdot (t, \mathbf{x}_t - \mathbf{x}_0) + \varphi(R) [\nabla_{(t, \mathbf{x}_t)} (t, \mathbf{x}_t - \mathbf{x}_0)] \quad (\text{product rule}) \\ &= \varphi'(R) \cdot \frac{(t, \mathbf{x}_t - \mathbf{x}_0)}{R} \cdot (t, \mathbf{x}_t - \mathbf{x}_0) + (d + 1)\varphi(R) \quad (\text{chain rule}) \\ &= \varphi'(R)R + (d + 1)\varphi(R) \quad (\text{definition of } R) \\ &= \frac{[\varphi(R)R^{d+1}]'}{R^d}. \quad (\text{reversing the product rule}) \end{aligned} \quad (14.10)$$

This implies that we must have $[\varphi(R)R^{d+1}]' = 0$ or $\varphi(R)R^{d+1} = C$, *i.e.* $\varphi(R) = CR^{-(d+1)}$. As such, an ansatz for Eq. (14.7) is

$$G(t, 0; \mathbf{x}_t, \mathbf{x}_0) = C \times \frac{(t, \mathbf{x}_t - \mathbf{x}_0)}{(t^2 + \|\mathbf{x}_t - \mathbf{x}_0\|^2)^{(d+1)/2}}. \quad (14.11)$$

14.4.2. Imposing the constraints

it can be seen from the above derivation, that under the assumption of isotropy, the gravitational field becomes the only possible ansatz. To further prove that this is a feasible solution, we must also check that it satisfies the other constraints, in which the key one is

$$\int G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0) d\mathbf{x}_t = C \times \int \frac{t}{(t^2 + \|\mathbf{x}_t - \mathbf{x}_0\|^2)^{(d+1)/2}} d\mathbf{x}_t = 1. \quad (14.12)$$

To prove this, we only need to show the integral is independent from \mathbf{x}_0 and t , and provided that a suitable normalisation constant C is chosen, we can always make the value of the integral to be unity. For any $t > 0$, we can make the substitution of $z = (\mathbf{x}_t - \mathbf{x}_0)/t$. Since $\mathbf{x}_t \in \mathbb{R}^d$, we also have $z \in \mathbb{R}^d$. With this, Eq. (14.12) becomes

$$\int G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0) d\mathbf{x}_t = C \times \int \frac{1}{(1 + \|z\|^2)^{(d+1)/2}} dz, \quad (14.13)$$

which we can see now the integral is completely independent from \mathbf{x}_0 and t , so we just need to choose a suitable normalisation constant C to make the integral to become unity. In the following discussions, we assume that such C has already been chosen.

As for the initial values, we need to prove the limit $\lim_{t \rightarrow 0^+} G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0) = \delta(\mathbf{x}_t - \mathbf{x}_0)$, which we only need to show this based on the definition of the Dirac function: (1) When $\mathbf{x}_t \neq \mathbf{x}_0$, the limit is obviously zero. (2) When $\mathbf{x}_t = \mathbf{x}_0$, the limit is obviously 1. (3) As we have already shown, the integral of $G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0)$ over \mathbf{x}_t is always unity. These three points are the basic properties, or simply put, the definition of the Dirac function, so we have make sure the function has the correct initial values.

14.4.3. Result analysis

According to Eq. (14.8), we now have

$$\mathbf{u}(t, \mathbf{x}_t) = C \times \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0)} \left[\frac{(t, \mathbf{x}_t - \mathbf{x}_0)}{(t^2 + \|\mathbf{x}_t - \mathbf{x}_0\|^2)^{(d+1)/2}} \right], \quad (14.14)$$

which we can further apply the relationship $\mathbb{E}_x[x] = \underset{\mu}{\operatorname{argmin}} \mathbb{E}_x[\|x - \mu\|^2]$ to construct a training target that is similar to score-matching. which will not be repeated here.

It has been mentioned previously, that $G_1(t, 0; x_t, x_0)$ is equivalent to $p_t(x_t|x_0)$, which we have now known the exact form

$$p_t(x_t|x_0) \propto \frac{t}{(t^2 + \|x_t - x_0\|^2)^{(d+1)/2}}, \quad (14.15)$$

when T is sufficiently large, the influence of x_0 upon $p_t(x_t|x_0)$ becomes negligible, in which case $p_t(x_t|x_0)$ becomes the prior distribution that is independent from x_0 :

$$p_{\text{prior}}(x_T) \propto \frac{T}{(T^2 + \|x_T\|^2)^{(d+1)/2}}. \quad (14.16)$$

This is the same prior distribution as derived in Section 13 but it comes more naturally under the current framework. Since now we have $p_t(x_t|x_0)$, in principle, we can also achieve the sampling of $x_t \sim p(x_t|x_0)$. From Eq. (14.13), we know that if we make the substitution of $z = (x_t - x_0)/t$, then we have

$$p(z) \propto \frac{1}{(1 + \|z\|^2)^{(d+1)/2}}, \quad (14.17)$$

From which we can first sample from $p(z)$ and then obtain x_t from $x_t = x_0 + zt$. As for the actual sampling of $p(z)$, since it only depends on the value of the modulus, we could first sample the modulus using the inverse cumulative function and then randomly sample an orientation to finish the entire sampling process, which is the same as sampling from the prior distribution. Nevertheless, when considering the remaining problems, we have discovered the following surprise!

14.4.4. Revisiting the problem

In Section 13, we have mentioned that the sampling procedure provided in the original paper was of the following form,

$$x_t = x_0 + \|\varepsilon_x\|(1 + \tau)^m \mathbf{u}, \quad t = \|\varepsilon_t\|(1 + \tau)^m, \quad (14.18)$$

in which $(\varepsilon_x, \varepsilon_t) \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_{(d+1) \times (d+1)})$ and $m \sim \mathcal{U}[0, M]$. \mathbf{u} is the unit vector that distribute uniformly on the d -dimensional hypersphere. τ , σ and M are all constants.

Initially, it was thought that such a sampling procedure was purely a subjective design from the original authors, with not much reasoning. Subsequently, we found a surprising coincident, that this sampling procedure is a way to achieve the sampling from the distribution given in Eq. (14.17). We shall prove this in the following discussions. First of all, we substitute the second half of Eq. (14.18) into the first half, which we get:

$$x_t = x_0 + t \frac{\|\varepsilon_x\|}{|\varepsilon_t|} \cdot \mathbf{u}, \quad (14.19)$$

that is of exactly the same form of $x_t = x_0 + tz$ derived above, and \mathbf{u} is an isotropic random variable. In this case, the question becomes whether $\frac{\|\varepsilon_x\|}{|\varepsilon_t|}$ follows the same probability distribution as $\|z\|$. The answer is yes. Paying attention to the face that when converting the probability distribution from the Cartesian to spherical coordinates, we need to multiply an extra term that is proportional to the radius to the power of $d - 1$, as such, according to Eq. (14.17), we have

$$p(\|z\|) \propto \frac{\|z\|^{d-1}}{(1 + \|z\|^2)^{(d+1)/2}}. \quad (14.20)$$

On the other hand, according to the choice of $(\varepsilon_x, \varepsilon_t) \sim \mathcal{N}(0, \mathbf{I}_{(d+1) \times (d+1)})$ (Here, because we are only interested in the ratio between two probability distributions, the standard deviations cancel out each other, thus, for simplicity, we can let $\sigma = 1$), so we have the following

$$p(\|\varepsilon_x\|) \propto \|\varepsilon_x\|^{d-1} e^{-\|\varepsilon_x\|^2/2}, \quad p(|\varepsilon_t|) \propto e^{-|\varepsilon_t|^2/2}. \quad (14.21)$$

Let $r = \|\varepsilon_x\|/|\varepsilon_t|$, such that $\|\varepsilon_x\| = r|\varepsilon_t|$, then take the usage of equality in probability, we have

$$\begin{aligned}
p(r)dr &= \mathbb{E}_{|\varepsilon_t| \sim p(|\varepsilon_t|)} [p(\|\varepsilon_x\| = r|\varepsilon_t|)d(r|\varepsilon_t|)] \\
&\propto \mathbb{E}_{|\varepsilon_t| \sim p(|\varepsilon_t|)} \left[r^{d-1} |\varepsilon_t|^d e^{-r^2 |\varepsilon_t|^2 / 2} dr \right] \\
&\propto \int_0^\infty r^{d-1} |\varepsilon_t|^d e^{-r^2 |\varepsilon_t|^2 / 2} e^{-|\varepsilon_t|^2 / 2} d|\varepsilon_t| dr \\
&\propto \int_0^\infty r^{d-1} |\varepsilon_t|^d e^{-(r^2+1)|\varepsilon_t|^2 / 2} d|\varepsilon_t| dr \\
&= \frac{r^{d-1}}{(1+r^2)^{(d+1)/2}} \int_0^\infty s^d e^{-s^2/2} ds dr \quad (s = |\varepsilon_t| \sqrt{r^2+1}) \\
&\propto \frac{r^{d-1}}{(1+r^2)^{(d+1)/2}}
\end{aligned} \tag{14.22}$$

This shows that $p(r)$ is fully equivalent to Eq. (14.20). As such, using $(\|\varepsilon_x\|/|\varepsilon_t|)\mathbf{u}$ has indeed provided an effective sampling strategy for \mathbf{z} , which is numerically much easier than using the inverse cumulative function. But this was not explicitly mentioned in the original paper.

14.5. Space-time separation

What we have done so far is to solve for (t, \mathbf{x}_t) under the assumption of isotropy in the $(d+1)$ -dimensional space. To some extent, this counts for the simplest possible solution, although some readers might find it difficult to accept, as this gravitational field-based generative model is seemingly more complicated mathematically. However, when solving the mathematical physics equations, it is very common that we try solving the problem first by starting from a simple isotropic ansatz.

Nevertheless, it is indeed difficult to rationalise the isotropic nature of (t, \mathbf{x}_t) when taking the space-time variables together as a whole. We are more comfortable in dealing with isotropy in space by separating out the time variable. Here, we shall seek for the solution under the separation of time and space.

14.5.1. An ansatz

What the above discussions imply here, are that the isotropy should now apply only to the spatial variable \mathbf{x}_t in the d -dimensional space, and we will separate $G(t, 0; \mathbf{x}_t, \mathbf{x}_0)$ into two parts $G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0)$ and $G_{>1}(t, 0; \mathbf{x}_t, \mathbf{x}_0)$ to help with our understanding. Here, $G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0)$ is only a scalar, in which case, isotropy implies that it only depends on the modulus $r = \|\mathbf{x}_t - \mathbf{x}_0\|$, which we shall denote it as $\phi_t(r)$. Whereas $G_{>1}(t, 0; \mathbf{x}_t, \mathbf{x}_0)$ is a d -dimensional vector, and isotropy implies that it points to \mathbf{x}_0 from all directions, and the modulus is only dependent upon $r = \|\mathbf{x}_t - \mathbf{x}_0\|$. As such, we can let

$$G_{>1}(t, 0; \mathbf{x}_t, \mathbf{x}_0) = \varphi_t(r)(\mathbf{x}_t - \mathbf{x}_0). \tag{14.23}$$

With this, we have [from Eq. (14.4)]:

$$\begin{aligned}
0 &= \frac{\partial}{\partial t} \phi_t(r) + \nabla_{\mathbf{x}_t} \left(\varphi_t(r)(\mathbf{x}_t - \mathbf{x}_0) \right) \\
&= \frac{\partial}{\partial t} \phi_t(r) + r \frac{\partial}{\partial r} \varphi_t(r) + d\varphi_t(r) \\
&= \frac{\partial}{\partial t} \phi_t(r) + \frac{1}{r^{d-1}} \frac{\partial}{\partial r} \left(\varphi_t(r) r^d \right).
\end{aligned} \tag{14.24}$$

Here, we have two undetermined functions $\phi_t(r)$ and $\varphi_t(r)$ but only one differential equation, which makes it easier to solve. Since the conditional constraint is applied to $G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0)$, that is, to the function $\phi_t(r)$ not $\varphi_t(r)$, in this case, we could, for simplicity, solve for $\varphi_t(r)$ with a given function of $\phi_t(r)$ that satisfies the constraints. With this, the solution for Eq. (14.24) is

$$\varphi_t(r) = -\frac{1}{r^d} \int \frac{\partial}{\partial t} \phi_t(r) r^{d-1} dr = -\frac{1}{r^d} \frac{\partial}{\partial t} \int \phi_t(r) r^{d-1} dr. \tag{14.25}$$

14.5.2. Gaussian diffusion

Here, we shall further prove that, the common ODE-based diffusion model that builds upon the Gaussian distribution is a special case for Eq. (14.25). For the assumption that is based on the Gaussian distribution, we have the following transition probability:

$$G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0) = p_t(\mathbf{x}_t | \mathbf{x}_0) = \frac{1}{(2\pi\sigma_t^2)^{d/2}} e^{-\|\mathbf{x}_t - \mathbf{x}_0\|^2 / 2\sigma_t^2}, \quad (14.26)$$

i.e. $\phi_t(r) = \frac{1}{(2\pi\sigma_t^2)^{d/2}} e^{-r^2/2\sigma_t^2}$, in which σ_t is a monotonically increasing function of the time variable t . It satisfies the boundary condition that $\sigma_0 = 0$ and σ_T being sufficiently large. $\sigma_0 = 0$ is set to satisfy the initial value condition, and σ_T is to be made sufficiently large in order to make the prior distribution becomes independent from the sample distribution. As for the unity integral constraint, this is naturally satisfied by the Gaussian distribution.

Substituting Eq. (14.26) into Eq. (14.25) we have the following solution

$$\varphi_t(r) = \frac{\dot{\sigma}_t}{(2\pi\sigma_t^2)^{d/2}\sigma_t} e^{-r^2/2\sigma_t^2} = \frac{\dot{\sigma}_t}{\sigma_t} \phi_t(r), \quad (14.27)$$

in which the integration over r has taken the usage of the incomplete Gamma function that the author has derived it with the help of Mathematica. With this solution, we have

$$\begin{aligned} \mathbf{u}_1(t, 0; \mathbf{x}_t, \mathbf{x}_0) &= \int p_t(\mathbf{x}_t | \mathbf{x}_0) p(\mathbf{x}_0) d\mathbf{x}_0 = p_t(\mathbf{x}_t) \quad [\text{from Eq. (14.26)}] \\ \mathbf{u}_{>1}(t, 0; \mathbf{x}_t, \mathbf{x}_0) &= \int \frac{\dot{\sigma}_t}{\sigma_t} (\mathbf{x}_t - \mathbf{x}_0) p_t(\mathbf{x}_t | \mathbf{x}_0) p_0(\mathbf{x}_0) d\mathbf{x}_0 \\ &= -\dot{\sigma}_t \sigma_t \int \nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t | \mathbf{x}_0) p_0(\mathbf{x}_0) d\mathbf{x}_0 \\ &= -\dot{\sigma}_t \sigma_t \nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t), \end{aligned} \quad (14.28)$$

in which, based on Eq. (14.5), we have

$$f_t(\mathbf{x}_t) = \frac{\mathbf{u}_{>1}(t, \mathbf{x}_t)}{\mathbf{u}_1(t, \mathbf{x}_t)} = -\dot{\sigma}_t \sigma_t \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t). \quad (14.29)$$

This result is identical to those derived in Section 12

14.6. Reverse engineering

Solving for $\varphi_t(r)$ with a given $\phi_t(r)$ is theoretically simple, but there are two major challenges: (1) It is difficult to construct a function $\phi_t(r)$ that is simultaneously satisfying both the initial value and the integral conditions, and (2) One cannot guarantee that the integration over r will always have a simple analytical form. As such, we need to come up with a method of reverse engineering to construct $\phi_t(r)$.

We know that $\phi_t(r)$ is the probability density under the Cartesian coordinates, which, upon transformation into the spherical harmonics, one must multiply it with $C_d r^{d-1}$, in which C_d is a d -dependent constant. Based on Eq. (14.5) that the final result is a ratio that will not be affected by this constant, so we can ignore it for simplicity. This leads straight to the integrand in Eq. (14.25). As such the integral in Eq. (14.25)

$$\int \phi_t(r) r^{d-1} dr \quad (14.30)$$

is a cumulative probability function. It is not always easy to calculate the cumulative probability function through integration, but going the other way around via differentiation is very easy. As such, we could build the cumulative probability function, and then solve for the corresponding $\phi_t(r)$ and $\varphi_t(r)$, bypassing the challenge of integration.

To build the cumulative probability function $\psi_t(r)$, it must satisfy the following conditions: (1) $\psi_t(0) = 0$ and $\psi_t(\infty) = 1$. (2) $\psi(r)$ is a monotonically increasing function of r . (3) $\forall r > 0, \lim_{t \rightarrow 0^+} \psi_t(r) = 1$. For people who are familiar with the activation function, it should not be difficult to construct such functions, it is basically the smoothed approximation to the step function such as $\tanh(r/t)$ and $1 - e^{-r/t}$. With the cumulative function $\psi_t(r)$, and based

on Eq. (14.25), we have

$$\phi_t(r) = \frac{1}{r^{d-1}} \frac{\partial}{\partial r} \psi_t(r), \quad \varphi_t(r) = -\frac{1}{r^d} \frac{\partial}{\partial t} (\phi_t(r) + \lambda_t), \quad (14.31)$$

in which λ_t is an arbitrary function of t , which can be set as zero in the usual situation. Certainly, all of the isotropic solutions are fundamentally equivalent to each other, including the results derived for the gravitational diffusions. They all could be included in the equations above. We can also derive different outcomes using coordinate transformation, as the cumulative function $\psi_t(r)$ is only a function of a scalar variable, and we can transform the cumulative probability function for different distributions, providing they are all well-behaved monotonically increasing functions.

14.7. Summary

To summarise, here, we have constructed a general framework of building ODE-based diffusion model. Theoretically speaking, the framework encompassed all ODE-based diffusion models, from which we can construct any type of new and funky diffusion models. For examples, the current derivations are all based on the isotropic assumption, which we could change $\varphi(R)$ into a more generic form of $\varphi(t; \mathbf{x}_t, \mathbf{x}_0)$ and solve the partial differential equation using the method of characteristic, to obtain a family of new models. Overall, this is really a factory of ODE-based diffusion models.

Some readers may ask, what is the purpose of deriving so many variants for diffusion models, if we just want a diffusion generative model that is practically workable. As a matter of fact, we want to discover and understand the fundamental principles behind the constructions of the diffusion generative models, such that we may discover other diffusion models that are more efficient. This is a never-ending process of perfectionism.

The numerical experimental results from the previously discussed gravitational diffusion model, as an ODE-based one, already led to a slightly better result than the conventional Gaussian diffusion model. This shows that, even they are all based on the isotropy assumption, the practical performances of these mathematically equivalent models can still be different from each other. Therefore, answering the question of what is a better design for the diffusion model would be a very meaningful research direction in the future.

15. The General Procedure of Constructing an ODE (Part II)

When writing the previous section, the author believed that we had already found a general procedure for constructing an ODE-based diffusion generative model. However, a reader has later provided a more efficient and straightforward alternative, which the author was truly amazed by the inspiration behind the model. After some discussions and thoughts, the author discovered that the fundamental idea behind this new approach is based on the **method of characteristics** for solving the first-order partial differential equations, in which we first construct a particular vector field to enforce the initial condition, and then the differential equation will be solved to match the boundary conditions that are to be satisfied by the end points. As such, the conditions for both the starting and end points that need to be matched are satisfied simultaneously. This is a really smart approach, which the author will summarise his thoughts in this section, as a sequel to the previous section.

Original blog post see <https://www.spaces.ac.cn/archives/9379>

15.1. Revisiting the previous results

We shall first briefly revise the results from the previous section. Suppose a random variable $\mathbf{x}_0 \in \mathbb{R}^d$ will transform itself continuously in space into \mathbf{x}_T by following an ODE of the form,

$$\frac{d\mathbf{x}_t}{dt} = \mathbf{f}_t(\mathbf{x}_t), \quad (15.1)$$

then the corresponding probability distribution $p_t(\mathbf{x}_t)$ at any given time t shall follow the following continuity equation:

$$\frac{\partial}{\partial t} p_t(\mathbf{x}_t) = -\nabla_{\mathbf{x}_t} \cdot (\mathbf{f}_t(\mathbf{x}_t) p_t(\mathbf{x}_t)). \quad (15.2)$$

By letting $\mathbf{u}(t, \mathbf{x}_t) = (p_t(\mathbf{x}_t), f_t(\mathbf{x}_t)p_t(\mathbf{x}_t)) \in \mathbb{R}^{d+1}$, then the above equation can be simplified into

$$\begin{cases} \nabla_{(t, \mathbf{x}_t)} \cdot \mathbf{u}(t, \mathbf{x}_t) = 0 \\ \mathbf{u}_1(0, \mathbf{x}_0) = p_0(\mathbf{x}_0), \quad \int \mathbf{u}_1(t, \mathbf{x}_t) d\mathbf{x}_t = 1. \end{cases} \quad (15.3)$$

To solve this equation, we applied the technique of Green's function, in which we first solve

$$\begin{cases} \nabla_{(t, \mathbf{x}_t)} \cdot G(t, 0; \mathbf{x}_t, \mathbf{x}_0) = 0 \\ G_1(0, 0; \mathbf{x}_t, \mathbf{x}_0) = \delta(\mathbf{x}_t - \mathbf{x}_0), \quad \int G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0) = 1. \end{cases} \quad (15.4)$$

From here, we have the following solution

$$\mathbf{u}(t, \mathbf{x}_t) = \int G(t, 0; \mathbf{x}_t, \mathbf{x}_0) p_0(\mathbf{x}_0) d\mathbf{x}_0 = \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0)} [G(t, 0; \mathbf{x}_t, \mathbf{x}_0)], \quad (15.5)$$

which satisfies our boundary conditions.

15.2. Geometric interpretations

The idea behind Green's function is simple. First of all, let us forget about synthesising a complex set of data, by just assuming we only have one data point \mathbf{x}_0 and try to figure out how we should synthesise this single data. This may look simple, in which we can simply take $\mathbf{x}_T \times 0 + \mathbf{x}_0$. However, this is an oversimplification, because we want a slowly-varying and continuous transformation, such that for any arbitrary \mathbf{x}_t at $t < T$, it will follow a smooth trajectory to \mathbf{x}_0 at $t = 0$.

Since we are only interested in building a diffusion generative model, in principle, we are not interested in the actual shape of the trajectory, as long as they all pass through \mathbf{x}_0 . In this case, we could choose any clusters of trajectories that we prefer, which will be collectively denoted as

$$\varphi_t(\mathbf{x}_t | \mathbf{x}_0) = \mathbf{x}_T. \quad (15.6)$$

Let us emphasise again that this represents a cluster of trajectories that take \mathbf{x}_0 as the starting point and the \mathbf{x}_T as the end point. The independent and dependent variables in each trajectory are t and \mathbf{x}_t , respectively. The starting point \mathbf{x}_0 is always fixed, whereas the final endpoint \mathbf{x}_T is variable. There is also no restriction on the shape of the trajectory, we could choose a straight line, a parabola or anything else.

Now we take the total derivatives to the both sides of Eq. (15.6). As \mathbf{x}_T can be changed arbitrarily, it is equivalent to the integration constant in the differential equation, and its derivative is simply 0, as such, we have,

$$\frac{\partial \varphi_t(\mathbf{x}_t | \mathbf{x}_0)}{\partial \mathbf{x}_t} \frac{d\mathbf{x}_t}{dt} + \frac{\partial \varphi_t(\mathbf{x}_t | \mathbf{x}_0)}{\partial t} = \mathbf{0}, \quad (15.7)$$

which gives,

$$\frac{d\mathbf{x}_t}{dt} = - \left(\frac{\partial \varphi_t(\mathbf{x}_t | \mathbf{x}_0)}{\partial \mathbf{x}_t} \right)^{-1} \frac{\partial \varphi_t(\mathbf{x}_t | \mathbf{x}_0)}{\partial t}.$$

Comparing this to Eq. (15.1), we obtain

$$f_t(\mathbf{x}_t | \mathbf{x}_0) = - \left(\frac{\partial \varphi_t(\mathbf{x}_t | \mathbf{x}_0)}{\partial \mathbf{x}_t} \right)^{-1} \frac{\partial \varphi_t(\mathbf{x}_t | \mathbf{x}_0)}{\partial t}. \quad (15.8)$$

Here, we have changed the original function of $f_t(\mathbf{x}_t)$ into $f_t(\mathbf{x}_t | \mathbf{x}_0)$ to signify that all the trajectories have a common starting point of \mathbf{x}_0 . This implies that, the trajectories from the force field built from this ODE will always pass through \mathbf{x}_0 . This makes sure the initial value condition for the Green's function is always satisfied.

Not quite sure about this. This probably comes from the total derivative of function $f(\mathbf{x}(t), t)$ with respect to t is given by $\frac{\partial f}{\partial \mathbf{x}} \frac{d\mathbf{x}}{dt} + \frac{\partial f}{\partial t}$.

15.3. The method of characteristics

Once we have the condition for the initial value, we could ask further on how to satisfy the condition for the final value, which we want \mathbf{x}_T to follow a simple distribution that is independent from \mathbf{x}_0 . The limitation in the derivations shown in the previous section is that we cannot guarantee the final value follows a simple distribution, and must be analysed case-by-case in a post-hoc fashion. Here, our design principle is to satisfy the condition that is imposed on the initial values directly through the design of $\varphi_t(\mathbf{x}_t|\mathbf{x}_0)$, which leaves us with plenty of rooms to satisfy the conditions imposed on the final value. Once these conditions are satisfied, the integration constraint will be naturally satisfied provided that the continuity equation [Eq. (15.2)] is also satisfied.

Mathematically speaking, what we want to do is to solve Eq. (15.2) with a given $f_t(\mathbf{x}_t|\mathbf{x}_0)$ and $P_T(\mathbf{x}_T)$. This is a first-order partial differential equation which may be solved using the method of characteristics. Firstly, we rewrite Eq. (15.2) into the following equivalent form:

$$\frac{\partial}{\partial t} p_t(\mathbf{x}_t|\mathbf{x}_0) + \nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t|\mathbf{x}_0) \cdot f_t(\mathbf{x}_t|\mathbf{x}_0) = -p_t(\mathbf{x}_t|\mathbf{x}_0) \nabla_{\mathbf{x}_t} f_t(\mathbf{x}_t|\mathbf{x}_0). \quad (15.9)$$

Similarly as before, we seek for its solution with a fixed starting point \mathbf{x}_0 , such that we have rewritten $p_t(\mathbf{x}_t)$ as $p_t(\mathbf{x}_t|\mathbf{x}_0)$ to signify this as the solution with starting point \mathbf{x}_0 .

The idea behind the method of characteristics is to first seek the solution of a partial differential equation along a specified trajectory, which transforms the partial differential equation into an ordinary differential equation, making it easier to solve. More specifically, we assume that \mathbf{x}_t is a function of t , and solve the partial differential equation along the trajectory that is defined by Eq. (15.1). As Eq. (15.1) holds, we could replace $f_t(\mathbf{x}_t|\mathbf{x}_0)$ in Eq. (15.9) by $\frac{d\mathbf{x}_t}{dt}$, such that the L.H.S. of Eq. (15.9) becomes the total derivative of $P_t(\mathbf{x}_t|\mathbf{x}_0)$, which gives us

$$\frac{d}{dt} p_t(\mathbf{x}_t|\mathbf{x}_0) = -p_t(\mathbf{x}_t|\mathbf{x}_0) \nabla_{\mathbf{x}_t} f_t(\mathbf{x}_t|\mathbf{x}_0). \quad (15.10)$$

Noted that, at this point, all \mathbf{x}_t should be replaced by a function with respect to t , which, in principle, can be solved from the equation for the trajectory [Eq. (15.6)]. With this substitution, both p and f in the above equation become purely the functions of t , such that the above equation is just linear ordinary differential equation of p , which should have the following solution:

$$p_t(\mathbf{x}_t|\mathbf{x}_0) = C \exp \left(\int_t^T \nabla_{\mathbf{x}_s} f_s(\mathbf{x}_s|\mathbf{x}_0) ds \right). \quad (15.11)$$

Substituting the condition for the final value $p_T(\mathbf{x}_T)$, we have $C = p_T(\mathbf{x}_T)$, such that

$$p_t(\mathbf{x}_t|\mathbf{x}_0) = p_T(\mathbf{x}_T) \exp \left(\int_t^T \nabla_{\mathbf{x}_s} f_s(\mathbf{x}_s|\mathbf{x}_0) ds \right). \quad (15.12)$$

If we further substitute into the equation for the trajectory [Eq. (15.6)], we obtain a function that will only include \mathbf{x}_0 , \mathbf{x}_t and t . This is the final result for the Green's function that we are trying to solve, i.e. $G_1(t, 0; \mathbf{x}_t, \mathbf{x}_0)$. Correspondingly, $G_{>1}(t, 0; \mathbf{x}_t, \mathbf{x}_0) = p_t(\mathbf{x}_t|\mathbf{x}_0) f_t(\mathbf{x}_t|\mathbf{x}_0)$.

15.4. The training target

From the Green's function, we obtain the following:

$$\begin{aligned} u_1(t, \mathbf{x}_t) &= \int p_t(\mathbf{x}_t|\mathbf{x}_0) p_0(\mathbf{x}_0) d\mathbf{x}_0 = p_t(\mathbf{x}_t) \\ u_{>1}(t, \mathbf{x}_t) &= \int f_t(\mathbf{x}_t|\mathbf{x}_0) p_t(\mathbf{x}_t|\mathbf{x}_0) p_0(\mathbf{x}_0) d\mathbf{x}_0, \end{aligned} \quad (15.13)$$

which leads to

$$\begin{aligned} f_t(\mathbf{x}_t) &= \frac{u_{>1}(t, \mathbf{x}_t)}{u_1(t, \mathbf{x}_t)} \\ &= \int f_t(\mathbf{x}_t|\mathbf{x}_0) \frac{p_t(\mathbf{x}_t|\mathbf{x}_0) p_0(\mathbf{x}_0)}{p_t(\mathbf{x}_t)} d\mathbf{x}_0 \end{aligned}$$

$$\begin{aligned}
&= \int f_t(\mathbf{x}_t|\mathbf{x}_0)p_t(\mathbf{x}_0|\mathbf{x}_t)d\mathbf{x}_0 \\
&= \mathbb{E}_{\mathbf{x}_0 \sim p_t(\mathbf{x}_0|\mathbf{x}_t)}[f_t(\mathbf{x}_t|\mathbf{x}_0)].
\end{aligned} \tag{15.14}$$

Based on the discussion in Section 5, in which we built the training function based on score-matching, we can build our training target here similarly as

$$\begin{aligned}
&\mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} \left[\mathbb{E}_{\mathbf{x}_0 \sim p_t(\mathbf{x}_0|\mathbf{x}_t)} \left[\|v_\theta(\mathbf{x}_t, t) - f_t(\mathbf{x}_t|\mathbf{x}_0)\|_2^2 \right] \right] \\
&= \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_t \sim p_t(\mathbf{x}_t|\mathbf{x}_0)p(\mathbf{x}_0)} \left[\|v_\theta(\mathbf{x}_t, t) - f_t(\mathbf{x}_t|\mathbf{x}_0)\|_2^2 \right],
\end{aligned} \tag{15.15}$$

which has the same form as the **conditional flow matching** given in the article entitled “Flow Matching for Generative Modelling”¹. As we shall see later, the results from this article can all be derived using the method proposed above. After model training, sampling can be achieved simply by solving the differential equation $\frac{d\mathbf{x}_t}{dt} = v_\theta(\mathbf{x}_t, t)$. It can be seen from Eq. (15.15) that, the only requirement that we have here is to have a distribution $p_t(\mathbf{x}_t|\mathbf{x}_0)$ that is easy to sample from.

15.5. A few examples

Probably the abstract results from above are not easy to be understood. In what follows, we will provide some examples to help developing some intuitive understanding.

15.5.1. Linear trajectory

As the simplest example, we assume that \mathbf{x}_T transforms into \mathbf{x}_0 by following a linear trajectory. For simplicity, we can also set $T = 1$. Without loss of generality, we can write the equation for \mathbf{x}_t as

$$\mathbf{x}_t = (\mathbf{x}_1 - \mathbf{x}_0)t + \mathbf{x}_0 \quad \Rightarrow \quad \frac{\mathbf{x}_t - \mathbf{x}_0}{t} + \mathbf{x}_0 = \mathbf{x}_1. \tag{15.16}$$

According to Eq. (15.8), we have

$$f_t(\mathbf{x}_t|\mathbf{x}_0) = \frac{\mathbf{x}_t - \mathbf{x}_0}{t}, \tag{15.17}$$

and now $\nabla_{\mathbf{x}_t} \cdot f_t(\mathbf{x}_t|\mathbf{x}_0) = \frac{d}{t}$, such that, according to Eq. (15.12), we have

$$p_t(\mathbf{x}_t|\mathbf{x}_0) = \frac{p_1(\mathbf{x}_1)}{t^d}. \tag{15.18}$$

Now substituting in the expression for \mathbf{x}_1 from Eq. (15.16), we have

$$p_t(\mathbf{x}_t|\mathbf{x}_0) = \frac{p_1\left(\frac{\mathbf{x}_t - \mathbf{x}_0}{t} + \mathbf{x}_0\right)}{t^d}. \tag{15.19}$$

More specifically, if $p_1(\mathbf{x}_1)$ follows the standard normal distribution, then the above equation should follow $p_t(\mathbf{x}|\mathbf{x}_0) \sim \mathcal{N}(\mathbf{x}_t; (1-t)\mathbf{x}_0, t^2\mathbf{I})$, which corresponds to the result from the Gaussian diffusion model. The new result under this framework, *is to allow us to choose more general prior distributions*, such as the uniform distribution. Also, when introducing the score-matching in Eq. (15.15), we already mentioned that all we need to know is how to sample from the distribution of $p_t(\mathbf{x}_t|\mathbf{x}_0)$, and Eq. (15.19) tells us that this is easily achievable if the prior distribution is on that is easy to sample from. This is because

$$\mathbf{x}_t \sim p_t(\mathbf{x}_t|\mathbf{x}_0) \quad \Leftrightarrow \quad \mathbf{x}_t = (1+t)\mathbf{x}_0 + t\boldsymbol{\varepsilon} \quad \boldsymbol{\varepsilon} \sim p_1(\boldsymbol{\varepsilon}). \tag{15.20}$$

15.5.2. A further generalisation

The above result can be further generalised into the following:

$$\mathbf{x}_t = \boldsymbol{\mu}_t(\mathbf{x}_t) + \sigma_t \mathbf{x}_1 \quad \Rightarrow \quad \mathbf{x}_1 = \frac{\mathbf{x}_t - \boldsymbol{\mu}_t(\mathbf{x}_t)}{\sigma_t}. \tag{15.21}$$

Here, $\boldsymbol{\mu}_t(\mathbf{x}_t)$ is any arbitrary function that satisfies the boundary conditions of $\boldsymbol{\mu}_0(\mathbf{x}_0) = \mathbf{x}_0$ and $\boldsymbol{\mu}_1(\mathbf{x}_0) = 0$, that is defined over the domain of $\mathbb{R}^d \rightarrow \mathbb{R}^d$. σ_t is a monotonically increasing

function from $\sigma_0 = 0$ to $\sigma_1 = 1$. According to Eq. (15.8), we have the following

$$f_t(x_t|x_0) = \dot{\mu}_t(x_0) + \frac{\sigma_t}{\sigma_t} \left(x_t - \mu_t(x_0) \right). \quad (15.22)$$

This is equivalent to Eq. (15) in the article “Flow Matching for Generative Modelling”. Now $\nabla_{x_t} f_t(x_t|x_0) = d\dot{\sigma}_t/\sigma_t$, and from eq.15.12, we have

$$p_t(x_t|x_0) = p_1(x_1)/\sigma_t^d. \quad (15.23)$$

Substituting in the expression for x_1 , we have the following final result:

$$p_t(x_t|x_0) = \frac{p_1 \left(\frac{x_t - \mu_t(x_0)}{\sigma_t} \right)}{\sigma_t^d}. \quad (15.24)$$

This is the general result for linear ODE based diffusion model. It include prior distributions that can be either Gaussian or non-Gaussian.

15.5.3. Even more complicated?

In all the examples above, the transformations (or the trajectories) from x_0 to x_t are all established based on a simple linear interpolation function of t . A natural question to ask is, then, can we consider more complicated trajectories? This is theoretically possible. However, using more complex trajectories implies the need to embed more assumptions in the model, and it is usually difficult for us to prove whether all the target data follow these assumptions. Therefore, it is usually not necessary to consider more complex trajectories. In addition, for complex trajectories, it becomes more difficult to derive analytical solutions, making it impossible to proceed both theoretically and practically.

Most importantly, we are now only dealing with the trajectories for synthesising a single sample. It is not difficult to show that, even though this trajectory is a straight line, for synthesising multiple samples, the trajectories will lead to a more complex curves. In this case, if we already made the trajectory for synthesising a single sample to be unnecessarily complicated, then for multiple samples, the complexities behind the trajectories will be extremely high, making the model more likely to become unstable.

15.6. Summary

In summary, we have followed the context from the previous section, and further discussed the idea behind building an ODE-based diffusion model. This time, starting from a geometric intuition, we build specific vector fields to ensure the results follows the distributions of the initial values, and satisfy those for the final values by solving the differential equations. This allows us to solve for the Green's function that satisfies the conditions that both the initial and final values should follow. More specifically, it allows us to use any simple distributions as the prior distributions, thus overcoming the reliance on the Gaussian distributions in constructing the diffusion models.

16. Wasserstein Distance \leq Score-Matching

The Wasserstein distance, is a distance measure that quantifies the differences between two probability distributions defined based on the optimum transport theory. Many readers first came across the Wasserstein distance were probably from the birth of WGAN in 2017, which opened up a new branch for understanding GAN from the perspective of the optimum transport theory. It also made this theory start to play more important role in machine learning. For a long time, GAN has been the leader in the generative modelling, until the rise of the diffusion models in the past few years, even though GAN still remains as a compelling generative model.

The fundamental model architectures behind the diffusion models and GAN are very different, therefore, the research on these two models have been relatively independent from each other. However, the article entitled “Score-Based Generative Modelling Secretly Minimizes the Wasserstein Distance”²⁰ published at the end of 2022 had broken this barrier. It proved that the loss function for the diffusion model derived from the score-matching can be written as the upper limit of the Wasserstein distance. This implies that, when we minimise the loss function for the diffusion model, what we do is effectively the same as WGAN, in which the Wasserstein distance between two probability distributions is also minimised.

Original blog post see <https://www.spaces.ac.cn/archives/9467>

16.1. Result analysis

More specifically, the results presented in the original paper, was directly aligned with the SDE-based diffusion generative models discussed previously in Section 5. The core result is the following inequality:

$$\mathcal{W}_2[p_0, q_0] \leq \int_0^T g_t^2 I_t \left(\mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) - s_\theta(\mathbf{x}_t, t)\|_2^2 \right] \right)^{\frac{1}{2}} dt + I_T \mathcal{W}_2[p_T, q_T], \quad (16.1)$$

in which I_t is a non-negative function of t , the meaning of which will be discussed later.

How should we understand this inequality? Firstly, a diffusion model can be understood as a motion from $t = T$ to $t = 0$ that is described by an SDE. The terms p_T and q_T on the far right-hand-side are the random distribution for sampling at T , where p_T is usually the standard normal distribution. In real practices, we have $p_T = q_T$ such that $\mathcal{W}_2[p_T, q_T] = 0$. It was left in the equation to give the most theoretically generalised result. Secondly, p_0 on the left-hand-side, is obtained from solving the following SDE:

$$d\mathbf{x}_t = \left[f_t(\mathbf{x}_t) - g_t^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \right] dt + g_t d\mathbf{w}, \quad (16.2)$$

in the reversed order, starting from a random point \mathbf{x}_T that is sampled from p_T . This is the target sample distribution that we would like to generate. On the other hand q_0 is the distribution obtained from solving the SDE:

$$d\mathbf{x}_t = \left[f_t(\mathbf{x}_t) - g_t^2 s_\theta(\mathbf{x}_t, t) \right] dt + g_t d\mathbf{w}, \quad (16.3)$$

that starts from a random point sampled from q_T to reach the corresponding distribution at $t = 0$, in which $s_\theta(\mathbf{x}_t, t)$ is a neural network approximation to the score function $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$. As such, q_0 is the sample distribution that is generated from the diffusion generative model, and $\mathcal{W}_2[p_0, q_0]$ represents the Wasserstein distance between the true and generated sample distributions. Lastly, the remainder is the integral term, in which the key part is the integrand:

$$\mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) - s_\theta(\mathbf{x}_t, t)\|_2^2 \right], \quad (16.4)$$

which is nothing but the “score-matching” loss for the diffusion generative model.

the above result shows that, when we train the diffusion model with the score-matching loss, we have indirectly minimised the Wasserstein distance between the sample and the generated distributions. Differently from WGAN, the latter optimises $\mathcal{W}_1[p_0, q_0]$ distance, whereas here, the $\mathcal{W}_2[p_0, q_0]$ distance is optimised.

Note: Strictly speaking, Eq. (16.4) is not the loss function for the diffusion model, the latter should really be the “conditional score-matching”, and its relationship with score-matching is as following:

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) - s_\theta(\mathbf{x}_t, t)\|_2^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\left\| \mathbb{E}_{\mathbf{x}_0 \sim p_t(\mathbf{x}_0|\mathbf{x}_t)} \left[\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) \right] - s_\theta(\mathbf{x}_t, t) \right\|_2^2 \right] \\ &\leq \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\mathbf{x}_0|\mathbf{x}_t)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) - s_\theta(\mathbf{x}_t, t)\|_2^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0), \mathbf{x}_t \sim p_t(\mathbf{x}_t|\mathbf{x}_0)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) - s_\theta(\mathbf{x}_t, t)\|_2^2 \right], \end{aligned} \quad (16.5)$$

in which the last line corresponds to the conditional score-matching loss function for the diffusion generative model. The first equality holds because $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = \mathbb{E}_{\mathbf{x}_0 \sim p_t(\mathbf{x}_0|\mathbf{x}_t)} [\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0)]$. The second inequality is based on the Jansen’s inequality. Finally, the third equation is from the Bayesian formula. This implies that the conditional score-matching is the upper limit to score-matching, therefore, it should also be the upper limit to the Wasserstein distance.

From Eq. (16.1), we can also better rationalise why we need to drop the coefficients before the modulus in the target function for training the diffusion model. Here, the Wasserstein

distance is a well-defined metric for measuring the similarity between the two probability distribution functions, while $g_t^2 I_t$ in Eq. (16.1) is a monotonically increasing function of t . This means that we need to properly increase the score-matching loss when t is small. In Section 5, we had derived the final form of the score-matching loss being

$$\frac{1}{\bar{\beta}_t^2} \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}_0(\mathbf{x}_0)} \left[\left\| \varepsilon_\theta(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \varepsilon, t) - \varepsilon \right\|_2^2 \right]. \quad (16.6)$$

Ignoring the coefficient $1/\bar{\beta}_t^2$ is equivalent to multiplying the loss function with $\bar{\beta}_t^2$, which is also a monotonically increasing function of t . As such, for simplicity, we can consider ignoring the coefficient is an approach to make it closer to the Wasserstein distance between the two probability distributions.

16.2. Preparatory works

Although the original article has provided the proof for Eq. (16.1), it requires many prerequisite knowledges such as the continuity equation, gradient flow, and so on. Especially it used one theory without proof, which as hidden in the Chapter 8 of a monograph on gradient flow, or Chapter 5 of a monograph on the optimum transport theory, which was too difficult for the author to understand. After some trial, the author had finally arrived a (partial) proof of Eq. (16.1), which only requires the usage of the definition for the Wasserstein distance, fundamentals of differential equations and the Cauchy inequality. This is much simpler than the original proof. After some modifications and improvements, it is presented as follows.

Before proceeding to the proof, let us first make some preparations, which we shall summarise the basic concepts and conclusions that we need to apply subsequently. Firstly, the Wasserstein distance is defined as

$$\mathcal{W}_\rho[p, q] = \left(\inf_{\gamma \in \Pi[p, q]} \iint \gamma(\mathbf{x}, \mathbf{y}) \|\mathbf{x} - \mathbf{y}\|^\rho d\mathbf{x} d\mathbf{y} \right)^{\frac{1}{\rho}}, \quad (16.7)$$

in which $\Pi[p, q]$ is the probability density function for all the joint probability distributions in which p, q are the corresponding marginal distributions. It describes the concrete proposals for the transportation of masses between the two distributions.

Here, we will only consider the case in which $\rho = 2$, which is the only case that gives a simple derivation in what follows. Note that in the above definition for the Wasserstein distance, it includes the computation of the definitive lower bound (inf). This implies that, for any $\gamma \in \Pi[p, q]$ that we can find, we will always have

$$\mathcal{W}_2[p, q] \leq \left(\iint \gamma(\mathbf{x}, \mathbf{y}) \|\mathbf{x} - \mathbf{y}\|^2 d\mathbf{x} d\mathbf{y} \right)^{\frac{1}{2}}. \quad (16.8)$$

This is the key behind our proof here. To prove this bound, we shall be applying the Cauchy inequality:

$$\begin{aligned} \text{in vector form} \quad & \mathbf{x} \cdot \mathbf{y} \leq \|\mathbf{x}\| \cdot \|\mathbf{y}\|, \\ \text{in the form of expectation} \quad & \mathbb{E}_{\mathbf{x}}[f(\mathbf{x}) \cdot g(\mathbf{x})] \leq \left(\mathbb{E}_{\mathbf{x}}[f^2(\mathbf{x})] \right)^{\frac{1}{2}} \cdot \left(\mathbb{E}_{\mathbf{x}}[g^2(\mathbf{x})] \right)^{\frac{1}{2}}. \end{aligned} \quad (16.9)$$

Additionally, in the following proof, we will also assume that the function $g_t(\mathbf{x})$ satisfies the single-sided Lipschitz constraint defined as

$$\left(g_t(\mathbf{x}) - g_t(\mathbf{y}) \right) \cdot (\mathbf{x} - \mathbf{y}) \leq \mathcal{L}_t \|\mathbf{x} - \mathbf{y}\|^2, \quad (16.10)$$

which is a weaker version of the usual Lipschitz constraint, *i.e.* a function that satisfies the full Lipschitz constraint should also satisfy the single-sided Lipschitz constraint. Basically this will give a more appropriate definition of the Wasserstein distance and stable gradient during model training.

16.3. The initial trial

Eq. (16.1) is an overly generalised result, analysing which will not help us in rationalising and understanding the meaning behind. Therefore, we must first simplify the problem to see whether we can prove a weaker result. How should we simplify the problem? First of all, Eq. (16.1) considered the differences in the prior distributions (p_T and q_T), which we shall

Whilst the proof here seems not too challenge to understand, it does take a lot of inspirations to come up with such a proof!

consider them as being identical to each other. Secondly, the reverse process underpinning the derivation of Eq. (16.1) was based on an SDE [Eq. (16.2)], which we shall first consider a definitive transformation that is described by an ODE.

More specifically, let us consider a sample z that is drawn from the distribution $q(z)$ at $t = T$, which we shall treat it as our initial value, and then we shall let it evolve along trajectories defined by two different ODEs:

$$\frac{d\mathbf{x}_t}{dt} = f_t(\mathbf{x}_t) \quad \frac{d\mathbf{y}_t}{dt} = g_t(\mathbf{y}_t). \quad (16.11)$$

Assuming that, at any given time t , the distributions of \mathbf{x}_t and \mathbf{y}_t are p_t and q_t , respectively. We shall try to estimate an upper bound for $\mathcal{W}_2[p_0, q_0]$.

Since both \mathbf{x}_t and \mathbf{y}_t are evolved from the same initial value of z with respect to their own ODEs, therefore, they are both well-defined functions of z , which should be more precisely denoted as $\mathbf{x}_t(z)$ and $\mathbf{y}_t(z)$, but we have omitted z for brevity. This means that, for a given \mathbf{x} , $\mathbf{x}_t \leftrightarrow \mathbf{y}_t$ establishes a transport pathway between p_t and q_t . As such, based on Eq. (16.8), we can write

$$\mathcal{W}_2^2[p_t, q_t] \leq \mathbb{E}_z [\|\mathbf{x}_t - \mathbf{y}_t\|_2^2] \triangleq \tilde{\mathcal{W}}_2^2[p_t, q_t]. \quad (16.12)$$

Now we shall try to find the bound for $\tilde{\mathcal{W}}_2^2[p_t, q_t]$. In order to link it with $f_t(\mathbf{x}_t)$ and $g_t(\mathbf{y}_t)$, we shall take the derivative of $\tilde{\mathcal{W}}_2^2[p_t, q_t]$ with respect to t :

$$\begin{aligned} \pm \frac{d(\tilde{\mathcal{W}}_2^2[p_t, q_t])}{dt} &= \pm 2\mathbb{E}_z \left[(\mathbf{x}_t - \mathbf{y}_t) \cdot \left(\frac{d\mathbf{x}_t}{dt} - \frac{d\mathbf{y}_t}{dt} \right) \right] \\ &= \pm 2\mathbb{E}_z [(\mathbf{x}_t - \mathbf{y}_t) \cdot (f_t(\mathbf{x}_t) - g_t(\mathbf{y}_t))] \\ &= \pm 2\mathbb{E}_z \underbrace{[(\mathbf{x}_t - \mathbf{y}_t) \cdot (f_t(\mathbf{x}_t) - g_t(\mathbf{x}_t))]}_{\text{Cauchy inequality (vector form)}} \pm 2\mathbb{E}_z \underbrace{[(\mathbf{x}_t - \mathbf{y}_t) \cdot (g_t(\mathbf{x}_t) - g_t(\mathbf{y}_t))]}_{\text{Lipschitz constraint}} \\ &\leq 2\mathbb{E}_z [\|\mathbf{x}_t - \mathbf{y}_t\| \cdot \|f_t(\mathbf{x}_t) - g_t(\mathbf{x}_t)\|] + 2\mathbb{E}_z [\mathcal{L}_t \|\mathbf{x}_t - \mathbf{y}_t\|^2] \\ &\leq 2 \left(\mathbb{E}_z [\|\mathbf{x}_t - \mathbf{y}_t\|^2] \right)^{\frac{1}{2}} \left(\mathbb{E}_z [\|f_t(\mathbf{x}_t) - g_t(\mathbf{y}_t)\|^2] \right)^{\frac{1}{2}} + 2\mathbb{E}_z [\mathcal{L}_t \|\mathbf{x}_t - \mathbf{y}_t\|^2] \\ &= 2\tilde{\mathcal{W}}_2[p_t, q_t] \left(\mathbb{E}_z [\|f_t(\mathbf{x}_t) - g_t(\mathbf{y}_t)\|^2] \right)^{\frac{1}{2}} + 2\mathcal{L}_t \tilde{\mathcal{W}}_2^2[p_t, q_t], \end{aligned} \quad (16.13)$$

in which, to derive the first inequality, we need to use the vector form of the Cauchy inequality and the Lipschitz constraint [Eq. (16.10)], and we applied the expectation form of the Cauchy inequality in deriving the second inequality. The \pm sign implies that the final inequality will hold regardless of whether we take the $+$ or $-$ sign. In what follows, we take the results from the $-$ side. Using the fact that $(w^2)' = 2ww'$, in which w on both sides of the equation cancel out each other, we get

$$-\frac{d(\tilde{\mathcal{W}}_2^2[p_t, q_t])}{dt} \leq \left(\mathbb{E}_z [\|f_t(\mathbf{x}_t) - g_t(\mathbf{y}_t)\|_2^2] \right)^{\frac{1}{2}} + \mathcal{L}_t \tilde{\mathcal{W}}_2^2[p_t, q_t]. \quad (16.14)$$

Using the variation of parameters, let us further assume that $\tilde{\mathcal{W}}_2[p_t, q_t] = C_t \exp \left(\int_t^T \mathcal{L}_s ds \right)$, and substitute this into the equation above, we get

$$-\frac{dC_t}{dt} \leq \exp \left(- \int_t^T \mathcal{L}_s ds \right) \left(\mathbb{E}_z [\|f_t(\mathbf{x}_t) - g_t(\mathbf{y}_t)\|_2^2] \right)^{\frac{1}{2}}. \quad (16.15)$$

Integrating both sides over the interval $[0, T]$, and using the initial condition $C_T = 0$ (identical prior distributions with distance between them being zero), we get:

$$C_0 \leq \int_0^T \exp \left(- \int_t^T \mathcal{L}_s ds \right) \left(\mathbb{E}_z [\|f_t(\mathbf{x}_t) - g_t(\mathbf{y}_t)\|_2^2] \right)^{\frac{1}{2}} dt, \quad (16.16)$$

Not sure how the \mathcal{L}_t term in Eq. (16.14) becomes eliminated here?

as such,

This bit is also a bit puzzling.

$$\tilde{\mathcal{W}}_2[p_0, q_0] \leq C_0 \exp \left(\int_0^T \mathcal{L}_s ds \right) = \int_0^T I_t \left(\mathbb{E}_z [\|f_t(\mathbf{x}_t) - g_t(\mathbf{y}_t)\|_2^2] \right)^{\frac{1}{2}} dt, \quad (16.17)$$

in which $I_t = \exp \left(\int_0^t \mathcal{L}_s ds \right)$. According to Eq. (16.12), this is also the upper limit for $\mathcal{W}_2[p_0, q_0]$. Finally, since the expectation value is a function of \mathbf{x}_t and \mathbf{x}_t is a function of \mathbf{z} , therefore, the expectation value taken over \mathbf{z} is equivalent to the expectation value over \mathbf{x}_t , such that,

$$\tilde{\mathcal{W}}_2[p_0, q_0] \leq \int_0^T I_t \left(\mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} [\|f_t(\mathbf{x}_t) - g_t(\mathbf{y}_t)\|_2^2] \right)^{\frac{1}{2}} dt. \quad (16.18)$$

16.4. Continuing the proof

As a matter of fact, the simplified version shown in the above inequality (16.18) does not have a fundamental difference compared to the more generalised form of (16.1). Its derivation has already included all the general ideas behind the full derivation of the generalised result. Now let us finalise the remaining part of the proof.

First, we need to generalise Eq. (16.18) to the case where the initial (prior) distributions are different from each other. Assuming that the two prior distributions are $p_T(z_1)$ and $q_T(z_2)$, we sample the initial value from $p_T(z_1)$ to evolve \mathbf{x}_t , and from $q_T(z_2)$ to evolve \mathbf{y}_t . As such, \mathbf{x}_t and \mathbf{y}_t are functions of z_1 and z_2 , respectively, not the function of an identical \mathbf{z} as we assumed before. In this case, we cannot straight away build a single transport protocol, but also need another transport protocol between z_1 and z_2 . Let us choose the optimum transport protocol $\gamma^*(z_1, z_2)$ between the two distributions of $p_T(z_1)$ and $q_T(z_2)$. Similar to Eq. (16.12), we now have

$$\mathcal{W}_2^2[p_t, q_t] \leq \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \sim \gamma^*(z_1, z_2)} [\|\mathbf{x}_t - \mathbf{y}_t\|_2^2] \triangleq \tilde{\mathcal{W}}_2^2[p_t, q_t]. \quad (16.19)$$

Because we did not change the definitions, except the expectation from \mathbb{E}_z to $\mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2}$, both the procedure shown in Eq. (16.13) and the inequalities (16.14) and (16.15) still hold. However, when we take the integral over $[0, T]$ on both sides of (16.15), we no longer have $C_T = 0$, but $C_T = \tilde{\mathcal{W}}_2[p_T, q_T] = \mathcal{W}_2[p_T, q_T]$, so the final result is

$$\tilde{\mathcal{W}}_2[p_0, q_0] \leq \int_0^T I_t \left(\mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} [\|f_t(\mathbf{x}_t) - g_t(\mathbf{y}_t)\|_2^2] \right)^{\frac{1}{2}} dt + I_T \mathcal{W}_2[p_T, q_T]. \quad (16.20)$$

At this point, let us return to the diffusion model. In Section 6, we have shown that, for a given forward process, it has a set of corresponding reverse processes that are described by the following SDE:

$$d\mathbf{x} = \left(f_t(\mathbf{x}_t) - \frac{1}{2} (g_t^2 + \sigma_t^2) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \right) dt + \sigma_t d\mathbf{w}, \quad (16.21)$$

in which σ_t is a time-dependent function of standard deviation that we are free to choose, and Eq. (16.2) corresponds to the case where $\sigma_t = g_t$. As we have only considered ODE above, we will consider the case in which $\sigma_t = 0$. In this case, (16.20) still holds. What we only need to do is to change $f_t(\mathbf{x}_t)$ into $f_t(\mathbf{x}_t) - \frac{1}{2} g_t^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ and change $g_t(\mathbf{x}_t)$ into $f_t(\mathbf{x}_t) - \frac{1}{2} g_t^2 s_\theta(\mathbf{x}_t, t)$. Substituting these back into Eq. (16.20), we recover the result (16.1) shown at the beginning. Of course, do not forget that the single-sided Lipschitz constraint (16.10) that we must impose on $g_t(\mathbf{x}_t)$, which now we could further elaborate our assumption on $f_t(\mathbf{x}_t)$ and $s_\theta(\mathbf{x}_t, t)$, which we will not discuss these results here.

16.5. Further generalisation and its difficulties

Supposedly, we should carry on the work and finish the proof for the case where $\sigma_t \neq 0$. Unfortunately, the idea presented in this section cannot be fully generalised to the case of SDE, which we shall provide some analysis on why this is the case in what follows. As a matter of face, for most readers, it should be sufficient to understand the case of ODE-based derivation, the full details for SDE case are not that critically important.

For simplicity, let us take Eq. (16.2) as an example. What we need to do is to estimate the difference between the distributions for the evolving trajectories as described by the following

two SDEs:

$$\begin{cases} dx_t = \left[f_t(x_t) - g_t^2 \nabla_{x_t} \log p_t(x_t) \right] dt + g_t dw \\ dy_t = \left[f_t(y_t) - g_t^2 s_\theta(y_t, t) \right] dt + g_t dw. \end{cases} \quad (16.22)$$

In another word, we want to know how much the final distribution will change if we replace the accurate score function $\nabla_{x_t} \log p_t(x_t)$ by the estimate $s_\theta(y_t, t)$. We use the same idea of transforming the SDE into an ODE first, from which we can reuse the proof before. According to Eq. (16.21), we know that the corresponding ODE to the first SDE is

$$dx_t = \left[f_t(x_t) - \frac{1}{2} g_t^2 \nabla_{x_t} \log p_t(x_t) \right] dt. \quad (16.23)$$

Deriving the ODE for the second SDE is a bit tricky, which we need to first change it into the form of $-g_t^2 \nabla_{y_t} \log q_t(y_t)$ before applying Eq. (16.21):

$$\begin{aligned} dy_t &= \left[\underbrace{f_t(y_t) - g_t^2 s_\theta(y_t, t) + g_t^2 \nabla_{y_t} \log q_t(y_t)}_{\text{treat as a whole}} - g_t^2 \nabla_{y_t} \log q_t(y_t) \right] dt + g_t dw \\ &= \left[f_t(y_t) - g_t^2 s_\theta(y_t, t) + g_t^2 \nabla_{y_t} \log q_t(y_t) - \frac{1}{2} g_t^2 \nabla_{y_t} \log q_t(y_t) \right] dt \\ &= \left[f_t(y_t) - g_t^2 s_\theta(y_t, t) + \frac{1}{2} g_t^2 \nabla_{y_t} \log q_t(y_t) \right] dt. \end{aligned} \quad (16.24)$$

Repeat the process shown in Eq. (16.13) for the above two ODEs, then the major difference is that we now have an extra term:

$$-\frac{1}{2} g_t^2 \mathbb{E}_z \left[(x_t - y_t) (\nabla_{x_t} p_t(x_t) - \nabla_{y_t} \log q_t(y_t)) \right]. \quad (16.25)$$

If this term is negative, then the result from Eq. (16.13) will still hold. So the key question is whether we can prove the following:

$$\mathbb{E}_z \left[(x_t - y_t) (\nabla_{x_t} p_t(x_t) - \nabla_{y_t} \log q_t(y_t)) \right] \geq 0.$$

However, we could prove that this does not generally hold with specific examples. In the original article, a similar term has also appeared, except that the expectation value was taken over the transport proposals between x_t and y_t . The article has provided a simple proof, but for the newcomers with limited knowledge on the optimum transport theory, we can only stop here.

More specifically, we cannot impose the single-sided Lipschitz constraint over the score functions, as we can easily find examples where the constraint is violated. As such, to prove the above inequality, we can only follow the methodology provided in the original article, without imposing other constraints but obeying the properties of the data distributions.

17. The General Procedure of Constructing an ODE (Part II)

After completing Section 15, the author believed that a general procedure for constructing an ODE-based diffusion model had been completed. However, later, the author was brought into the attention of a new ICLR2023 article entitled “*Flow Straight and Fast, Learning to Generate and Transfer Data with Rectified Flow*”²¹, which provided another proposal for building an ODE-based diffusion model that is extremely simple and elegant. This will be discussed in this section.

17.1. The intuitive results

As we know, the diffusion model is modelling a stochastic process that evolves x_T to x_0 , whereas an ODE-based diffusion model fixes such an evolution process, by enforcing the evolution to occur by following an ODE:

$$\frac{dx_t}{dt} = f_t(x_t). \quad (17.1)$$

Original blog post see <https://www.spaces.ac.cn/archives/9497>

The core of building an ODE-based diffusion model, is to design a function $f_t(x_t)$ that describes the trajectory of the evolving sample, and simultaneously, constitutes a transformation between two given distributions $p_0(x_0)$ and $p_T(x_T)$. In other words, if we randomly sample x_T from $p_T(x_T)$ and evolve it according to the above ODE, we would like to ensure the outcome x_0 follows the sample distribution $p_0(x_0)$. The idea from the original paper was very simple, we choose $x_0 \sim p_0(x_0)$ and $x_T \sim p_T(x_T)$ randomly, and assume they transform into each other following the trajectory of

$$x_t = \varphi_t(x_t, x_0). \quad (17.2)$$

Such a trajectory is defined by a known function that we are free to design. Theoretically, it only needs to satisfy the following boundary conditions:

$$x_0 = \varphi_0(x_0, x_T), \quad x_T = \varphi_T(x_0, x_T). \quad (17.3)$$

while being a continuous function over $[x_0, x_T]$. From here, one can write out the differential equation that the trajectory function must satisfy:

$$\frac{dx_t}{dt} = \frac{d\varphi_t(x_0, x_T)}{dt}. \quad (17.4)$$

However, this differential equation does not have a practical value. This is because we want to generate x_0 from a given x_T , but the right-hand-side of Eq. (17.4) is also a function of x_0 , unless we have already known x_0 , in this case, we would not need a diffusion model in the first place. The ODE will only be practically useful if its right-hand-side is of a similar form as in Eq. (17.1), which is only a function of x_t . (Theoretically, we can also make it dependent on x_T , but we will not consider this case, in general.) In this case, an intuitive and innovative ideal was formed, that is, to learn a function $v_\theta(x_t, t)$ to regress the right-hand-side of the above ODE. To do this, we optimise the following target:

$$\mathbb{E}_{x_0 \sim p_0(x_0), x_T \sim p_T(x_T)} \left[\left\| v_\theta(x_t, t) - \frac{d\varphi_t(x_0, x_T)}{dt} \right\|_2^2 \right]. \quad (17.5)$$

When $v_\theta(x_t, t)$ is a sufficiently accurate approximation for $\frac{d\varphi_t(x_0, x_T)}{dt}$, we can think Eq. (17.4) will still hold even if we replace its right-hand-side by $v_\theta(x_t, t)$, from which we obtain a practically workable diffusion ODE:

$$\frac{dx_t}{dt} = v_\theta(x_t, t). \quad (17.6)$$

17.2. A simple example

As a simple example, we assume $T = 1$ and the sample evolves along a linear trajectory:

$$x_t = \varphi_t(x_0, x_1) = (x_1 - x_0)t + x_0. \quad (17.7)$$

In this case, we have

$$\frac{\partial \varphi_t(x_0, x_1)}{\partial t} = x_1 - x_0, \quad (17.8)$$

which leads to the following training target

$$\mathbb{E}_{x_0 \sim p_0(x_0), x_T \sim p_T(x_T)} \left[\left\| v_\theta((x_1 - x_0)t + x_0, t) - (x_1 - x_0) \right\|_2^2 \right], \quad (17.9)$$

or equivalently written as

$$\mathbb{E}_{x_0, x_t \sim p(x_0)p_t(x_t|x_0)} \left[\left\| v_\theta(x_t, t) - \frac{x_t - x_0}{t} \right\|_2^2 \right], \quad (17.10)$$

which completes the whole process of constructing a diffusion model. This result is identical to the result of linear trajectory as derived in Section 14, and it is the key model being investigated in the original article, which was dubbed at the ‘‘Rectified Flow’’. From this example

of linear trajectory, it can be seen that the process of constructing an ODE-based diffusion model using this idea only contains a few lines of derivations, which is almost incredible in the sense that it overthrows the perception that the diffusion model is very complicated.

17.3. Further proofs

Nevertheless, the results presented above can only be considered as an intuitive guess, because we have not proven that, theoretically, the outcomes of training the model based on Eq. (17.5), which leads to the ODE shown in Eq. (17.6), will indeed lead to a transformation between the two probability distributions of $p_0(\mathbf{x}_0)$ and $p_T(\mathbf{x}_T)$. To achieve this goal, the author's initial thought was to prove that the optimum solution for the target distribution from Eq. (17.5) satisfies the following equation of continuity:

$$\frac{\partial p_t(\mathbf{x}_t)}{\partial t} = -\nabla_{\mathbf{x}_t} \left(p_t(\mathbf{x}_t) \mathbf{v}_\theta(\mathbf{x}_t, t) \right). \quad (17.11)$$

If this is the case, then based on the correspondence between ODE and the equation of continuity derived in Section 12. Eq. (17.6) is indeed describing the transformation between two probability distributions of $p_0(\mathbf{x}_0)$ and $p_T(\mathbf{x}_T)$.

However, such a procedure is a bit cumbersome. Fundamentally, the equation of continuity itself is derived from the ODE through the use of the following relationship:

This is basically another way of writing Eq. (12.3).

$$\mathbb{E}_{\mathbf{x}_t + \Delta t} [\phi(\mathbf{x}_{t+\Delta t})] = \mathbb{E}_{\mathbf{x}_t} [\phi(\mathbf{x}_t + f_t(\mathbf{x}_t)\Delta t)]. \quad (17.12)$$

This means that Eq. (17.12) is a more fundamental starting point, for which we only need to prove the optimum solution for Eq. (17.5) satisfies this relationship, which contains \mathbf{x}_t as the only independent variable. For this, we write out the following, (in which we have short-handed $\varphi_t(\mathbf{x}_0, \mathbf{x}_T)$ as φ_t):

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_t + \Delta t} [\phi(\mathbf{x}_{t+\Delta t})] &= \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_T} [\phi(\varphi_{t+\Delta t})] \\ &= \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_T} \left[\phi(\mathbf{x}_t) + \Delta t \frac{\partial \varphi_t}{\partial t} \cdot \nabla_{\varphi} \phi(\varphi_t) \right] \\ &= \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_T} [\phi(\mathbf{x}_t)] + \Delta t \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_T} \left[\frac{\partial \varphi_t}{\partial t} \cdot \nabla_{\mathbf{x}_t} \phi(\mathbf{x}_t) \right] \\ &= \mathbb{E}_{\mathbf{x}_t} [\phi(\mathbf{x}_t)] + \Delta t \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_T} \left[\frac{\partial \varphi_t}{\partial t} \cdot \nabla_{\mathbf{x}_t} \phi(\mathbf{x}_t) \right]. \end{aligned} \quad (17.13)$$

Here, the first equality is derived from Eq. (17.2), the second equality comes from the first-order Taylor expansion, and the third equality also comes from Eq. (17.2). Finally, in the last equality, the expectation taken over the variables \mathbf{x}_0 and \mathbf{x}_T is equivalent to the expectation taken over \mathbf{x}_t , since the latter is a well-defined function of \mathbf{x}_0 and \mathbf{x}_T .

Here, we can see that $\frac{\partial \varphi_t}{\partial t}$ is also a function of \mathbf{x}_0 and \mathbf{x}_T . We can make a further assumption that the function defined in Eq. (17.2) with respect to \mathbf{x}_T is an invertible function, meaning we can solve for $\mathbf{x}_T = \psi(\mathbf{x}_0, \mathbf{x}_t)$ from Eq. (17.2). Substituting this into $\frac{\partial \varphi_t}{\partial t}$ enables us to write the above expectation as a function of \mathbf{x}_0 and \mathbf{x}_t . With this, we have,

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_t + \Delta t} [\phi(\mathbf{x}_{t+\Delta t})] &= \mathbb{E}_{\mathbf{x}_t} [\phi(\mathbf{x}_t)] + \Delta t \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_T} \left[\frac{\partial \varphi_t}{\partial t} \cdot \nabla_{\mathbf{x}_t} \phi(\mathbf{x}_t) \right] \\ &= \mathbb{E}_{\mathbf{x}_t} [\phi(\mathbf{x}_t)] + \Delta t \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_t} \left[\frac{\partial \varphi_t}{\partial t} \cdot \nabla_{\mathbf{x}_t} \phi(\mathbf{x}_t) \right] \\ &= \mathbb{E}_{\mathbf{x}_t} [\phi(\mathbf{x}_t)] + \Delta t \mathbb{E}_{\mathbf{x}_t} \left[\underbrace{\mathbb{E}_{\mathbf{x}_0 | \mathbf{x}_t} \left[\frac{\partial \varphi_t}{\partial t} \right]}_{\text{function of } \mathbf{x}_t} \cdot \nabla_{\varphi} \phi(\varphi_t) \right] \\ &= \mathbb{E}_{\mathbf{x}_t} \left[\phi \left(\mathbf{x}_t + \Delta t \mathbb{E}_{\mathbf{x}_0 | \mathbf{x}_t} \left[\frac{\partial \varphi_t}{\partial t} \right] \right) \right]. \end{aligned} \quad (17.14)$$

The second equation holds because now $\frac{\partial \varphi_t}{\partial t}$ is also a function of \mathbf{x}_0 and \mathbf{x}_t , in this case, for the second expectation value, it can be rewritten as being taken over \mathbf{x}_0 and \mathbf{x}_t . The third

equality is equivalent to breaking down the joint distribution $p(x_0, x_t) = p(x_0|x_t)p(x_t)$, in this case, it signifies that x_0 and x_t are not independent from each other, such that we have to write $x_0|x_t$ to denote their inter-dependency. Note that $\frac{\partial \varphi_t}{\partial t}$ is a function of x_0 and x_t , such that when we take the expectation value of it over x_0 , we are left with x_t as the only independent variable. We shall see later that this is the function of x_t that we are looking for! The last equality combines the two terms above with Taylor expansion.

Now we have derived the result

$$\mathbb{E}_{x_t+\Delta t} \left[\phi(x_{t+\Delta t}) \right] = \mathbb{E}_{x_t} \left[\phi \left(x_t + \Delta t \mathbb{E}_{x_0|x_t} \left[\frac{\partial \varphi_t}{\partial t} \right] \right) \right], \quad (17.15)$$

which holds for any testing function ϕ . This implies that

$$x_{t+\Delta t} = x_t + \Delta t \mathbb{E}_{x_0|x_t} \left[\frac{\partial \varphi_t}{\partial t} \right] \Rightarrow \frac{\partial x_t}{\partial t} = \mathbb{E}_{x_0|x_t} \left[\frac{\partial \varphi_t}{\partial t} \right]. \quad (17.16)$$

This is the ODE that we are looking for. According to

$$\mathbb{E}_x[x] = \underset{\mu}{\operatorname{argmin}} \mathbb{E}_x[\|x - \mu\|^2], \quad (17.17)$$

it shows that the right-hand-side of Eq. (17.16) is exactly the optimum solution of the training target given in Eq. (17.5). this proves that the optimum solution for Eq. (17.5), which leads to Eq. (17.6), is indeed providing a transformation between two distributions of $p_0(x_0)$ and $p_T(x_T)$.

18. Score-Matching = Conditional Score-Matching

In the previous discussions, we have encountered the terms “*score-matching*” and “*conditional score-matching*” multiple times. These are the concepts that appear frequently in the literature of diffusion and energy models. Particularly in many articles where the training targets for the diffusion models were based on the “score-matching”, but for all the mainstreamed diffusion models, such as DDPM, the training target is actually established based on the conditional score-matching.

Therefore, what is the relationship between “*score-matching*” and “*conditional score-matching*”? Are they equivalent to each other? This section will try to answer these questions in details.

18.1. Score-matching

First of all, score-matching refers to the following target function for model training:

$$\mathbb{E}_{x_t \sim p_t(x_t)} \left[\|\nabla_{x_t} \log p_t(x_t) - s_\theta(x_t, t)\|_2^2 \right], \quad (18.1)$$

in which θ are the learnable parameters. Obviously, the objective of score-function is to train a model $s_\theta(x_t, t)$ to approximate $\nabla_{x_t} \log p_t(x_t)$, in which we call the term $\nabla_{x_t} \log p_t(x_t)$ as the *score*.

In the diffusion generative model, $p_t(x_t)$ is given by the following expression:

$$p_t(x_t) = \int p_t(x_t|x_0)p_0(x_0)dx_0 = \mathbb{E}_{x_0 \sim p_0(x_0)}[p_t(x_t|x_0)]. \quad (18.2)$$

Here, $p_t(x_t|x_0)$ is usually constructed as a simple distribution with analytical form (such as a conditional normal distribution). $p_0(x_0)$ is a distribution that is fixed and given, which is usually the distribution of the training data, from which we can sample real data but its analytical form is generally unknown. From Eq. (18.2), we can derive the following:

$$\begin{aligned} \nabla_{x_t} \log p_t(x_t) &= \frac{\nabla_{x_t} p_t(x_t)}{p_t(x_t)} \\ &= \frac{\int \nabla_{x_t} p_t(x_t|x_0)p_0(x_0)dx_0}{\int p_t(x_t|x_0)p_0(x_0)dx_0} \\ &= \frac{\mathbb{E}_{x_0 \sim p_0(x_0)}[\nabla_{x_t} p_t(x_t|x_0)]}{\mathbb{E}_{x_0 \sim p_0(x_0)}[p_t(x_t|x_0)]}. \end{aligned} \quad (18.3)$$

Original blog post see <https://www.spaces.ac.cn/archives/9509>

Based on our assumption, both $\nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t|\mathbf{x}_0)$ and $p_t(\mathbf{x}_t|\mathbf{x}_0)$ have known analytical expressions, so in principle, we could estimate $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ from sampling \mathbf{x}_0 . However, this involves taking the quotient of two expectation values, which leads to a biased estimate in general. As such, one will need a large batch-size when training the model using Eq. (18.1) to overcome the bias in order to reach a better outcome.

18.2. Conditional score-matching

In reality, most diffusion models apply the following conditional score-matching as the training target:

$$\mathbb{E}_{\mathbf{x}_0, \mathbf{x}_t \sim p_0(\mathbf{x}_0)p_t(\mathbf{x}_t|\mathbf{x}_0)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) - s_\theta(\mathbf{x}_t, t)\|_2^2 \right] \quad (18.4)$$

Based on our assumptions, we have an analytical form for $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0)$, so that above training target is practically computable, which requires sampling many pairs of $(\mathbf{x}_0, \mathbf{x}_t)$ to estimate the target function. More importantly, this is an unbiased estimate, so its value is not strongly dependent upon the batch size, making it more workable in model training.

In order to reveal the relationship between score-matching and conditional score-matching, we will also need another expression for the score function:

$$\begin{aligned} \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) &= \frac{\nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t)}{p_t(\mathbf{x}_t)} \\ &= \frac{\int p_0(\mathbf{x}_0) \nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t|\mathbf{x}_0) d\mathbf{x}_0}{p_t(\mathbf{x}_t)} \\ &= \frac{\int p_0(\mathbf{x}_0) p_t(\mathbf{x}_t|\mathbf{x}_0) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) d\mathbf{x}_0}{p_t(\mathbf{x}_t)} \\ &= \int p_t(\mathbf{x}_0|\mathbf{x}_t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) d\mathbf{x}_0 \\ &= \mathbb{E}_{\mathbf{x}_0 \sim p_t(\mathbf{x}_0|\mathbf{x}_t)} \left[\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) \right]. \end{aligned} \quad (18.5)$$

18.3. The inequality relationship

First of all, we can very quickly prove that the conditional score-matching is an upper bound for score-matching, meaning when we minimise the loss of conditional score-matching, we are also minimising the score-matching loss to some extent. This is not difficult to prove, which has already been demonstrated in Section 16, which shall be represented as following:

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) - s_\theta(\mathbf{x}_t, t)\|_2^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\left\| \mathbb{E}_{\mathbf{x}_0 \sim p_t(\mathbf{x}_0|\mathbf{x}_t)} \left[\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) \right] - s_\theta(\mathbf{x}_t, t) \right\|_2^2 \right] \\ &\leq \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \mathbb{E}_{\mathbf{x}_0 \sim p_t(\mathbf{x}_0|\mathbf{x}_t)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) - s_\theta(\mathbf{x}_t, t)\|_2^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0), \mathbf{x}_t \sim p_t(\mathbf{x}_t|\mathbf{x}_0)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) - s_\theta(\mathbf{x}_t, t)\|_2^2 \right], \end{aligned} \quad (18.6)$$

in which the first equality comes from Eq. (18.5), the second inequality is derived based on the Jansen's inequality. Finally, the last equality is from the Baye's rule.

18.4. The equality relationship

Here, we shall further prove that the difference between the conditional score-matching and score-matching is a term that is independent from the trainable parameters.

To prove this, we start from score-matching, which we have

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) - s_\theta(\mathbf{x}_t, t)\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|^2 + \|s_\theta(\mathbf{x}_t, t)\|^2 - 2s_\theta(\mathbf{x}_t, t) \cdot \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \right]. \end{aligned} \quad (18.7)$$

Whereas for the conditional score-matching, we have

$$\mathbb{E}_{\mathbf{x}_0, \mathbf{x}_t \sim p_0(\mathbf{x}_0)p_t(\mathbf{x}_t|\mathbf{x}_0)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) - s_\theta(\mathbf{x}_t, t)\|^2 \right]$$

$$\begin{aligned}
&= \mathbb{E}_{\mathbf{x}_0, \mathbf{x}_t \sim p_0(\mathbf{x}_0)p_t(\mathbf{x}_t|\mathbf{x}_0)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0)\|^2 + \|s_\theta(\mathbf{x}_t, t)\|^2 - 2s_\theta(\mathbf{x}_t, t) \cdot \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) \right] \\
&\quad (18.8) \\
&= \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t), \mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} \left[\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0)\|^2 + \|s_\theta(\mathbf{x}_t, t)\|^2 - 2s_\theta(\mathbf{x}_t, t) \cdot \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) \right] \\
&= \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0)\|^2] + \|s_\theta(\mathbf{x}_t, t)\|^2 \right. \\
&\quad \left. - 2s_\theta(\mathbf{x}_t, t) \cdot \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0)\|^2] \right] \\
&= \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0)\|^2] + \|s_\theta(\mathbf{x}_t, t)\|^2 - 2s_\theta(\mathbf{x}_t, t) \cdot \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \right].
\end{aligned}$$

Taking the difference between the two, the result is

$$\mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \left[\mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)} [\|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0)\|^2] - \|\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|^2 \right], \quad (18.9)$$

which is independent from the training parameters. This means that minimising the target function from conditional score-matching is fully equivalent to minimising the target function from score-matching. This result was first demonstrated in the article “*A Connection between Score Matching and Denoising Autoencoders*”¹¹.

Since the two approaches are theoretically equivalent to each other, does it mean our previous claim that score-matching requires larger batch-size for training than the conditional score-matching will not hold? This is not true. As estimating the score $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ using Eq. (18.3) for evaluating the score-matching loss function will still contain bias which needs to be overcome with the use of large batch-size. However, when we further expanded and simplified the target function [Eq. (18.1)], we have effectively changed the biased estimate into an unbiased one, making it less dependent on the batch size. This is saying, although the two training targets are theoretically equivalent, from a statistical point of view, they are two different statistical quantities and are only exactly equivalent to each other when the sample size approach infinity.

References

- ¹Y. Lipman, R. T. Chen, H. Ben-Hamu, M. Nickel, and M. Le, “Flow matching for generative modeling,” arXiv preprint arXiv:2210.02747 (2022).
- ²A. Razavi, A. Van den Oord, and O. Vinyals, “Generating diverse high-fidelity images with VQ-VAE-2,” Adv. Neural Info. Proc. Sys. **32** (2019).
- ³J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” Adv. Neural Info. Proc. Sys. **33**, 6840 (2020).
- ⁴J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, “Deep unsupervised learning using nonequilibrium thermodynamics,” in *International conference on machine learning* (pmlr, 2015) p. 2256.
- ⁵C. Luo, “Understanding diffusion models: A unified perspective,” arXiv:2208.11970 (2022).
- ⁶M. Zhang, J. Lucas, J. Ba, and G. E. Hinton, “Lookahead optimizer: k steps forward, 1 step back,” Adv. Neural Info. Proc. Sys. **32** (2019).
- ⁷J. Song, C. Meng, and S. Ermon, “Denoising diffusion implicit models,” arXiv:2010.02502 (2020).
- ⁸Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, “Score-based generative modeling through stochastic differential equations,” arXiv:2011.13456 (2020).
- ⁹B. D. Anderson, “Reverse-time diffusion equation models,” Stochastic Processes and their Applications **12**, 313 (1982).
- ¹⁰A. Hyvärinen and P. Dayan, “Estimation of non-normalized statistical models by score matching,” J. Mach. Learn. Res. **6** (2005).
- ¹¹P. Vincent, “A connection between score matching and denoising autoencoders,” Neural Comput. **23**, 1661 (2011).
- ¹²A. Q. Nichol and P. Dhariwal, “Improved denoising diffusion probabilistic models,” in *Inter. Conf. Machine Learning* (PMLR, 2021) p. 8162.
- ¹³F. Bao, C. Li, J. Zhu, and B. Zhang, “Analytic-DPM: an analytic estimate of the optimal reverse variance in diffusion probabilistic models,” arXiv:2201.06503 (2022).
- ¹⁴F. Bao, C. Li, J. Sun, J. Zhu, and B. Zhang, “Estimating the optimal covariance with imperfect mean in diffusion probabilistic models,” arXiv:2206.07309 (2022).
- ¹⁵P. Dhariwal and A. Nichol, “Diffusion models beat GANs on image synthesis,” Adv. Neural Info. Proc. Sys. **34**, 8780 (2021).
- ¹⁶X. Liu, D. H. Park, S. Azadi, G. Zhang, A. Chopikyan, Y. Hu, H. Shi, A. Rohrbach, and T. Darrell, “More control for free! Image synthesis with semantic diffusion guidance,” in *Proceedings of the IEEE/CVF winter conference on applications of computer vision* (2023) p. 289.
- ¹⁷J. Ho and T. Salimans, “Classifier-free diffusion guidance,” arXiv:2207.12598 (2022).
- ¹⁸A. Bansal, E. Borgnia, H.-M. Chu, J. Li, H. Kazemi, F. Huang, M. Goldblum, J. Geiping, and T. Goldstein, “Cold diffusion: Inverting arbitrary image transforms without noise,” Adv. Neural Info. Proc. Sys. **36**, 41259 (2023).
- ¹⁹Y. Xu, Z. Liu, M. Tegmark, and T. Jaakkola, “Poisson flow generative models,” Adv. Neural Info. Proc. Sys. **35**, 16782 (2022).
- ²⁰D. Kwon, Y. Fan, and K. Lee, “Score-based generative modeling secretly minimizes the wasserstein distance,” Adv. Neu. Info. Proc. Sys. **35**, 20205 (2022).
- ²¹X. Liu, C. Gong, and Q. Liu, “Flow straight and fast: Learning to generate and transfer data with rectified flow,” arXiv:2209.03003 (2022).