

Score-Based Generative Models

Pramook Khungurn

April 23, 2022

In 2021, I read a paper on denoising diffusion models, which proposes a new type of generative models [HJA20]. Now, it turns out that there are parallel works by Yang Song and Stefano Ermon on the so-called “score-based models,” which is later discovered to be deeply connected (in other words, pretty much equivalent) to the former approach [SE19].¹ The approaches share the advantages of being very stable to train and being capable of generating high quality samples, and they also share the disadvantage of being slow when sampling. Researchers have claimed that these models beat GANs in image generation [DN21], and so my interest is piqued.

I read the introductory blog post by Yang Song [Son21], but it seems that I lack the background to understand this body of work. This note aims to fill this understanding gap by summarizing relevant research papers.

1 Preliminary

- We are given n data items $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ that are sampled i.i.d. from a probability distribution $p_{\text{data}}(\mathbf{x})$, which is unknown to us.
- We are interested in modeling $p_{\text{data}}(\mathbf{x})$ by finding a model $p_{\theta}(\mathbf{x})$ with parameters θ that best approximates it.
 - To reduce levels of subscription, we will sometimes write $p_{\theta}(\mathbf{x})$ as $p(\mathbf{x}; \theta)$.
- For the models in this note, we cannot compute the probability $p_{\theta}(\mathbf{x})$ directly.
- However, we would still be able to sample from it, which is something that is of practical use.
- Hence, the focus would be on (1) how to estimate the parameters θ , and (2) how to sample from the model given the parameters.

2 Score Matching [Hyvärinen 2005]

- First, however, we need to take a detour from generative modeling and study a related problem: parameter estimation of unnormalized models.
- For some probabilistic models, the form of the probability distribution is known up to the normalization constant:

$$p(\mathbf{x}) \propto q(\mathbf{x})$$

So, we have that

$$p(\mathbf{x}) = \frac{q(\mathbf{x})}{Z}.$$

¹FYI, my time in grad school overlapped with that of Stefano’s, but we never interacted.

where

$$Z = \int p(\mathbf{x}) \, d\mathbf{x}$$

is the normalization constant.

- A common class of such model is the **energy-based model**, where

$$p(\mathbf{x}) \propto e^{-E(\mathbf{x})}.$$

Here, $E(\mathbf{x})$ is called the **energy function**, and the normalization constant

$$Z = \int e^{-E(\mathbf{x})} \, d\mathbf{x}$$

is called the **partition function**. Energy-based models show up a lot in statistical physics.

- Another class of such models is the **graphical model** where

$$p(\mathbf{x}) \propto \prod_{\mathbf{a} \in \mathcal{F}} p_{\mathbf{a}}(\mathbf{x}).$$

Here, we assume that $\mathbf{x} = (x_1, x_2, \dots, x_n)$, \mathcal{F} is a set of subsets of $\{1, 2, \dots, n\}$, and $p_{\mathbf{a}}(\mathbf{x})$ is a function of components of \mathbf{x} whose set of indices are exactly \mathbf{a} . You can read more about such models in another note of mine [Khu21].

- We are interested in such probability distributions with parameters. In other words,

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{q(\mathbf{x}; \boldsymbol{\theta})}{Z(\boldsymbol{\theta})}$$

where

$$Z(\boldsymbol{\theta}) = \int q(\mathbf{x}; \boldsymbol{\theta}) \, d\mathbf{x}.$$

- We are particularly interested in estimating $\boldsymbol{\theta}$ from sampled data $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$.
- The standard approach would be to perform maximum likelihood estimation (MLE):

$$\arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p(\mathbf{x}_i; \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \left\{ \sum_{i=1}^N \log q(\mathbf{x}_i; \boldsymbol{\theta}) - N \log Z(\boldsymbol{\theta}) \right\}$$

- The general problem is that computing $Z(\boldsymbol{\theta})$ is often infeasible.
- The common way to deal with this is to estimate $Z(\boldsymbol{\theta})$ with Monte Carlo integration.
 - Markov chain Monte Carlo (MCMC) methods are often employed to generate samples that yield low variances.
 - Nevertheless, MCMC methods are slow because they need to generate many samples before convergence.
- In 2005, Aapo Hyvärinen proposed **score matching** as a way to estimate $\boldsymbol{\theta}$ without explicitly dealing with $Z(\boldsymbol{\theta})$ [Hyv05].
- The idea is to “match” the “score function” instead of doing the optimization on the probabilities directly.

- The (Stein) **score function** of a probability distribution $p(\mathbf{x}; \boldsymbol{\theta})$ is the gradient with respect to \mathbf{x} of its logarithm:

$$\boldsymbol{\Psi}(\mathbf{x}; \boldsymbol{\theta}) = \begin{bmatrix} \Psi_1(\mathbf{x}; \boldsymbol{\theta}) \\ \Psi_2(\mathbf{x}; \boldsymbol{\theta}) \\ \vdots \\ \Psi_n(\mathbf{x}; \boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \partial(\log p(\mathbf{x}; \boldsymbol{\theta}))/\partial x_1 \\ \partial(\log p(\mathbf{x}; \boldsymbol{\theta}))/\partial x_2 \\ \vdots \\ \partial(\log p(\mathbf{x}; \boldsymbol{\theta}))/\partial x_n \end{bmatrix} = \nabla_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\theta}).$$

Viewing the distribution p as a function with signature $\mathbb{R}^d \rightarrow \mathbb{R}$, we have that $\boldsymbol{\Psi}$ has signature $\mathbb{R}^d \rightarrow \mathbb{R}^d$.

- Note that the good thing about the score function is that it allows us to bypass the partition function:

$$\begin{aligned} \nabla_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\theta}) &= \nabla_{\mathbf{x}} \log \frac{q(\mathbf{x}; \boldsymbol{\theta})}{Z(\boldsymbol{\theta})} = \nabla_{\mathbf{x}} (\log q(\mathbf{x}; \boldsymbol{\theta}) - \log Z(\boldsymbol{\theta})) = \nabla_{\mathbf{x}} \log q(\mathbf{x}; \boldsymbol{\theta}) - \nabla_{\mathbf{x}} \log Z(\boldsymbol{\theta}) \\ &= \nabla_{\mathbf{x}} \log q(\mathbf{x}; \boldsymbol{\theta}). \end{aligned}$$

This is because $Z(\boldsymbol{\theta})$ is a constant with respect to \mathbf{x} .

- Recall that we are given $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ sampled from an unknown distribution p_{data} . We can define the score function of the data distribution:

$$\boldsymbol{\Psi}_{\text{data}}(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}).$$

- The “matching” in score matching is trying to find $\boldsymbol{\theta}$ that makes $\boldsymbol{\Psi}(\mathbf{x}; \boldsymbol{\theta})$ as close as possible to $\boldsymbol{\Psi}_{\text{data}}(\mathbf{x})$. Operationally, Hyvärinen proposes minimizing the expected squared Euclidean distance between the scores:

$$J(\boldsymbol{\theta}) = \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}} [\|\boldsymbol{\Psi}(\mathbf{x}; \boldsymbol{\theta}) - \boldsymbol{\Psi}_{\text{data}}(\mathbf{x})\|^2] = \frac{1}{2} \int p_{\text{data}}(\mathbf{x}) \|\boldsymbol{\Psi}(\mathbf{x}; \boldsymbol{\theta}) - \boldsymbol{\Psi}_{\text{data}}(\mathbf{x})\|^2 d\mathbf{x}.$$

The estimator is thus given by:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}).$$

- The function $J(\boldsymbol{\theta})$ is commonly known as the **Fisher divergence** between two distributions. Given two distributions p_0 and p_1 , the Fisher divergence is defined as:

$$F(p_0 \| p_1) = E_{\mathbf{x} \sim p_0} [\|\nabla_{\mathbf{x}} \log p_0(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_1(\mathbf{x})\|^2] = \int p_0(\mathbf{x}) \|\nabla_{\mathbf{x}} \log p_0(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_1(\mathbf{x})\|^2 d\mathbf{x}.$$

In other words, we are minimizing $F(p_{\text{data}}(\mathbf{x}) \| p(\mathbf{x}; \boldsymbol{\theta}))$.

- Recall again that we are only given $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$. We don't know what p_{data} is. So, how do we do the optimization then? The good news is that we can rewrite the Fisher divergence in a form that does not involve p_{data} instance the expectation operator.

- **Theorem 1.** *We have that*

$$\begin{aligned} J(\boldsymbol{\theta}) &= \int p_{\text{data}}(\mathbf{x}) \sum_{i=1}^d \left[\frac{\partial \Psi_i(\mathbf{x}; \boldsymbol{\theta})}{\partial x_i} + \frac{1}{2} (\Psi_i(\mathbf{x}; \boldsymbol{\theta}))^2 \right] d\mathbf{x} + C \\ &= \int p_{\text{data}}(\mathbf{x}) \sum_{i=1}^d \left[\frac{\partial^2 (\log p(\mathbf{x}; \boldsymbol{\theta}))}{\partial x_i^2} + \frac{1}{2} \left(\frac{\partial (\log p(\mathbf{x}; \boldsymbol{\theta}))}{\partial x_i} \right)^2 \right] d\mathbf{x} + C \end{aligned} \quad (1)$$

where C is a constant that does not depend on $\boldsymbol{\theta}$. The theorem holds under the following conditions:

- Ψ is differentiable.
- $p_{\text{data}}(\mathbf{x})$ is differentiable.
- $E_{\mathbf{x} \sim p_{\text{data}}}[\|\Psi(\mathbf{x}; \boldsymbol{\theta})\|^2]$ and $E_{\mathbf{x} \sim p_{\text{data}}}[\|\Psi_{\text{data}}(\mathbf{x})\|^2]$ are finite for every $\boldsymbol{\theta}$.
- $p_{\text{data}}(\mathbf{x})\Psi(\mathbf{x}; \boldsymbol{\theta})$ has bounded support.

The proof can be found in Hyvärinen’s paper [Hyv05], and it is based on applying integration by parts.

- Note that the RHS of (1) can be rewritten as:

$$J(\boldsymbol{\theta}) = E_{\mathbf{x} \sim p_{\text{data}}} \left[\nabla_{\mathbf{x}} \cdot \Psi(\mathbf{x}; \boldsymbol{\theta}) + \frac{1}{2} \|\Psi(\mathbf{x}; \boldsymbol{\theta})\|^2 \right] = E_{\mathbf{x} \sim p_{\text{data}}} \left[\Delta_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\theta}) + \frac{1}{2} \|\nabla_{\mathbf{x}} \log p(\mathbf{x}; \boldsymbol{\theta})\|^2 \right].$$

Here, $\nabla_{\mathbf{x}} \cdot$ is the **divergence operator**

$$\nabla_{\mathbf{x}} \cdot \mathbf{f} = \sum_{i=1}^d \frac{\partial f_i}{\partial x_i},$$

and $\Delta_{\mathbf{x}} = \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{x}}$ is the **Laplace operator**

$$\Delta_{\mathbf{x}} f = \nabla_{\mathbf{x}} \cdot \nabla_{\mathbf{x}} f = \sum_{i=1}^d \frac{\partial^2 f}{\partial x_i^2}.$$

- Theorem 1 allows us to approximate $J(\boldsymbol{\theta})$ by Monte Carlo integrations using the samples we have:

$$J(\boldsymbol{\theta}) \approx \hat{J}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{j=1}^N \left[\nabla_{\mathbf{x}} \cdot \Psi(\mathbf{x}_j; \boldsymbol{\theta}) + \frac{1}{2} \|\Psi(\mathbf{x}_j; \boldsymbol{\theta})\|^2 \right]. \quad (2)$$

This means that we can now optimize for $\boldsymbol{\theta}$.

- Hyvärinen also shows that optimizing for $J(\boldsymbol{\theta})$ would yield the right distribution.

Theorem 2. Assume that there is a unique $\boldsymbol{\theta}^*$ such that $p_{\text{data}}(\mathbf{x}) = p(\mathbf{x}; \boldsymbol{\theta}^*)$ and that $q(\mathbf{x}; \boldsymbol{\theta}) > 0$ for all \mathbf{x} and $\boldsymbol{\theta}$. Then, $J(\boldsymbol{\theta}) = 0$ if and only if $\boldsymbol{\theta} = \boldsymbol{\theta}^*$.

Corollary 3. Under the assumption of Theorem 2, the estimator $\hat{J}(\boldsymbol{\theta})$ is consistent. In other words, it converges in probability towards the true value of $J(\boldsymbol{\theta})$ when the sample size approaches infinity.

3 Denoising Score Matching [Vincent 2011]

- Minimizing the Fisher divergence using the estimate in Equation (2), however, is still not practical. This is because we need to compute

$$\nabla_{\mathbf{x}} \cdot \Psi(\mathbf{x}_j; \boldsymbol{\theta}) = \sum_{i=1}^d \frac{\partial^2 (\log p(\mathbf{x}_j; \boldsymbol{\theta}))}{\partial x_i^2},$$

which requires computing second-order derivatives. While this can certainly be arranged using modern deep learning framework, it would entail computing the Hessian of large neural network $q(\cdot; \boldsymbol{\theta})$, which is not practical.

- In 2011, Pascal Vincent discovered that, by changing the target distribution p_{data} a little, we can rewrite $J(\boldsymbol{\theta})$ into a function that does not involve second-order derivatives [Vin11].

- We change p_{data} by convolving it with an isotropic Gaussian noise. Let us denote this corrupted distribution by p_{data}^σ where σ is the standard deviation of the Gaussian. The sampling process of this distribution is as follows.

1. Sample $\mathbf{x} \sim p_{\text{data}}$.
2. Sample $\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{x}, \sigma^2 I)$ and return $\tilde{\mathbf{x}}$ as the output of the sampling process.

In this way, we have that

$$p_{\text{data}}^\sigma(\tilde{\mathbf{x}}) = \int k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}) p_{\text{data}}(\mathbf{x}) d\mathbf{x}.$$

where

$$k(\tilde{\mathbf{x}}|\mathbf{x}) = \frac{1}{(\sqrt{2\pi}\sigma)^d} \exp\left(-\frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|^2}{2\sigma^2}\right)$$

denotes the Gaussian kernel function.

- We note that p_{data}^σ would not be very different from p_{data} if σ is small enough.
- Let us perform score matching on p_{data}^σ . We have that the Fisher divergence is given by

$$\begin{aligned} J^\sigma(\theta) &= \frac{1}{2} E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\|\Psi(\tilde{\mathbf{x}}; \theta) - \Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}})\|^2 \right] \\ &= \frac{1}{2} E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\langle \Psi(\tilde{\mathbf{x}}; \theta) - \Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}}), \Psi(\tilde{\mathbf{x}}; \theta) - \Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \rangle \right] \\ &= \frac{1}{2} E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\|\Psi(\tilde{\mathbf{x}}; \theta)\|^2 - 2\langle \Psi(\tilde{\mathbf{x}}; \theta), \Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \rangle + \|\Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}})\|^2 \right] \\ &= E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\frac{\|\Psi(\tilde{\mathbf{x}}; \theta)\|^2}{2} \right] - E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\langle \Psi(\tilde{\mathbf{x}}; \theta), \Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \rangle \right] + E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\frac{\|\Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}})\|^2}{2} \right] \end{aligned}$$

Because $\|\Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}})\|^2$ is an expression that does not involve (θ) , we can treat as a constant that is irrelevant to the optimization process. As a result, we may write

$$J^\sigma(\theta) = E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\frac{\|\Psi(\tilde{\mathbf{x}}; \theta)\|^2}{2} \right] - E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\langle \Psi(\tilde{\mathbf{x}}; \theta), \Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \rangle \right] + C_1$$

Looking at the middle term, we have that

$$\begin{aligned} E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\langle \Psi(\tilde{\mathbf{x}}; \theta), \Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \rangle \right] &= \int_{\tilde{\mathbf{x}}} p_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \langle \Psi(\tilde{\mathbf{x}}; \theta), \Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \rangle d\tilde{\mathbf{x}} \\ &= \int_{\tilde{\mathbf{x}}} p_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \left\langle \Psi(\tilde{\mathbf{x}}; \theta), \frac{\partial(\log p_{\text{data}}^\sigma(\tilde{\mathbf{x}}))}{\partial \tilde{\mathbf{x}}} \right\rangle d\tilde{\mathbf{x}} \\ &= \int_{\tilde{\mathbf{x}}} \left\langle \Psi(\tilde{\mathbf{x}}; \theta), p_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \frac{\partial(\log p_{\text{data}}^\sigma(\tilde{\mathbf{x}}))}{\partial \tilde{\mathbf{x}}} \right\rangle d\tilde{\mathbf{x}}. \end{aligned}$$

Using the fact that

$$f(\mathbf{x}) \frac{\partial(\log f(\mathbf{x}))}{\partial \mathbf{x}} = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}$$

for any differentiable function f , we have that

$$\begin{aligned} p_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \frac{\partial(\log p_{\text{data}}^\sigma(\tilde{\mathbf{x}}))}{\partial \tilde{\mathbf{x}}} &= \frac{\partial p_{\text{data}}^\sigma(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}} = \frac{\partial}{\partial \tilde{\mathbf{x}}} \int_{\mathbf{x}} k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}) p_{\text{data}}(\mathbf{x}) d\mathbf{x} = \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \frac{\partial k^\sigma(\tilde{\mathbf{x}}|\mathbf{x})}{\partial \tilde{\mathbf{x}}} d\mathbf{x} \\ &= \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} d\mathbf{x} \end{aligned}$$

Plugging the integral back, we have that

$$\begin{aligned}
E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\langle \Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}), \Psi_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \rangle \right] &= \int_{\tilde{\mathbf{x}}} \left\langle \Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}), \int p_{\text{data}}(\mathbf{x}) k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} d\mathbf{x} \right\rangle d\tilde{\mathbf{x}} \\
&= \int_{\tilde{\mathbf{x}}} \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \left\langle \Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}), \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} \right\rangle d\mathbf{x} d\tilde{\mathbf{x}} \\
&= \int_{\mathbf{x}} \int_{\tilde{\mathbf{x}}} p_{\text{data}}(\mathbf{x}) k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \left\langle \Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}), \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} \right\rangle d\tilde{\mathbf{x}} d\mathbf{x} \\
&= E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^\sigma(\cdot|\mathbf{x})} \left[\left\langle \Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}), \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} \right\rangle \right].
\end{aligned}$$

As a result,

$$J^\sigma(\boldsymbol{\theta}) = E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\frac{\|\Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta})\|^2}{2} \right] - E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^\sigma(\cdot|\mathbf{x})} \left[\left\langle \Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}), \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} \right\rangle \right] + C_1.$$

Looking at the first term, we have that

$$\begin{aligned}
E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\frac{\|\Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta})\|^2}{2} \right] &= \int_{\tilde{\mathbf{x}}} p_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \frac{\|\Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta})\|^2}{2} d\tilde{\mathbf{x}} \\
&= \int_{\tilde{\mathbf{x}}} \left(\int_{\mathbf{x}} k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}) p_{\text{data}}(\mathbf{x}) d\mathbf{x} \right) \frac{\|\Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta})\|^2}{2} d\tilde{\mathbf{x}} \\
&= \int_{\mathbf{x}} \int_{\tilde{\mathbf{x}}} p_{\text{data}}(\mathbf{x}) k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \frac{\|\Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta})\|^2}{2} d\tilde{\mathbf{x}} d\mathbf{x} \\
&= E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^\sigma(\cdot|\mathbf{x})} \left[\frac{\|\Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta})\|^2}{2} \right].
\end{aligned}$$

So,

$$J^\sigma(\boldsymbol{\theta}) = E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^\sigma(\cdot|\mathbf{x})} \left[\frac{\|\Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta})\|^2}{2} - \left\langle \Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}), \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} \right\rangle \right] + C_1.$$

Now, let

$$C_2 = E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^\sigma(\cdot|\mathbf{x})} \left[\frac{1}{2} \left\| \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} \right\|^2 \right].$$

We have that

$$\begin{aligned}
J^\sigma(\boldsymbol{\theta}) &= E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^\sigma(\cdot|\mathbf{x})} \left[\frac{\|\Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta})\|^2}{2} - \left\langle \Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}), \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} \right\rangle + \frac{1}{2} \left\| \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} \right\|^2 \right] - C_2 + C_1 \\
&= \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^\sigma(\cdot|\mathbf{x})} \left[\left\| \Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}) - \frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} \right\|^2 \right] - C_2 + C_1 \\
&= \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^\sigma(\cdot|\mathbf{x})} \left[\left\| \Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}) - \nabla \log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \right\|^2 \right] - C_2 + C_1.
\end{aligned}$$

Letting $\tilde{J}^\sigma(\boldsymbol{\theta})$ denote the expectation term on the RHS, the above equation becomes:

$$J^\sigma(\boldsymbol{\theta}) = \tilde{J}^\sigma(\boldsymbol{\theta}) - C_2 + C_1. \quad (3)$$

Because C_1 and C_2 are constants with respect to $\boldsymbol{\theta}$, we have that

$$\arg \min_{\boldsymbol{\theta}} J^\sigma(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \tilde{J}^\sigma(\boldsymbol{\theta}).$$

- The above lengthy derivation tells us that, instead of minimizing $J^\sigma(\boldsymbol{\theta})$, we may minimize $\tilde{J}^\sigma(\boldsymbol{\theta})$ instead. We can see that the new objective is much nicer because

$$\frac{\partial(\log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}))}{\partial \tilde{\mathbf{x}}} = \frac{\partial}{\partial \tilde{\mathbf{x}}} \log \left(\frac{1}{(\sqrt{2\pi}\sigma)^2} \exp \left(-\frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|^2}{2\sigma^2} \right) \right) = \frac{\partial}{\partial \tilde{\mathbf{x}}} \left(-\frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|^2}{2\sigma^2} \right) = -\frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma^2}. \quad (4)$$

So,

$$\begin{aligned} \tilde{J}^\sigma(\boldsymbol{\theta}) &= \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^\sigma(\cdot|\mathbf{x})} \left[\left\| \boldsymbol{\Psi}(\tilde{\mathbf{x}}; \boldsymbol{\theta}) + \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma^2} \right\|^2 \right] \\ &= \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} \left[\left\| \boldsymbol{\Psi}(\mathbf{x} + \sigma\boldsymbol{\xi}; \boldsymbol{\theta}) + \frac{\sigma\boldsymbol{\xi}}{\sigma^2} \right\|^2 \right] \\ &= \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} \left[\left\| \boldsymbol{\Psi}(\mathbf{x} + \sigma\boldsymbol{\xi}; \boldsymbol{\theta}) + \frac{\boldsymbol{\xi}}{\sigma} \right\|^2 \right] \\ &= \frac{\sigma^2}{2} E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} \left[\left\| \sigma\boldsymbol{\Psi}(\mathbf{x} + \sigma\boldsymbol{\xi}; \boldsymbol{\theta}) + \boldsymbol{\xi} \right\|^2 \right]. \end{aligned}$$

Dropping the σ^2 factor does not change the optimization problem, so we typically write the loss function as

$$\begin{aligned} \tilde{J}^\sigma(\boldsymbol{\theta}) &= \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} \left[\left\| \sigma\boldsymbol{\Psi}(\mathbf{x} + \sigma\boldsymbol{\xi}; \boldsymbol{\theta}) + \boldsymbol{\xi} \right\|^2 \right] \\ &\approx \frac{1}{2N} \sum_{j=1}^N E_{\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} \left[\left\| \sigma\boldsymbol{\Psi}(\mathbf{x}_j + \sigma\boldsymbol{\xi}; \boldsymbol{\theta}) + \boldsymbol{\xi} \right\|^2 \right], \end{aligned}$$

which does not have any second-order derivatives.

- Let us summarize the above discussion into a theorem.

Theorem 4. *The score matching objective can be recasted as a denoising objective. In other words, we have that*

$$\arg \min_{\boldsymbol{\theta}} J^\sigma(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \tilde{J}^\sigma(\boldsymbol{\theta}).$$

where

$$\begin{aligned} J^\sigma(\boldsymbol{\theta}) &= \frac{1}{2} E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^\sigma} \left[\left\| \boldsymbol{\Psi}(\tilde{\mathbf{x}}; \boldsymbol{\theta}) - \boldsymbol{\Psi}_{\text{data}}^\sigma(\tilde{\mathbf{x}}) \right\|^2 \right], \\ \tilde{J}^\sigma(\boldsymbol{\theta}) &= \frac{1}{2} E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^\sigma(\cdot|\mathbf{x})} \left[\left\| \boldsymbol{\Psi}(\tilde{\mathbf{x}}; \boldsymbol{\theta}) - \nabla \log k^\sigma(\tilde{\mathbf{x}}|\mathbf{x}) \right\|^2 \right] \\ &= \frac{\sigma^2}{2} E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} \left[\left\| \sigma\boldsymbol{\Psi}(\mathbf{x} + \sigma\boldsymbol{\xi}; \boldsymbol{\theta}) + \boldsymbol{\xi} \right\|^2 \right] \\ &\approx \frac{\sigma^2}{2N} \sum_{j=1}^N E_{\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} \left[\left\| \sigma\boldsymbol{\Psi}(\mathbf{x}_j + \sigma\boldsymbol{\xi}; \boldsymbol{\theta}) + \boldsymbol{\xi} \right\|^2 \right]. \end{aligned}$$

- What the new objective suggests us to do is as follows.
 - For each data point \mathbf{x}_j , we corrupt it by adding a Gaussian noise $\sigma\boldsymbol{\xi}$ where $\boldsymbol{\theta}$ is sampled according to $\mathcal{N}(\mathbf{0}, I)$ to get $\tilde{\mathbf{x}} = \mathbf{x}_j + \sigma\boldsymbol{\xi}$.
 - We feed the corrupted data $\tilde{\mathbf{x}}$ to our network, and let it produce the score $\boldsymbol{\Psi}(\tilde{\mathbf{x}}; \boldsymbol{\theta}) = \nabla \log q(\tilde{\mathbf{x}}; \boldsymbol{\theta})$.

- We want $\sigma\Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta})$ to be as close as possible to $-\boldsymbol{\xi}$. In other words, if we add $\sigma^2\Psi(\tilde{\mathbf{x}}; \boldsymbol{\theta}) \approx -\sigma\boldsymbol{\xi}$ to $\tilde{\mathbf{x}}$, we should get the original signal \mathbf{x}_j back. Phrasing this another way, **we want to denoise $\tilde{\mathbf{x}}$ back to \mathbf{x}_j** .

As a result, we match the score of p_{data}^σ (and therefore p_{data} if σ is small enough) by training a denoiser!

- Note that the result in this section bears a similarity to the following result in statistics.

Theorem 5 (Tweedie’s formula). *Suppose \mathbf{x} and \mathbf{y} are random variables that satisfy the distribution $\mathbf{y} \sim \mathcal{N}(\mathbf{x}, \sigma^2 I)$. Then,*

$$E[\mathbf{x}|\mathbf{y}] = \mathbf{y} + \sigma^2 \nabla \log p(\mathbf{y})$$

where $p(\mathbf{y})$ is the probability density of \mathbf{y} .

4 Sampling from a Score-Based Model

- Suppose we have optimized $\boldsymbol{\theta}$ to match the score of p_{data} (or p_{data}^σ if we use denoising score matching). How do we sample from data points from $p_{\boldsymbol{\theta}}$ then?
- We actually do not have access to $p_{\boldsymbol{\theta}}$ because we do not know the normalizing constant $Z(\boldsymbol{\theta})$. We only know $q_{\boldsymbol{\theta}}$, which is unnormalized. This forces us to use techniques such as Markov Chain Monte Carlo (MCMC), which can sample from unnormalized distributions.
- There are many variants of MCMC algorithms [Khu22], but we should use those that exploit the score $\nabla \log q_{\boldsymbol{\theta}}(\cdot)$, which is readily available (because we have just spent a lot of time matching it). Two candidates naturally emerge:
 - **Hamiltonian Monte Carlo (HMC)**.
 - **Langevin algorithm**, for which there are two easy variants.
 - * **Unadjusted Langevin algorithm (ULA)**.
 - * **Metropolis-adjusted Langevin algorithm (MALA)**.
- ULA is the simplest of the bunch. It requires picking a small constant $\varepsilon > 0$ as a parameter, and it goes as follows:
 - Start from a random point \mathbf{x}_0 .
 - **for** $m = 1, 2, 3, \dots$
 - * Sample $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)$.
 - * Compute $\mathbf{x}_m \leftarrow \mathbf{x}_{m-1} + \frac{\varepsilon}{2} \nabla \log q_{\boldsymbol{\theta}}(\mathbf{x}_{m-1}) + \sqrt{\varepsilon} \boldsymbol{\xi}$.
 - end for**

If our target distribution is not too nasty, then the Markov chain should approach a stationary distribution after being simulated for M iterations for some big value of M . We can then take $\mathbf{x}_M, \mathbf{x}_{M+1}$, and so on as the output samples.

- ULA has the following advantages:
 - It only requires the ability to compute the score $\nabla \log q_{\boldsymbol{\theta}}(\mathbf{x}_{m-1})$. Hence, we can make our model take in the input sample \mathbf{x} and output the score directly without having to worry about what the (unnormalized) probability of \mathbf{x} would be.
 - This also greatly simplify training because, in denoising score matching, we do not have to compute any gradients with respect to \mathbf{x} . The only gradient we have to compute is that of $\tilde{J}^\sigma(\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$, which would be used in parameter optimization.

- However, ULA has several problems:
 - It may not converge to a stationary distribution.
 - If it does converge, the stationary distribution might not be p_{θ} .

Nevertheless, according to Song and Ermon, people ignore these problems and use it anyway [SE19].

- MALA and HMC, while guaranteeing the correct stationary distribution, are more complicated.
 - They require the ability to evaluate $q_{\theta}(\mathbf{x})$ when computing the acceptance probability.
 - During sampling, we have to evaluate

$$\nabla \log q_{\theta}(\mathbf{x}) = \frac{\partial(\log q_{\theta}(\mathbf{x}))}{\partial \mathbf{x}}.$$

- Moreover, after evaluating the gradient above, we have to use it in the calculation of the loss $\tilde{J}^{\sigma}(\theta)$, which we must compute the gradient

$$\frac{\partial}{\partial \theta} \tilde{J}^{\sigma}(\theta)$$

with respect to θ for optimization. This can be done with modern deep learning libraries, but the implementation would be slower than that of ULA.

- We have learned how to optimize a score-based model and also learned how to sample from it. This gives us a form of generative models.

5 Noise-Conditional Score-Based Model [Song and Ermon 2019]

- Vanilla score-based models are not practical generative models because of several reasons. The Song–Ermon paper has toy examples that illustrate these points.
 - **The manifold hypothesis.**
 - * Real-world datasets have some sort of constraints on their members, and this phenomenon is often thought of as each dataset residing in a low-dimensional manifold inside the ambient space.
 - * This means that the score is zero almost everywhere in the ambient space. Vanilla score matching thus cannot work on real-world datasets.
 - * Convoluting p_{data} with a Gaussian of small variance extends the support to the whole ambient space and allows for denoising score matching.
 - * However, we are not in the clear yet.
 - **Low data density regions.**
 - * In practice, we only have a finite number of samples of the data. In areas where the p_{data} is low, we might not have enough samples to accurately estimate the scores in that area.
 - * When two modes of the data distribution are separated by low density regions, Langevin dynamics will not be able to correctly recover the relative weights of these two modes in reasonable time and so might not converge to the true distribution.
- Song and Ermon observed that convoluting the target distribution with a Gaussian (i.e., working with p_{data}^{σ} instead of p_{data}) solves both problems simultaneously.
 - The support of any Gaussian distribution is the whole ambient space, so the data no longer resides in a low-dimensional manifold.

- If σ is high enough of the Gaussian is high enough, it can fill areas where the data distribution is low.
- We see that, as σ becomes larger, the perturbed distribution p_{data}^σ becomes more tractable but different from p_{data} . So, if we have a sequence of standard deviations

$$\sigma_1 > \sigma_2 > \dots > \sigma_L,$$

then we have that

$$p_{\text{data}}^{\sigma_1}, p_{\text{data}}^{\sigma_2}, \dots, p_{\text{data}}^{\sigma_L}$$

is a sequence of less and less tractable distributions that converge to p_{data} .

- The idea, then, is to do Langevin dynamics on the above sequence of probability distributions rather than on p_{data} or on a single p_{data}^σ .
 - The algorithm would have L phases, and the i th phase would use the score of $p_{\text{data}}^{\sigma_i}$.
 - The hope is that the steps near the beginning would allow us to explore all significant regions of p_{data} , and the steps near the end would home us in on one of the modes.
- The practice of starting with a tractable distribution and successively refine it to a target distribution is not new. It can be found in the method of **simulated annealing** in the context of optimization [KGV83], and **annealed importance sampling** in statistics [Nea98].

5.1 Training and Sampling

- Song and Ermon propose to implement the above idea with the **noise-conditional score networks** (NCSN). This is a single network $\mathbf{s}_\theta(\mathbf{x}, \sigma)$ to estimate the score $\nabla \log p_{\text{data}}^\sigma(\mathbf{x})$ of the data distribution after being convolved with an isotropic Gaussian with standard deviation σ .
- For training and sampling, the paper fixes a sequence of standard deviations $\{\sigma_i\}_{i=1}^L$. It chooses one such that

$$\frac{\sigma_1}{\sigma_2} = \frac{\sigma_2}{\sigma_3} = \dots = \frac{\sigma_{L-1}}{\sigma_L} > 1.$$

σ_1 should be large enough to fill low density areas, and σ_L should be small enough to minimize the effect on data.

- The NCSN is trained with the following objective

$$\mathcal{L}(\theta) = \frac{1}{L} \sum_{i=1}^L \lambda(\sigma_i) \tilde{J}^{\sigma_i}(\theta) = \frac{1}{L} \sum_{i=1}^L \lambda(\sigma_i) \frac{\sigma_i^2}{2} E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\sigma_i \mathbf{s}_\theta(\mathbf{x} + \sigma_i \boldsymbol{\xi}, \sigma) + \boldsymbol{\xi}\|^2].$$

- The paper chooses the weight $\lambda(\sigma) = 2/\sigma^2$ in order to make sure that the magnitude of $\lambda(\sigma)\ell(\theta, \sigma)$ does not depend on σ . So, the loss becomes

$$\mathcal{L}(\theta) = \frac{1}{L} \sum_{i=1}^L E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\sigma_i \mathbf{s}_\theta(\mathbf{x} + \sigma_i \boldsymbol{\xi}, \sigma) + \boldsymbol{\xi}\|^2]. \quad (5)$$

- The sampling algorithm, called **annealed Langevin dynamics**, is as follows.

Initialize $\tilde{\mathbf{x}}_0$.

for $i \leftarrow 1$ to L **do**

```

 $\alpha_i \leftarrow \varepsilon \cdot \sigma_i^2 / \sigma_L^2$ 
for  $t \leftarrow 1$  to  $T$  do
    Draw  $\xi_t \sim \mathcal{N}(\mathbf{0}, I)$ .
     $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \xi_t$ 
end for
 $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$ 
end for
return  $\tilde{\mathbf{x}}_0$ .

```

Here ε and T are hyperparameters of the algorithm where ε is the step size, and T is the number of time steps Langevin dynamics is simulated for each standard deviation.

- The paper chooses $\alpha_i \propto \sigma_i^2$ because it wants to make the signal to noise ratio $\frac{\|\alpha_i \mathbf{s}_\theta(\mathbf{x}, \sigma_i)\|}{2\|\sqrt{\alpha_i} \xi\|}$ constant.

5.2 Image Inpainting

- It is very easy to modify the annealed Langevin dynamics algorithm to do inpainting instead of image generation.
- The modified algorithm takes an image \mathbf{x} and a mask \mathbf{m} of the part of \mathbf{x} that should be preserved.
- The inpainting algorithm is as follows:

```

Initialize  $\tilde{\mathbf{x}}_0$ .
for  $i \leftarrow 1$  to  $L$  do
     $\alpha_i \leftarrow \varepsilon \cdot \sigma_i^2 / \sigma_L^2$ 
    Draw  $\tilde{\xi} \sim \mathcal{N}(\mathbf{0}, \sigma_i^2 I)$ .
     $\mathbf{y} \leftarrow \mathbf{x} + \tilde{\xi}$ .
    for  $t \leftarrow 1$  to  $T$  do
        Draw  $\xi_t \sim \mathcal{N}(\mathbf{0}, I)$ .
         $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \xi_t$ 
         $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_t \otimes (\mathbf{1} - \mathbf{m}) + \mathbf{y} \otimes \mathbf{m}$ .
    end for
     $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$ 
end for
return  $\tilde{\mathbf{x}}_0$ .

```

5.3 Network Architecture

- The network architecture used in the paper relies on three importance components.
 - Conditional instance normalization.
 - Dilated convolutions.
 - U-Net architecture.
- Conditional instance normalization.

- Let \mathbf{x} denote a feature tensor with C channels, and let \mathbf{x}_k denote its k th channel. Let μ_k and s_k denote the mean and the standard deviation of elements of \mathbf{x}_k . In standard instance normalization, the k th channel of the output tensor \mathbf{z} is given by

$$\mathbf{z}_k = \gamma[k] \frac{\mathbf{x}_k - \mu_k}{s_k} + \beta[k]$$

where $\gamma = (\gamma[1], \gamma[2], \dots, \gamma[C])$ and $\beta = (\beta[1], \beta[2], \dots, \beta[C])$ are learnable parameters.

- The paper observed that instance normalization removes all information about μ_k , and so may lead to color shifts when generating images. To alleviate this problem, it tries to introduce μ_k back as follows:

$$\mathbf{z}_k = \gamma[k] \frac{\mathbf{x}_k - \mu_k}{s_k} + \beta[k] + \alpha[k] \frac{\mu_k - m}{v}$$

where m is the mean of the μ_k 's among the channels, v is the standard deviation of the μ_k 's, and $\alpha = (\alpha[1], \dots, \alpha[C])$ is another learnable parameter.

- The “conditional” part of the conditional instance normalization refers to the fact that the paper uses different α , β , and γ for each of the σ_i . Hence, the equation becomes:

$$\mathbf{z}_k = \gamma[i, k] \frac{\mathbf{x}_k - \mu_k}{s_k} + \beta[i, k] + \alpha[i, k] \frac{\mu_k - m}{v}.$$

- Dilate convolution.
 - The paper replaces all subsampling layers in its networks with dilated convolution, except the first subsampling layer.
- U-Net architecture.
 - The overall architecture is that of RefineNet [LMSR16], which is a U-Net that incorporates ResNet’s design.
 - Customizations:
 - * 4-cascaded RefineNet.
 - * Pre-activation residual blocks.
 - * Replace all batch normalization with conditional instance normalization.
 - * Replace max pooling layer in Refine Blocks with average pooling in order to get smoother images.
 - * Add a conditional instance normalization before each convolution and average pooling in the Refine Blocks.
 - * Activation function is ELU.
 - * Use dilated convolutions instead of subsampling layers in residual blocks, except the first one.
 - * Dilation is increased by a factor of 2 when going to the next cascade.

5.4 Experimental and Training Setup

- The paper used NCSNs to generate images from MNIST, CelebA, and CIFAR-10. All images are of size 32×32 . Each pixel has value in range $[0, 1]$.
- The paper chooses $L = 10$, $\sigma_1 = 1$, and $\sigma_{10} = 0.01$.
- For sampling the paper chooses $T = 10$ and $\varepsilon = 2 \times 10^{-5}$. The initial value of \mathbf{x}_0 is sampled from a uniform distribution.
- All models were trained with Adama with learning rate of 0.001 for 200000 iterations. The batch size was 128.

6 Improved Training [Song and Ermon 2020]

- Song and Ermon published a followed up paper [SE20] to the one in the last section [SE19]. It does not present a new algorithm but rather introduces various improvements to the model.
- It turns out that there are several problems with the model in [SE19].
 - The hyperparameter settings such as $L = 10$, $\sigma_1 = 1$, and $\sigma_{10} = 0.01$ were all ad hoc.
 - The authors couldn't train models to general images with resolution higher than 32×32 .
- To address the above problems, the authors propose several techniques to improve training. These include:
 - A heuristic to choose the noise scale.
 - A simple improvement to the model architecture.
 - A recommendation to use exponential moving average (EMA) of model parameters while training.
- With all the techniques in the last item, the authors were able to train score-based models to on FFHQ and LSUN datasets to generate images with resolutions as large as 256×256 .
- We will now talk about each of the techniques in more details.

6.1 Technique 1: How to choose the initial noise scale

- In [SE19], σ_1 is set to 1 perhaps because the pixel values are in the $[0, 1]$ range.
- However, in real world images, it is typical to have two data points \mathbf{x}_1 and \mathbf{x}_2 where $\|\mathbf{x}_1 - \mathbf{x}_2\| \gg 1$. Hence, even with the perturbed data distribution $p_{\text{data}}^{\sigma_1}$, it would be hard for Langevin dynamics to reach \mathbf{x}_2 when starting from \mathbf{x}_1 , and this can lead to mode collapse.
- Hence, the paper propose the following technique:

Technique 1. Choose σ_1 to be as large as the maximum Euclidean distance between all pairs of training data points.

6.2 Technique 2: How to choose other noise scales

- After choosing σ_1 , we next have to choose $\sigma_2, \dots, \sigma_L$.
- In [SE19], the paper choose σ_L to be a low value (0.01). The other noise scales are chosen so that the ratio σ_{i-1}/σ_i is a constant.
- What we really want from the noise scales, however, is that the probability density $p_{\text{data}}^{\sigma_{i-1}}$ and $p_{\text{data}}^{\sigma_i}$ should be similar enough that we can transition from the former to the latter with few problems.
- The authors performed a simplified analysis where p_{data}^{σ} is assumed to be $\mathcal{N}(\mathbf{0}, \sigma^2 I)$. They showed that, if we let $r = \|\mathbf{x}\|$, then it follows that, when the data is D -dimensional,

$$p_{\text{data}}^{\sigma}(r) = \frac{1}{2^{D/2-1}\Gamma(D/2)} \frac{r^{D-1}}{\sigma^D} \exp\left(\frac{-r^2}{2\sigma^2}\right)$$

where $\Gamma(\cdot)$ denotes the gamma function.

- Moreover, they showed that

$$r - \sqrt{D}\sigma \rightarrow \mathcal{N}(0, \sigma^2/2)$$

as $D \rightarrow \infty$. Hence, for high resolution images, it is OK to assume that

$$p_{\text{data}}^\sigma(r) \approx \mathcal{N}(\sqrt{D}\sigma, \sigma^2/2).$$

- To simplify the notation, let m_i denotes $\sqrt{D}\sigma_i$ and s_i denotes $\sqrt{\sigma_i^2/2}$. We have that

$$p_{\text{data}}^{\sigma_i}(r) \approx \mathcal{N}(m_i, s_i^2).$$

- Now, we make sure that $p_{\text{data}}^{\sigma_{i-1}}$ and $p_{\text{data}}^{\sigma_i}$ are similar. So, we want to make sure that they overlap as much as possible. In other words, $p_{\text{data}}^{\sigma_i}$ should have high mass in areas where $p_{\text{data}}^{\sigma_{i-1}}$ has high mass.
- $p_{\text{data}}^{\sigma_{i-1}}(r)$ has high mass in the interval $\mathcal{I}_{i-1} = [m_{i-1} - 3s_{i-1}, m_{i-1} + 3s_{i-1}]$. The mass of $p_{\text{data}}^{\sigma_i}$ on the above interval is given by

$$p_{\text{data}}^{\sigma_i}(r \in \mathcal{I}_{i-1}) = \text{erf}(\sqrt{2D}(\gamma_i - 1) + 3\gamma_i) - \text{erf}(\sqrt{2D}(\gamma_i - 1) - 3\gamma_i)$$

where $\gamma_i = \sigma_{i-1}/\sigma_i$. We want the above quantity to be reasonably large.

- **Technique 2.** Choose $\sigma_2, \dots, \sigma_L$ to be a geometric progression with common ratio γ such that

$$\text{erf}(\sqrt{2D}(\gamma - 1) + 3\gamma) - \text{erf}(\sqrt{2D}(\gamma - 1) - 3\gamma) \approx 0.5.$$

- I think how to choose σ_L has not change: we fix it to 0.01 or something lower.

6.3 Technique 3: How to condition with the noise scale

- In [SE19], the noise scale is used to calculate the biases and scales of the instance normalization units.
- In [SE20], they significantly simplify the architecture based on the observation that $\|\mathbf{s}_\theta(\mathbf{x}, \sigma)\| \approx 1/\sigma$.
- **Technique 3.** Parameterize the NCSN with $\mathbf{s}_\theta(\mathbf{x}, \sigma) = \mathbf{s}_\theta(\mathbf{x})/\sigma$ where \mathbf{s}_θ is an unconditional network.
- The authors found that this new architecture achieved similar training losses compared to the architecture employed in [SE19].

6.4 Technique 4: How to configure Langevin dynamics

- In [SE19], for each noise scale, we run Langevin dynamics for T timesteps with step size $\alpha = \varepsilon \cdot \sigma_i^2/\sigma_L^2$.
- The [SE19] paper uses $\varepsilon = 2 \times 10^{-5}$ and $T = 100$.
- The authors performed another simplified analysis to understand the behavior of annealed Langevin dynamics.

- They assumed that $p_{\text{data}}^{\sigma_i} = \mathcal{N}(\mathbf{0}, \sigma_i^2 I)$.
- They would like to understand the behavior of running Langevin dynamics for T iterations under a single noise level σ_i . In this setting, we start with $\mathbf{x}_0 \sim p_{\text{data}}^{\sigma_{i-1}} = \mathcal{N}(\mathbf{0}, \sigma_{i-1}^2 I)$. Then, for T times, we then run the update step

$$\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t + \alpha \nabla_{\mathbf{x}} \log p_{\text{data}}^{\sigma_i}(\mathbf{x}_t) + \sqrt{2\alpha} \boldsymbol{\xi} \quad (6)$$

where $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)$.²

²Note that (6) is slightly different from $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \boldsymbol{\xi}_t$, which was used in the annealed Langevin dynamics algorithm in Section 5.1. However, one can see that they are essentially the same as we simply replaced α_i with 2α .

- It follows that $\mathbf{x}_T \sim N(\mathbf{0}, s_T^2 I)$ where

$$\frac{s_T^2}{\sigma_i^2} = \left(1 - \frac{\varepsilon}{\sigma_L^2}\right)^{2T} \left(\gamma^2 - \frac{2\varepsilon}{\sigma_L^2 - \sigma_L^2(1 - \varepsilon/\sigma_L^2)^2}\right) + \frac{2\varepsilon}{\sigma_L^2 - \sigma_L^2(1 - \varepsilon/\sigma_L^2)^2}$$

where $\gamma = \sigma_{i-1}/\sigma_i$.

- The ideal behavior should be that $\mathbf{x}_T \sim N(\mathbf{0}, \sigma_i^2 I)$. Hence, we should strive to make $s_T^2/\sigma_i^2 = 1$.
- **Technique 4.** Choose T as large as allowed by a computing budget and then select ε that makes s_T^2/σ_i^2 as close as possible to 1.

6.5 Technique 5: How to improve stability

- The authors observed that the images generated by the model of [SE19] have unstable visual quality: the FID score fluctuated wildly during training. Moreover, the generated images show color shifts which varies wildly based on the number of training iterations.
- To improve the stability of the generated images, the authors recommended using exponential moving average (EMA) of the network weights when generating images at test time.
- **Technique 5.** Apply exponential moving average to network parameters when sampling.

7 Generalization to Stochastic Differential Equations [Song et al. 2021]

- In this new paper, the authors called the approach in the [SE19] paper **score matching with Langevin dynamics** (SMLD).
- The **denoising diffusion probability model** (DDPM) [SWMG15, HJA20] is another class of generative models whose generative process involves inverting corruption (in our case, denoising) done to a data item.
 - I wrote a note on DDPMs at [Khu20].
- By denoising, the DDPM implicitly computes the score, and so both DDPM and SMLD can be collectively called **score-gased generative models**.
- In [SSDK⁺21], the authors propose a framework that unify SMLD and DDPM through the lense of stochastic differential equations. This enables new sampling methods and extensions to the existing models.
- The framework.
 - There is a continuous-time stochastic process called the **forward process** that corrupts the data distribution p_{data} into a distribution that has a well-known form (most of the time, a multi-variate Gaussian distribution).
 - The forward process is described by a stochastic differential equation (SDE) that has it as a solution.
 - Given a forward process, there is a **reverse process** that runs the process backwards in time. It is the solution of a reverse-time SDE [And82].
 - Generation can then be casted as simulating the reverse-time SDE, starting from the well-known distribution. This can be done by any SDE solver.

- The expression for the reverse SDE has the time-dependent score (i.e., $\nabla \log p_t(\mathbf{x})$) inside. As a result, any algorithm that claims to use this framework must be able to estimate this score.
- The above framework yielded the following benefits.
 - The framework allows employment of any general-purpose SDE solver to do sampling. The paper proposes two algorithms.
 - * The **Predictor–Corrector** (PC) samplers which has a Metropolis–Hastings-like correction step. Examples include the Langevin algorithm and HMC.
 - * The **deterministic** samplers that are based on the probability flow ordinary differential equation (ODE). This approach allows:
 - Fast sampling via black-box ODE solvers.
 - Flexible data manipulation via latent codes.
 - Exact likelihood computation.
 - * The generation process can be manipulated by conditioning on information not available during training. This allows:
 - class-conditional generation,
 - image inpainting,
 - colorization, and
 - solving other inverse problems.
 - The paper also give a better architecture for SMLD that achieved new state-of-the-art Inception scores and FID scores.

7.1 Recap: SMLD and DDPM

7.1.1 SMLD

- Similar to the previous sections, let $p_{\text{data}}^\sigma(\tilde{\mathbf{x}})$ denote the distribution obtained by convolving p_{data} with the spherical Gaussian $\mathcal{N}(\mathbf{0}, \sigma^2 I)$.

$$p_{\text{data}}^\sigma(\tilde{\mathbf{x}}) = \int p_{\text{data}}(\mathbf{x}) \mathcal{N}(\tilde{\mathbf{x}}; \mathbf{x}, \sigma^2 I) d\mathbf{x}.$$

- We specify a sequence of noise scales $\sigma_{\min} = \sigma_1 < \sigma_2 < \dots < \sigma_n = \sigma_{\max}$.
 - σ_{\min} is small enough that $p_{\text{data}}^{\sigma_{\min}}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x})$.
 - σ_{\max} is big enough that $p_{\text{data}}^{\sigma_{\max}}(\mathbf{x}) \approx \mathcal{N}(\mathbf{x}; \mathbf{0}, \sigma_{\max}^2 I)$.
- We train a noise-conditional score network (NCSN) $\mathbf{s}_\theta(\mathbf{x}, \sigma)$ to estimate the score $\nabla \log p_{\text{data}}^\sigma(\mathbf{x})$.
- The parameters of the NCSN can found by optimizing the optimization problem

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\sigma_i \mathbf{s}_\theta(\mathbf{x} + \sigma_i \boldsymbol{\xi}, \sigma_i) + \boldsymbol{\xi}\|^2]$$

where the loss function comes straight from (5).

- To prepare for discussion in the next section, we shall rewrite the loss function. First, note that

$$\begin{aligned} & E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\sigma_i \mathbf{s}_\theta(\mathbf{x} + \sigma_i \boldsymbol{\xi}, \sigma_i) + \boldsymbol{\xi}\|^2] \\ &= \sigma_i^2 E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} \left[\left\| \mathbf{s}_\theta(\mathbf{x} + \sigma_i \boldsymbol{\xi}, \sigma_i) + \frac{\boldsymbol{\xi}}{\sigma_i} \right\|^2 \right] \\ &= \sigma_i^2 E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^{\sigma_i}(\cdot | \mathbf{x})} [\|\mathbf{s}_\theta(\tilde{\mathbf{x}}, \sigma_i) - \nabla \log k^{\sigma_i}(\tilde{\mathbf{x}} | \mathbf{x})\|^2]. \end{aligned}$$

So, the optimization problem is equivalent to:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N \sigma_i^2 E_{\mathbf{x} \sim p_{\text{data}}, \tilde{\mathbf{x}} \sim k^{\sigma_i}(\cdot|\mathbf{x})} [\|\mathbf{s}_{\boldsymbol{\theta}}(\tilde{\mathbf{x}}, \sigma_i) - \nabla \log k^{\sigma_i}(\tilde{\mathbf{x}}|\mathbf{x})\|^2]. \quad (7)$$

- To sample from an SMLD model, we start by sampling $\mathbf{x}_n^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_{\max}^2 I)$. Suppose that $\mathbf{x}_i^{(0)}$ has been sampled. We then run the Langevin algorithm for M steps using the following update rule:

$$\mathbf{x}_i^{(j+1)} \leftarrow \mathbf{x}_i^{(j)} + \varepsilon_i \mathbf{s}_{\boldsymbol{\theta}^*}(\mathbf{x}_i^{(j)}, \sigma_i) + \sqrt{2\varepsilon_i} \boldsymbol{\xi}$$

where $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)$. Having sampled $\mathbf{x}_i^{(M)}$, we then set $\mathbf{x}_{i-1}^{(0)} \leftarrow \mathbf{x}_i^{(M)}$. We continue this process until we have sampled $\mathbf{x}_0^{(M)}$, which is the output of the sampling algorithm.

- Observe that, after optimizing for $\boldsymbol{\theta}^*$, we would have that

$$E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\boldsymbol{\xi} + \sigma \mathbf{s}_{\boldsymbol{\theta}^*}(\mathbf{x} + \sigma_i \boldsymbol{\xi}, \sigma_i)\|^2] \approx 0.$$

In other words,

$$E_{\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\boldsymbol{\xi}\|^2] \approx E_{\mathbf{x}_0 \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\sigma \mathbf{s}_{\boldsymbol{\theta}^*}(\mathbf{x} + \sigma_i \boldsymbol{\xi}, \sigma_i)\|^2] = E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^{\sigma_i}} [\|\sigma \mathbf{s}_{\boldsymbol{\theta}^*}(\tilde{\mathbf{x}}, \sigma_i)\|^2].$$

If $\boldsymbol{\xi} \in \mathbb{R}^d$, we have that $E_{\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\boldsymbol{\xi}\|^2] = d$. So,

$$E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^{\sigma_i}} [\|\sigma \mathbf{s}_{\boldsymbol{\theta}^*}(\tilde{\mathbf{x}}, \sigma_i)\|^2] \approx d,$$

or

$$E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^{\sigma_i}} [\|\nabla \log p_{\text{data}}^{\sigma_i}(\tilde{\mathbf{x}})\|^2] \approx E_{\tilde{\mathbf{x}} \sim p_{\text{data}}^{\sigma_i}} [\|\mathbf{s}_{\boldsymbol{\theta}^*}(\tilde{\mathbf{x}}, \sigma_i)\|^2] \propto \frac{1}{\sigma^2}.$$

7.1.2 DDPM

- According to [HJA20], we specify a sequence of *variance* scales $0 < \beta_1 < \beta_2 < \dots < \beta_n < 1$. In particular, we pick $n = 1000$, and β_1, \dots, β_n forms an arithmetic progression with $\beta_1 = 10^{-4}$ and $\beta_n = 0.02$.
- The **forward process** is a Markov chain $(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N)$ where $\mathbf{x}_0 \sim p_{\text{data}}$ and

$$\mathbf{x}_i \sim \mathcal{N}(\sqrt{1 - \beta_i} \mathbf{x}_{i-1}, \beta_i I).$$

- It can be shown that

$$\mathbf{x}_i \sim \mathcal{N}(\sqrt{\bar{\alpha}_i} \mathbf{x}_0, (1 - \bar{\alpha}_i) I)$$

where $\alpha_i = 1 - \beta_i$, and $\bar{\alpha}_i = \alpha_1 \alpha_2 \dots \alpha_i$.

- With n and $\alpha_1, \dots, \alpha_n$ being chosen as in [HJA20], we would have that $\sqrt{\bar{\alpha}_n} \approx 0$, and $1 - \bar{\alpha}_n \approx 1$. Hence, \mathbf{x}_n 's distribution would be approximately $\mathcal{N}(\mathbf{0}, I)$. In other words, the forward process corrupts p_{data} by gradually adding noise and transforms it to the standard Gaussian distribution.
- We now would like to create another Markov chain called the **reverse process** $(\mathbf{x}_n, \mathbf{x}_{n-1}, \dots, \mathbf{x}_0)$ where $\mathbf{x}_n \sim \mathcal{N}(\mathbf{0}, I)$ and \mathbf{x}_{i-1} is obtained from \mathbf{x}_i in such a way that reverts the forward transition from \mathbf{x}_{i-1} to \mathbf{x}_i . In this way, \mathbf{x}_0 would be distributed according to a good approximation of p_{data} .

- Here's a heuristics to derive the reverse process's transition. Recall that

$$\mathbf{x}_i \sim \mathcal{N}(\sqrt{1 - \beta_i} \mathbf{x}_{i-1}, \beta_i I).$$

By Tweedie's formula (Theorem 5), we have that

$$\begin{aligned} E[\sqrt{1 - \beta_i} \mathbf{x}_{i-1} | \mathbf{x}_i] &= \mathbf{x}_i + \beta_i \nabla \log p(\mathbf{x}_i) \\ E[\mathbf{x}_{i-1} | \mathbf{x}_i] &= \frac{\mathbf{x}_i + \beta_i \nabla \log p(\mathbf{x}_i)}{\sqrt{1 - \beta_i}} \end{aligned}$$

So, the distribution of \mathbf{x}_i should be

$$\mathbf{x}_{i-1} \sim \mathcal{N}\left(\frac{\mathbf{x}_i + \beta_i \nabla \log p(\mathbf{x}_i)}{\sqrt{1 - \beta_i}}, \beta_i I\right).$$

(In theory, the covariance matrix of $\mathbf{x}_{i-1} | \mathbf{x}_i$ is not $\beta_i I$ but something quite close. According to [HJA20], $\beta_i I$ was chosen through experimentation.)

- We will train a network $\mathbf{s}_\theta(\mathbf{x}_i, i)$ to estimate the score $\nabla \log p(\mathbf{x}_i)$.
- The optimizing problem to find parameters of the network is the same as (7). However, we need to recast the variables so that they fall in the DDPM context.

For SMLD, at noise scale σ_i , we first sample $\mathbf{x} \sim p_{\text{data}}$, and then we sample $\tilde{\mathbf{x}}$ according to the transition kernel $k^{\sigma_i}(\tilde{\mathbf{x}} | \mathbf{x}) = \mathcal{N}(\tilde{\mathbf{x}}; \mathbf{x}, \sigma_i^2 I)$.

For DDPM, at noise scale β_i , we sample $\mathbf{x}_0 \sim p_{\text{data}}$, and then we sample $\mathbf{x}_i \sim \mathcal{N}(\sqrt{\bar{\alpha}_i} \mathbf{x}_0, (1 - \bar{\alpha}_i) I)$. The transition kernel $k_i(\mathbf{x}_i | \mathbf{x})$ is given by $\mathcal{N}(\mathbf{x}_i; \sqrt{\bar{\alpha}_i} \mathbf{x}_0, (1 - \bar{\alpha}_i) I)$.

So, the optimization problem is:

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N (1 - \bar{\alpha}_i) E_{\mathbf{x} \sim p_{\text{data}}, \mathbf{x}_i \sim k_i(\cdot, \mathbf{x})} [\|\mathbf{s}_\theta(\mathbf{x}_i, i) - \nabla \log k_i(\mathbf{x}_i | \mathbf{x})\|^2] \quad (8)$$

- Now,

$$\nabla \log k_i(\mathbf{x}_i | \mathbf{x}) = \nabla \log \left(\frac{1}{(\sqrt{2\pi(1 - \bar{\alpha}_i)})^d} \exp \left(-\frac{\|\mathbf{x}_i - \sqrt{\bar{\alpha}_i} \mathbf{x}_0\|^2}{2(1 - \bar{\alpha}_i)} \right) \right) = -\frac{\mathbf{x}_i - \sqrt{\bar{\alpha}_i} \mathbf{x}_0}{1 - \bar{\alpha}_i}.$$

So, the loss function can be rewritten as

$$\begin{aligned} \theta^* &= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N (1 - \bar{\alpha}_i) E_{\mathbf{x} \sim p_{\text{data}}, \mathbf{x}_i \sim k_i(\cdot, \mathbf{x})} \left[\left\| \mathbf{s}_\theta(\mathbf{x}_i, i) + \frac{\mathbf{x}_i - \sqrt{\bar{\alpha}_i} \mathbf{x}_0}{1 - \bar{\alpha}_i} \right\|^2 \right] \\ &= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N (1 - \bar{\alpha}_i) E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} \left[\left\| \mathbf{s}_\theta(\sqrt{\bar{\alpha}_i} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_i} \boldsymbol{\xi}, i) + \frac{\sqrt{1 - \bar{\alpha}_i} \boldsymbol{\xi}}{1 - \bar{\alpha}_i} \right\|^2 \right] \\ &= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N E_{\mathbf{x} \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\sqrt{1 - \bar{\alpha}_i} \mathbf{s}_\theta(\sqrt{\bar{\alpha}_i} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_i} \boldsymbol{\xi}, i) + \boldsymbol{\xi}\|^2]. \end{aligned}$$

- We shall now derive the loss function to train $\mathbf{s}_\theta(\cdot, \cdot)$ with. Once again, recall that

$$\mathbf{x}_i \sim \mathcal{N}(\sqrt{\bar{\alpha}_i} \mathbf{x}_0, (1 - \bar{\alpha}_i) I).$$

So, by Tweedie's formula,

$$E[\sqrt{\bar{\alpha}_i} \mathbf{x}_0 | \mathbf{x}_i] = \mathbf{x}_i + (1 - \bar{\alpha}_i) \nabla \log p(\mathbf{x}_i).$$

Going loose with mathematics, we may say that

$$\begin{aligned} \sqrt{\bar{\alpha}_i} \mathbf{x}_0 &\approx \mathbf{x}_i + (1 - \bar{\alpha}_i) \nabla \log p(\mathbf{x}_i) \\ \mathbf{0} &\approx \mathbf{x}_i - \sqrt{\bar{\alpha}_i} \mathbf{x}_0 + (1 - \bar{\alpha}_i) \nabla \log p(\mathbf{x}_i) \\ \mathbf{0} &\approx \mathbf{x}_i - \sqrt{\bar{\alpha}_i} \mathbf{x}_0 + (1 - \bar{\alpha}_i) \mathbf{s}_{\theta}(\mathbf{x}_i, i). \end{aligned}$$

Because $\mathbf{x}_i \sim \mathcal{N}(\sqrt{\bar{\alpha}_i} \mathbf{x}_0, (1 - \bar{\alpha}_i)I)$, we have that $\mathbf{x}_i = \sqrt{\bar{\alpha}_i} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_i} \boldsymbol{\xi}$ where $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)$. As a result,

$$\begin{aligned} \mathbf{0} &\approx \sqrt{1 - \bar{\alpha}_i} \boldsymbol{\xi} + (1 - \bar{\alpha}_i) \mathbf{s}_{\theta}(\sqrt{\bar{\alpha}_i} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_i} \boldsymbol{\xi}, i) \\ \mathbf{0} &\approx \boldsymbol{\xi} + \sqrt{1 - \bar{\alpha}_i} \mathbf{s}_{\theta}(\sqrt{\bar{\alpha}_i} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_i} \boldsymbol{\xi}, i) \end{aligned}$$

In other words, a good parameter $\boldsymbol{\theta}$ is one such that the above "equation" holds. As result, the optimal parameter $\boldsymbol{\theta}^*$ can be found by solving the following optimization problem:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N E_{\mathbf{x}_0 \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\boldsymbol{\xi} + \sqrt{1 - \bar{\alpha}_i} \mathbf{s}_{\theta}(\sqrt{\bar{\alpha}_i} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_i} \boldsymbol{\xi}, i)\|^2].$$

- To sample from a DDPM, we simply run the reverse process's Markov chain. We start by sampling \mathbf{x}_N according to $\mathcal{N}(\mathbf{0}, I)$. Supposing that we have that sampled \mathbf{x}_i , we have that

$$\mathbf{x}_{i-1} \sim \mathcal{N}\left(\frac{\mathbf{x}_i + \beta_i \nabla \log p(\mathbf{x}_i)}{\sqrt{1 - \beta_i}}, \beta_i I\right) \approx \mathcal{N}\left(\frac{\mathbf{x}_i + \beta_i \mathbf{s}_{\theta^*}(\mathbf{x}_i, i)}{\sqrt{1 - \beta_i}}, \beta_i I\right),$$

and so \mathbf{x}_{i-1} can be sampled by:

$$\mathbf{x}_{i-1} \leftarrow \frac{\mathbf{x}_i + \beta_i \mathbf{s}_{\theta^*}(\mathbf{x}_i, i)}{\sqrt{1 - \beta_i}} + \sqrt{\beta_i} \boldsymbol{\xi}$$

where $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)$. This sampling process is called **ancestral sampling** because it is identical to ancestral sampling (i.e., sampling according to the DAG) of the graphical model

$$p(\mathbf{x}_0) = p(\mathbf{x}_N) \prod_{i=1}^n p(\mathbf{x}_{i-1} | \mathbf{x}_i).$$

- Observe that, after optimizing for $\boldsymbol{\theta}^*$, we would have that

$$E_{\mathbf{x}_0 \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\boldsymbol{\xi} + \sqrt{1 - \bar{\alpha}_i} \mathbf{s}_{\theta^*}(\sqrt{\bar{\alpha}_i} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_i} \boldsymbol{\xi}, i)\|^2] \approx 0.$$

In other words,

$$\begin{aligned} E_{\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\boldsymbol{\xi}\|^2] &\approx E_{\mathbf{x}_0 \sim p_{\text{data}}, \boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\sqrt{1 - \bar{\alpha}_i} \mathbf{s}_{\theta^*}(\sqrt{\bar{\alpha}_i} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_i} \boldsymbol{\xi}, i)\|^2] \\ &= E_{\mathbf{x}_i} [\|\sqrt{1 - \bar{\alpha}_i} \mathbf{s}_{\theta^*}(\mathbf{x}_i, i)\|^2]. \end{aligned}$$

If $\boldsymbol{\xi} \in \mathbb{R}^d$, we have that $E_{\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)} [\|\boldsymbol{\xi}\|^2] = d$. So,

$$E_{\mathbf{x}_i} [\|\mathbf{s}_{\theta^*}(\mathbf{x}_i, i)\|^2] \approx \frac{d}{1 - \bar{\alpha}_i},$$

or

$$E_{\mathbf{x}_i} [\|\nabla \log p(\mathbf{x}_i)\|^2] \approx E_{\mathbf{x}_i} [\|\mathbf{s}_{\theta^*}(\mathbf{x}_i, i)\|^2] \propto \frac{1}{1 - \bar{\alpha}_i}.$$

7.2 Generalization with SDEs

- We would like to construct a stochastic process $\{\mathbf{x}(t) : 0 \leq t \leq T\}$ such that
 - $\mathbf{x}(0) \sim p_0 = p_{\text{data}}$, and
 - $\mathbf{x}(T) \sim p_T$, where p_T has a tractable form.

The distribution p_T is typically an unstructured distribution with no information about p_0 such as a Gaussian distribution with fixed mean and variance.

- The stochastic process can be modeled as a the solution of the SDE

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t) dt + g(t) d\mathbf{W},$$

where

- \mathbf{W} is the standard Wiener process in \mathbb{R}^d ,
- $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a vector function called the **drift coefficient**, and
- $g : \mathbb{R} \rightarrow \mathbb{R}$ is a scalar function called the **diffusion coefficient**.
- Let $p_t(\mathbf{x})$ denote the distribution of $\mathbf{x}(0)$, and let $p_{st}(\mathbf{x}(t)|\mathbf{x}(s))$ denote the transition kernel from $\mathbf{x}(s)$ and $\mathbf{x}(t)$ where $0 \leq s < t \leq T$.
- According to Anderson [And82], there is another stochastic process where time runs backward, starting from $\mathbf{x}(T) \sim p_T$ and ending at $\mathbf{x}(0) \sim p_0$. The stochastic process is the solution of the following SDE:

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla \log p_t(\mathbf{x})] dt + g(t) d\overline{\mathbf{W}}$$

where $\overline{\mathbf{W}}$ is a standard Wiener process where time flows backward, and dt is an infinitesimal negative timestep.

- Hence, if we know how to evaluate $\nabla \log p_t(\mathbf{x})$ for all t , we can sample by simulating the backward SDE starting from p_T and going to p_0 .
- We shall train a score network $\mathbf{s}_\theta(\mathbf{x}, t)$ that estimates the score $\nabla \log p_t(\mathbf{x})$. The optimization problem for finding the parameters of the network is the continuous generalization of (7) and (8).

$$\theta^* = \arg \min_{\theta} E_{t \sim \text{Uniform}(0, T), \mathbf{x}(0) \sim p_{\text{data}}, \mathbf{x}(t) \sim p_{0t}(\cdot | \mathbf{x}(0))} [\lambda(t) \|\mathbf{s}_\theta(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_{0t}(\mathbf{x}(t) | \mathbf{x}(0))\|^2] \quad (9)$$

where $\lambda : [0, T] \rightarrow \mathbb{R}^+$ is a positive weighting function. As in both SMLD and DDPM, we typically choose $\lambda(t)$ so that

$$\lambda(t) \propto \frac{1}{E_{\mathbf{x}(t)} [\|\nabla_{\mathbf{x}(t)} \log p_t(\mathbf{x}(t))\|^2]}.$$

- In order to train, we need to know the transition kernel $p_{0t}(\mathbf{x}(t) | \mathbf{x}(0))$.
 - For both SMLD and DDPM, the kernels are Gaussians.
 - If $\mathbf{f}(\mathbf{x}, t)$ is an affine function of \mathbf{x} , then the solution is a Gaussian process, which means that the kernel is also a Gaussian distribution [SS19].
 - The paper also discusses how to deal with other types of kernels. In such a case, denoising score matching cannot be used.
- In the next two subsections, we will show that the SDE framework is a generalization of both SMLD and DDPM by deriving the SDEs of both approaches.

7.2.1 SMLD's SDE

- First, we shall cast SMLD's noise scales in terms of a Markov chain. This will allow us to say that the chain is a discretization of an SMD.
- Let $\mathbf{x}_i \sim p_{\text{data}}^{\sigma_i}$. We have that the transition kernel from \mathbf{x}_0 to \mathbf{x}_i is given by

$$k_{0i}(\mathbf{x}_i|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_i; \mathbf{x}_0, \sigma_i^2).$$

- Now, consider the sequence $(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N)$. To cast it as a Markov chain, we need to find transition kernels between consecutive terms. It follows that you can write

$$\mathbf{x}_i = \mathbf{x}_{i+1} + \sqrt{\sigma_i^2 - \sigma_{i-1}^2} \boldsymbol{\xi}_i$$

where $\boldsymbol{\xi}_i \sim \mathcal{N}(0, I)$, and we set $\sigma_0 = 0$. To see if this is the case, we have that

$$\mathbf{x}_i - \mathbf{x}_0 = \sum_{j=0}^i \sqrt{\sigma_j^2 - \sigma_{j-1}^2} \boldsymbol{\xi}_j. \quad (10)$$

Because the $\boldsymbol{\xi}_j$ s are Gaussians, it follows that $\mathbf{x}_i - \mathbf{x}_0$ is also a Gaussian. The mean and covariance matrix are given by

$$\begin{aligned} E[\mathbf{x}_i - \mathbf{x}_0] &= \sum_{j=1}^i \sqrt{\sigma_j^2 - \sigma_{j-1}^2} E[\boldsymbol{\xi}_j] = \mathbf{0} \\ \text{Cov}(\mathbf{x}_i - \mathbf{x}_0) &= \sum_{j=1}^i (\sigma_j^2 - \sigma_{j-1}^2) \text{Cov}(\boldsymbol{\xi}_j) = \sigma_i^2 \text{Cov}(\boldsymbol{\xi}_i) = \sigma_i^2 I. \end{aligned}$$

This means that

$$k_{0i}(\mathbf{x}_i|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_i; \mathbf{x}_0, \sigma_i^2 I)$$

as we have wanted.

- To turn the Markov chain interpretation into an SDE, we first perform a cosmetic change. Let $t = i/n$ and $\Delta t = 1/n$. Moreover, we shall now parameterize the Markov chain with t instead of i , so \mathbf{x}_i becomes $\mathbf{x}(t)$, and σ_i as $\sigma(t)$, and so on. Equation (10) can be rewritten as:

$$\begin{aligned} \mathbf{x}(t + \Delta t) &= \mathbf{x}(t) + \sqrt{\sigma(t + \Delta t)^2 - \sigma(t)^2} \boldsymbol{\xi} \\ \mathbf{x}(t + \Delta t) - \mathbf{x}(t) &\approx \sqrt{\frac{d[\sigma(t)^2]}{dt} \Delta t} \boldsymbol{\xi}(t) = \sqrt{\frac{d[\sigma(t)^2]}{dt}} (\sqrt{\Delta t} \boldsymbol{\xi}) \end{aligned}$$

Note that we dropped $\boldsymbol{\xi}$'s dependence on t because it doesn't have any to begin with. Now, according to the definition of the standard Wiener process, we have that

$$\mathbf{W}(t + \Delta t) - \mathbf{W}(t) = \sqrt{\Delta t} \boldsymbol{\xi}.$$

Hence,

$$\mathbf{x}(t + \Delta t) - \mathbf{x}(t) \approx \sqrt{\frac{d[\sigma(t)^2]}{dt}} (\mathbf{W}(t + \Delta t) - \mathbf{W}(t)).$$

Taking the limit as $n \rightarrow \infty$ (in other words, $\Delta t \rightarrow 0$), we have that

$$d\mathbf{x} = \sqrt{\frac{d[\sigma(t)^2]}{dt}} d\mathbf{W}.$$

The paper calls this SDE the **variance exploding** (VE) equation.

- According to [SS19], the VE SDE is a **linear SDE** because both the coefficient functions of dt and $d\mathbf{W}$ on the RHS are affine functions of \mathbf{x} .
- Moreover, the solution to a linear SDE is a Gaussian process, and its property is completely determined by the mean

$$\boldsymbol{\mu}(t) = E[\mathbf{x}(t)]$$

and the covariance matrix

$$\Sigma(t) = E[(\mathbf{x}(t) - \boldsymbol{\mu}(t))(\mathbf{x}(t) - \boldsymbol{\mu}(t))^T].$$

- Now, suppose we are an SDE of the form

$$d\mathbf{x} = f(t)\mathbf{x} dt + g(t) d\mathbf{W}.$$

Then, according to Equation (6.2) from [SS19], the mean and the covariance matrix of the solution satisfy the following two equations:

$$\frac{d\boldsymbol{\mu}(t)}{dt} = f(t)\boldsymbol{\mu}(t) \tag{11}$$

$$\frac{d\Sigma(t)}{dt} = 2f(t)\Sigma(t) + (g(t))^2 I. \tag{12}$$

- Applying (11) to the VP SDE, we have that

$$\frac{d\boldsymbol{\mu}(t)}{dt} = \mathbf{0},$$

which means that $\boldsymbol{\mu}(t) = \boldsymbol{\mu}(0)$ for all $t \geq 0$.

- Applying (12) to the VP SDE, we have that

$$\frac{d\Sigma(t)}{dt} = \frac{d[\sigma(t)]^2}{dt} I.$$

In other words, $\Sigma(t) = \Sigma(0) + (\sigma^2(t) - \sigma^2(0))I$ for all $t \geq 0$.

- Having derived $\boldsymbol{\mu}(t)$ and $\Sigma(t)$, we conclude that the transition kernel of the VE SDE is given by

$$p_{0t}(\mathbf{x}(t)|\mathbf{x}(0)) = \mathcal{N}(\mathbf{x}(t); \mathbf{x}(0), (\sigma^2(t) - \sigma^2(0))I).$$

This is done by assuming that $\mathbf{x}(0)$ is a constant (so $\Sigma(0) = 0$).

7.2.2 DDPM's SDE

- Unfortunately, the derivation in the [SSDK⁺21] was not at all rigorous. So, here I am trying to write a sensible derivation.
- The paper says that the following SDE generalizes DDPM:

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x} dt + \sqrt{\beta(t)} d\mathbf{W}.$$

The paper calls this the **variance preserving** (VP) SDE.

- Again, the VP SDE above is a linear SDE.

- Applying (11), the mean satisfies the differential equation

$$\frac{d\boldsymbol{\mu}(t)}{dt} = -\frac{1}{2}\beta(t)\boldsymbol{\mu}(t).$$

This gives

$$\boldsymbol{\mu}(t) = \boldsymbol{\mu}(0)e^{-\frac{1}{2}\int_0^t \beta(s) ds}.$$

- Applying (12), the covariance matrix satisfies the differential equation

$$\frac{d\Sigma(t)}{dt} = 2\left(-\frac{1}{2}\beta(t)\right)\Sigma(t) + \beta(t)I = -\beta(t)\Sigma(t) + \beta(t)I.$$

We can solve it as follows:

$$\begin{aligned} \frac{d\Sigma(t)}{dt} + \beta(t)\Sigma(t) &= \beta(t)I \\ e^{\int_0^t \beta(s) ds} \frac{d\Sigma(t)}{dt} + \beta(t)e^{\int_0^t \beta(s) ds} \Sigma(t) &= \beta(t)e^{\int_0^t \beta(s) ds} I \\ \frac{d}{dt} \left[e^{\int_0^t \beta(s) ds} \Sigma(t) \right] &= \frac{d}{dt} \left[e^{\int_0^t \beta(s) ds} \right] I \\ e^{\int_0^t \beta(s) ds} \Sigma(t) - \Sigma(0) &= (e^{\int_0^t \beta(s) ds} - 1)I \\ \Sigma(t) &= e^{-\int_0^t \beta(s) ds} \Sigma(0) + (1 - e^{-\int_0^t \beta(s) ds})I. \end{aligned}$$

- As a result, we have that the transition kernel is given by:

$$p_{0t}(\mathbf{x}(t)|\mathbf{x}(0)) = \mathcal{N}(\mathbf{x}(t); e^{-\frac{1}{2}\int_0^t \beta(s) ds} \mathbf{x}(0), (1 - e^{-\int_0^t \beta(s) ds})I).$$

- Note that the transition kernel for the Markov chain of DDPM is

$$k_i(\mathbf{x}_i|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_i; \sqrt{\bar{\alpha}_i} \mathbf{x}_0, (1 - \bar{\alpha}_i)I).$$

This is not exactly the same because $\alpha_i = 1 - \beta_i$. However, the two kernels are very similar in form. The mean is a product of square roots, and the variance is one minus the product of many noise scales.

- So, I guess it is fine as a generalization, but the derivation in the paper gave me so many WTF moments that I decided not to follow it.

7.2.3 Sub-VP SDE

- The paper also proposes another SDE called the **sub-VP SDE**:

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x} dt + \sqrt{\beta(t)(1 - e^{-2\int_0^t \beta(s) ds})} d\mathbf{W}.$$

- Its mean is the same as that of the VP SDE because the differential equation for the mean is the same.
- However, the variance equation becomes

$$\frac{d\Sigma(t)}{dt} = -\beta(t)\Sigma(t) + \beta(t)(1 - e^{-2\int_0^t \beta(s) ds})I.$$

Solving the equation, we have that

$$\begin{aligned}
e^{\int_0^t \beta(s) ds} \frac{d\Sigma(t)}{dt} + \beta(t) e^{\int_0^t \beta(s) ds} \Sigma(t) &= \beta(t) (e^{\int_0^t \beta(s) ds} - e^{-\int_0^t \beta(s) ds}) I \\
\frac{d}{dt} \left[e^{\int_0^t \beta(s) ds} \Sigma(t) \right] &= \frac{d}{dt} \left[e^{\int_0^t \beta(s) ds} + e^{-\int_0^t \beta(s) ds} \right] I \\
e^{\int_0^t \beta(s) ds} \Sigma(t) - \Sigma(0) &= (e^{\int_0^t \beta(s) ds} + e^{-\int_0^t \beta(s) ds} - 2) I \\
\Sigma(t) &= e^{-\int_0^t \beta(s) ds} \Sigma(0) + (1 + e^{-2\int_0^t \beta(s) ds} - 2e^{-\int_0^t \beta(s) ds}) I \\
\Sigma(t) &= e^{-\int_0^t \beta(s) ds} \Sigma(0) + (1 - e^{-\int_0^t \beta(s) ds})^2 I.
\end{aligned}$$

- As a result, the transition kernel is given by

$$p_{0t}(\mathbf{x}(t)|\mathbf{x}(0)) = \mathcal{N}(\mathbf{x}(t); e^{-\frac{1}{2}\int_0^t \beta(s) ds} \mathbf{x}(0), (1 - e^{-\int_0^t \beta(s) ds})^2 I).$$

- Because

$$(1 - e^{-\int_0^t \beta(s) ds})^2 \leq 1 - e^{-\int_0^t \beta(s) ds},$$

it follows that

$$\Sigma_{\text{subVP}}(t) \leq \Sigma_{\text{VP}}(t)$$

for all $t \geq 0$. Hence, the solution of the sub-VP SDE has less variance than the solution of the VP SDE. This is one of the advantage of using it instead of the VP SDE.

7.3 Sampling by Solving the Reverse SDE

- The discussion in the previous section gives us a number of SDEs. If we can compute and sample according to the transition kernel $p_{0t}(\mathbf{x}(t)|\mathbf{x}(0))$ of such an SDE, then we can train a score network \mathbf{s}_{θ^*} according to (9).
- After we have the score network, we can use it to construct the reverse-time SDE. We then simulate the SDE numerically in order to generate samples from $p_0 = p_{\text{data}}$.
- The paper discusses three approaches to solving the reverse-time SDE:
 - Using general-purpose numerical SDE solvers.
 - Using predictor-corrector samplers (i.e., score-based MCMC).
 - Using neural ODEs.

We will discuss these approaches in turn.

7.3.1 General-Purpose Numerical SDE Solvers

- You can look up the solvers in [KP13]. All of them can be used.
- In the paper, though, the aim of mentioning this approach is to (1) to cast the ancestral sampling procedure of DDPM in terms of a numerical SDE solver, and (2) to derive an ancestral sampling procedure for SMLD.
- However, I will not derive ancestral sampling for SMLD in this note. This is because I think it has nothing to do with the framework at all. The paper simply use the derivation in [HJA20] without connecting it to any SDEs.

DDPM's Ancestral Sampling as Numerically Solving an SDE

- Suppose we have an SDE of the form

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t) dt + g(t) d\mathbf{W}.$$

Then, a way to discretize it would be

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \mathbf{f}_i(\mathbf{x}_i) + g_i \boldsymbol{\xi}_i$$

where $\boldsymbol{\xi}_i \sim \mathcal{N}(\mathbf{0}, I)$ for all $i = 0, 1, \dots, N-1$. Note that, here, we have absorbed the discretization times into \mathbf{f}_i and g_i .

- The reverse-time SDE of the above SDE is given by

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g^2(t) \nabla \log p_t(\mathbf{x})] dt + g(t) d\bar{\mathbf{W}}.$$

The paper proposes to discretize this in the same way as we did above, so we have

$$\mathbf{x}_i = \mathbf{x}_{i+1} - f_{i+1}(\mathbf{x}_{i+1}) + g_{i+1}^2 \mathbf{s}_{\theta}(\mathbf{x}_{i+1}, i+1) + g_{i+1} \boldsymbol{\xi}_{i+1}. \quad (13)$$

Here, we use the score network to estimate the score $\nabla \log p_t(\mathbf{x})$.

- The paper calls using the discretization in (13) to solve the reverse-time SDE the **reverse diffusion sampler**.
- The paper wants to show that the DDPM update equation,

$$\mathbf{x}_i = \frac{\mathbf{x}_{i+1} + \beta_{i+1} \mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, i+1)}{\sqrt{1 - \beta_{i+1}}} + \sqrt{\beta_{i+1}} \boldsymbol{\xi}_{i+1}, \quad (14)$$

can also be written in the form of (13).

- However, we cannot do it exactly because of there is division by $\sqrt{1 - \beta_{i+1}}$ to both the \mathbf{x}_{i+1} and the $\beta_{i+1} \mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, i+1)$ term.
- The paper, instead, shows that (14) can be written in the form of (13) when we assume that β_{i+1} is infinitesimally small.
- So, let us assume then that β_{i+1} is infinitesimally small. In other words, we may assume that $\beta_{i+1}^2 \approx 0$. Taking the Taylor series of $1/\sqrt{1 - \beta_{i+1}}$ around $\beta_{i+1} = 0$, we have that

$$\frac{1}{\sqrt{1 - \beta_{i+1}}} = 1 + \frac{1}{2} \beta_{i+1} + O(\beta_{i+1}^2),$$

which allows us to perform the following approximation:

$$\frac{1}{\sqrt{1 - \beta_{i+1}}} \approx 1 + \frac{1}{2} \beta_{i+1}.$$

We now have that

$$\begin{aligned} \mathbf{x}_i &= \frac{1}{\sqrt{1 - \beta_{i+1}}} (\mathbf{x}_{i+1} + \beta_{i+1} \mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, i+1)) + \sqrt{\beta_{i+1}} \boldsymbol{\xi}_{i+1} \\ &\approx \left(1 + \frac{1}{2} \beta_{i+1}\right) (\mathbf{x}_{i+1} + \beta_{i+1} \mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, i+1)) + \sqrt{\beta_{i+1}} \boldsymbol{\xi}_{i+1} \\ &= \left(1 + \frac{1}{2} \beta_{i+1}\right) \mathbf{x}_{i+1} + \beta_{i+1} \mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, i+1) + \frac{1}{2} \beta_{i+1}^2 \mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, i+1) + \sqrt{\beta_{i+1}} \boldsymbol{\xi}_{i+1} \\ &\approx \left(1 + \frac{1}{2} \beta_{i+1}\right) \mathbf{x}_{i+1} + \beta_{i+1} \mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, i+1) + \sqrt{\beta_{i+1}} \boldsymbol{\xi}_{i+1} \\ &= \mathbf{x}_{i+1} - \left(-\frac{1}{2} \beta_{i+1} \mathbf{x}_i\right) + \beta_{i+1} \mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, i+1) + \sqrt{\beta_{i+1}} \boldsymbol{\xi}_{i+1}. \end{aligned}$$

This tells us that the DDPM update equation is an approximate discretization of the reverse-time SDE

$$d\mathbf{x} = \left(-\frac{1}{2}\beta(t)\mathbf{x} - \beta(t)\nabla \log p_t(\mathbf{x}) \right) dt + \sqrt{\beta(t)} d\bar{\mathbf{W}},$$

which corresponds to the ordinary SDE

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x} dt + \sqrt{\beta(t)} d\mathbf{W},$$

which is the VP SDE.

7.3.2 Predictor-Corrector Samplers

- The **predictor-corrector sampler** is an attempt to cast the SMLD’s sampling algorithm into the paper’s framework of solving the reserve-time SDE.
- The idea is to combine a numerical SDE solver with score-based MCMC approaches such as the Langevin algorithm or HMC.
 - At each time step, the numerical SDE solver gives an estimate of the sample at the next time step. This is the “predictor” step.
 - After getting a sample from the predictor, we run several iterations of our score-based MCMC algorithm to make the sample conform more to the probability distribution at that time step. This is the “corrector” step.
- Here’s the pseudocode for the sampler.

```

Sample  $\mathbf{x}_n \sim p_n$  where  $p_n$  is a simple distribution.
for  $i \leftarrow n - 1$  to 1 do
    Predict  $\mathbf{x}_i$  from  $\mathbf{x}_{i+1}$  with a step from a numerical SDE solver.
    for  $j \leftarrow 1$  do  $m$  do
        Update  $\mathbf{x}_i$  with a score-based MCMC step.
    end for
end for

```

- The predictor-corrector sampler is a generalization of both the SMLD’s sampler and the DDPM’s sampler.
 - For the SMLD’s sampler, the predictor step simply sets $\mathbf{x}_i \leftarrow \mathbf{x}_{i+1}$. The corrector step applies the update step of the Langevin algorithm.
 - For the DDPM’s sampler, the predictor step is Equation 14. The corrector step does nothing and leaves \mathbf{x}_i unchanged.
- The paper also proposes new algorithms: the predictor-corrector samplers for the VE and VP SDEs. The pseudocodes are given below.

- Predictor-corrector sampler for the VE SDE.

```

Sample  $\mathbf{x}_n \sim \mathcal{N}(\mathbf{0}, \sigma_{\max}^2 I)$ .
for  $i \leftarrow n - 1$  to 1 do
    Sample  $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, I)$ .
     $\mathbf{x}_i \leftarrow \mathbf{x}_{i+1} + (\sigma_{i+1}^2 - \sigma_i^2)\mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, \sigma_{i+1}) + \sqrt{\sigma_{i+1}^2 - \sigma_i^2}\boldsymbol{\xi}$ 
    for  $j \leftarrow 1$  do  $m$  do

```

```

    Sample  $\xi \sim \mathcal{N}(\mathbf{0}, I)$ .
     $\mathbf{x}_i \leftarrow \mathbf{x}_i + \varepsilon_i \mathbf{s}_{\theta^*}(\mathbf{x}_i, \sigma_i) + \sqrt{2\varepsilon_i} \xi$ .
  end for
end for
– Predictor–corrector sampler for the VP SDE.
  Sample  $\mathbf{x}_n \sim \mathcal{N}(\mathbf{0}, I)$ .
  for  $i \leftarrow n - 1$  to 1 do
    Sample  $\xi \sim \mathcal{N}(\mathbf{0}, I)$ .
     $\mathbf{x}_i \leftarrow \frac{\mathbf{x}_{i+1} + \beta_{i+1} \mathbf{s}_{\theta^*}(\mathbf{x}_{i+1}, i+1)}{\sqrt{1 - \beta_{i+1}}} + \sqrt{\beta_{i+1}} \xi$ 
    for  $j \leftarrow 1$  to  $m$  do
      Sample  $\xi \sim \mathcal{N}(\mathbf{0}, I)$ .
       $\mathbf{x}_i \leftarrow \mathbf{x}_i + \varepsilon_i \mathbf{s}_{\theta^*}(\mathbf{x}_i, \sigma_i) + \sqrt{2\varepsilon_i} \xi$ .
    end for
  end for
end for

```

7.3.3 Neural ODE Approach

References

- [And82] Brian D.O. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- [DN21] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. *CoRR*, abs/2105.05233, 2021.
- [HJA20] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *CoRR*, abs/2006.11239, 2020.
- [Hyv05] Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(24):695–709, 2005.
- [KGV83] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220(4598):671–680, 1983.
- [Khu20] Pramook Khungurn. Generative diffusion models. <https://pkhungurn.github.io/notes/notes/ml/generative-diffusion-models/index.html>, 2020. Accessed: 2022-03-15.
- [Khu21] Pramook Khungurn. Introduction to conditional random fields. <https://pkhungurn.github.io/notes/notes/ml/conditional-random-fields-intro/index.html>, 2021. Accessed: 2022-02-07.
- [Khu22] Pramook Khungurn. A primer on markov chain monte carlo algorithms. <https://pkhungurn.github.io/notes/notes/scicomp/mcmc-primer/mcmc-primer.pdf>, 2022. Accessed: 2022-02-08.
- [KP13] Peter E. Kloeden and Eckhard Platen. *Numerical Solution of Stochastic Differential Equations*. Springer, 2013.
- [LMSR16] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian D. Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. *CoRR*, abs/1611.06612, 2016.
- [Nea98] Radford M. Neal. Annealed importance sampling, 1998.

- [SE19] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *CoRR*, abs/1907.05600, 2019.
- [SE20] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *CoRR*, abs/2006.09011, 2020.
- [Son21] Yang Song. Generative modeling by estimating gradients of the data distribution. <https://yang-song.github.io/blog/2021/score/>, 2021. Accessed: 2022-02-07.
- [SS19] Simo Särkkä and Arno Solin. *Applied Stochastic Differential Equations*. Cambridge University Press, 2019.
- [SSDK⁺21] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- [SWMG15] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. *CoRR*, abs/1503.03585, 2015.
- [Vin11] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural Computation*, 23(7):1661–1674, 2011.