

Diffusion models

SDE-based perspectives

Julius Berner

October 12, 2022

University of Vienna

Introduction

Task

Sample from a high-dimensional distribution Y_0 .

Sampling from high-dimensional distributions

Task

Sample from a high-dimensional distribution Y_0 .

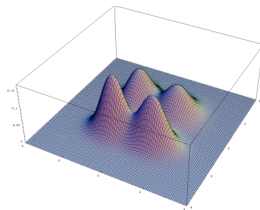
Y_0 can be given in the form of:

1. **samples** $Y_0^{(i)}$ (Y_0 (images, text, sound, ...)).



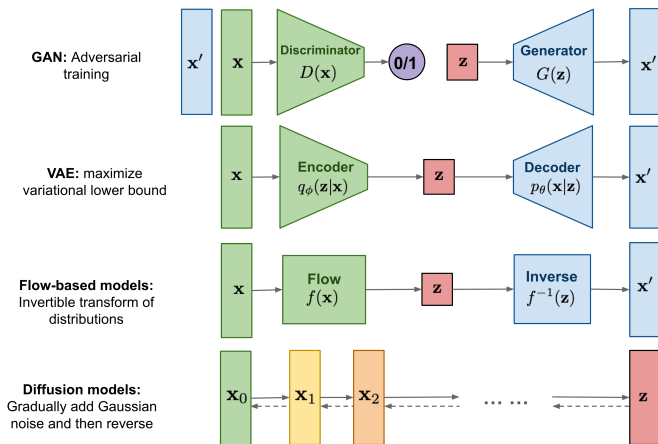
https://en.wikipedia.org/wiki/File:Cat_poster_1.jpg

2. an (unnormalized) **density** with $p_{Y_0} = \frac{1}{Z} e^{-\beta E}$ (e.g., in Bayesian statistics, computational physics and chemistry).



<https://en.wikipedia.org/wiki/File:Bimodal-bivariate-smal1.png>

Overview of generative models



<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

History: The development of diffusion models builds upon *(denoising) diffusion probabilistic modeling* [Sohl-Dickstein et al., 2015, Ho et al., 2020] and *score matching with Langevin dynamics* [Song and Ermon, 2019].

State-of-the art in generative modeling and likelihood estimation of high-dimensional image data [Nichol and Dhariwal, 2021, Kingma et al., 2021].

State-of-the art in generative modeling and likelihood estimation of high-dimensional image data [Nichol and Dhariwal, 2021, Kingma et al., 2021].

Figure 1: Sampling conditioned on the text prompt “a photograph of an astronaut riding a horse” using the stable diffusion model [Rombach et al., 2021].

Engineering Perspective

Diffusion process Y_t : Gradually add coordinate-wise Gaussian noise, i.e., conditioned on d -dimensional data Y_0 , we have that

$$Y_t = Y_0 + \int_0^t N; \quad N \sim N(0; I); \quad t \in [0; T]:$$

[Nichol and Dhariwal, 2021]

Diffusion process

Diffusion process Y_t : Gradually add coordinate-wise Gaussian noise, i.e., conditioned on d -dimensional data Y_0 , we have that

$$Y_t = Y_0 + \int_0^t N; \quad N \sim N(0; I); \quad t \in [0; T]:$$

[Nichol and Dhariwal, 2021]

Typical noise schedules for β_t and $\sigma_t = \frac{\beta_t - \beta_0}{\beta_1 - \beta_0}$:

[Nichol and Dhariwal, 2021]

Noise prediction objective (with batch-size n):

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \mathbb{E} \| \Phi(Y_t^{(i)}; t^{(i)}) - Y_0^{(i)} \|^2;$$

where Φ is typically a U-Net (with sinusoidal positional embeddings for t) and

- $Y_t^{(i)} = \sqrt{\alpha_t} Y_0^{(i)} + \sqrt{1 - \alpha_t} N^{(i)}$ (noisy image)
- $N^{(i)} \sim N(0; I)$ (standardized noise)
- $t^{(i)} \sim U([0; T])$ (time)
- $Y_0^{(i)} \sim Y_0$ (data)

are i.i.d. samples.

Noise prediction objective (with batch-size n):

$$\mathcal{L}(\theta) = \sum_{i=1}^n \mathbb{E} \left[\left\| N^{(i)} - \Phi(Y_t^{(i)}; t^{(i)}) \right\|^2 \right];$$

where Φ is typically a U-Net (with sinusoidal positional embeddings for t) and

- $Y_t^{(i)} = \sqrt{\alpha_t} Y_0^{(i)} + \sqrt{1 - \alpha_t} N^{(i)}$ (noisy image)
- $N^{(i)} \sim N(0; I)$ (standardized noise)
- $t^{(i)} \sim U([0; T])$ (time)
- $Y_0^{(i)} \sim Y_0$ (data)

are i.i.d. samples.

This is a **reparametrization of a denoising objective**, which works better in practice.

After training, we can approximately denoise Y_t as follows:

$$Y_0 \approx \frac{Y_t - \sqrt{1 - \alpha_t} \Phi(Y_t; t)}{\sqrt{\alpha_t}};$$

Architecture of typical U-Nets

Bayes' theorem yields the following formula for Y_s (conditioned on Y_0 and Y_t with $s < t$):

$$Y_s = \Theta_{t;s}(Y_0; Y_t; N); \quad N \sim N(0; 1);$$

where

$$\Theta_{t;s}(Y_0; Y_t; N) = \underbrace{\frac{\frac{2}{s} t}{\frac{2}{t} s} Y_t + s \frac{\frac{2}{t} \frac{2}{s}}{\frac{2}{s} t} Y_0}_{\text{mean}} + \frac{s}{\frac{2}{s} \frac{2}{s} \frac{4}{t}} \underbrace{\frac{2}{s} \frac{2}{t}}_{\text{standard deviation}} N;$$

Bayes' theorem yields the following formula for Y_s (conditioned on Y_0 and Y_t with $s < t$):

$$Y_s = \Theta_{t;s}(Y_0; Y_t; N); \quad N \sim N(0; 1);$$

where

$$\Theta_{t;s}(Y_0; Y_t; N) = \underbrace{\frac{\frac{2}{s} t}{\frac{2}{t} s} Y_t + \frac{s}{\frac{2}{t} s} \frac{\frac{2}{t} \frac{2}{s}}{\frac{2}{t}} Y_0}_{\text{mean}} + \underbrace{\frac{s}{\frac{2}{s}} \frac{\frac{2}{t} \frac{4}{s}}{\frac{2}{t} \frac{2}{s}} N}_{\text{standard deviation}}$$

Idea: Use the NN prediction for Y_0 and perform ancestral sampling.

1. Sample $X_T \sim N(0; 1)$ (approximately distributed as Y_T).
2. Iterate:

$$X_{t-1} := \Theta_{t;t-1} \left(\frac{X_t}{\underbrace{\frac{2}{t}}_{\text{denoising}}}; X_t; N^{(t)} \right)!$$

with i.i.d. $N^{(t)} \sim N(0; 1)$.

3. Output X_0 (approximately distributed as the data Y_0).

Bayes' theorem yields the following formula for Y_s (conditioned on Y_0 and Y_t with $s < t$):

$$Y_s = \Theta_{t;s}(Y_0; Y_t; N); \quad N \sim N(0; 1);$$

where

$$\Theta_{t;s}(Y_0; Y_t; N) = \underbrace{\frac{\frac{2}{s} t}{\frac{2}{t} s} Y_t + \frac{s}{\frac{2}{t} s} \frac{\frac{2}{t} s}{\frac{2}{t} s} Y_0}_{\text{mean}} + \underbrace{\frac{s}{\frac{2}{s}} \frac{\frac{2}{t} s}{\frac{2}{t} s} N}_{\text{standard deviation}}$$

Idea: Use the NN prediction for Y_0 and perform ancestral sampling.

1. Sample $X_T \sim N(0; 1)$ (approximately distributed as Y_T).
2. Iterate:

$$X_{t-1} := \Theta_{t;t-1} \left(\frac{X_t}{\underbrace{\frac{2}{t}}_{\text{denoising}}}; X_t; N^{(t)} \right)!$$

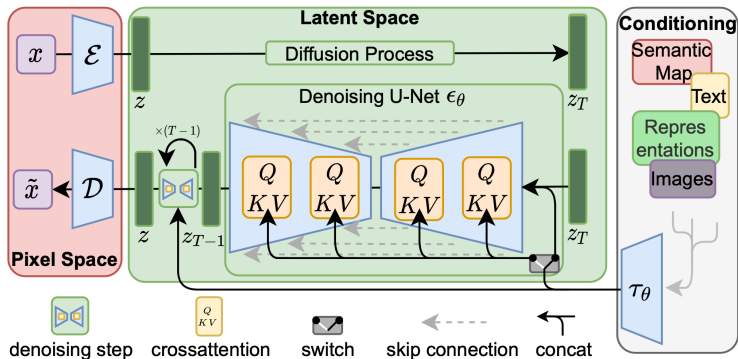
with i.i.d. $N^{(t)} \sim N(0; 1)$.

3. Output X_0 (approximately distributed as the data Y_0).

This can be viewed as variational auto-encoder with fixed encoder.

Stable diffusion model

Use diffusion in latent space of a pre-trained (regularized) auto-encoder and condition the U-Net on features given by a pre-trained domain-specific encoder (e.g., a transformer for text prompts):



[Rombach et al., 2021]

SDE-based perspective

Stochastic differential equations (SDEs)

Consider solutions Y to SDEs of the form

$$dY_s = \underbrace{\mu(Y_s)}_{\text{drift}} ds + \underbrace{\sigma(Y_s)}_{\text{diffusion}} dB_s;$$

where B_s is a d -dimensional Brownian motion.

Stochastic differential equations (SDEs)

Consider solutions Y to SDEs of the form

$$dY_s = \underbrace{\left(\frac{\partial Y_s}{\partial z} \right)}_{\text{drift}} ds + \underbrace{\left(\frac{\partial Y_s}{\partial z} \right)}_{\text{diffusion}} dB_s;$$

where B_s is a d -dimensional Brownian motion.

Intuition via Euler-Maruyama scheme $\hat{Y}_{t_{k+1}} \approx Y_{t_k}$:

$$\hat{Y}_{t_{k+1}} = \hat{Y}_{t_k} + \left(\frac{\partial \hat{Y}_{t_k}}{\partial z} \right) (t_{k+1} - t_k) + \left(\frac{\partial \hat{Y}_{t_k}}{\partial z} \right) \left(B_{t_{k+1}} - B_{t_k} \right);$$

$N(0; t_{k+1} - t_k)$

Figure 2: SDE solution and Euler-Maruyama scheme with $t_k = \frac{kT}{N}$ and $N = 4; 8$.

Consider a (time-varying) Ornstein–Uhlenbeck process

$$dY_s = \alpha(s)Y_s ds + \sigma(s)dB_s;$$

which diffuses the data Y_0 .

[Song et al., 2020]

Consider a (time-varying) Ornstein–Uhlenbeck process

$$dY_s = \theta(s)Y_s ds + \sigma(s)dB_s;$$

which diffuses the data Y_0 .

[Song et al., 2020]

Note that, conditioned on Y_0 , the solution Y_s is normally distributed. For the choices

$$\theta(s) = \frac{\theta(s)}{\sigma(s)} \quad \text{and} \quad \sigma^2(s) = 2 \int_0^s \theta(u) \sigma^2(u) du$$

we recover $p_{Y_s|Y_0}(j|Y_0) = N(\sigma(s)Y_0; \sigma^2(s))$.

We can reverse the diffusion (proven via the Fokker-Planck equation):

Reverse-time generative SDE/ODE [Anderson, 1982, Song et al., 2020]

The solutions to the SDE

$$dX_s = \left(r - \frac{1}{2} \sigma^2 \right) X_s ds + \sigma X_s dB_s; \quad X_0 = Y_T;$$

and the ODE

$$dX_s = \frac{1}{2} \sigma^2 X_s ds; \quad X_0 = Y_T;$$

both satisfy that $X_s = Y_{T-s}$, where $p_{Y_{T-s}}$ is the density of Y_{T-s} .

We can reverse the diffusion (proven via the Fokker-Planck equation):

Reverse-time generative SDE/ODE [Anderson, 1982, Song et al., 2020]

The solutions to the SDE

$$dX_s = \left[r \log p_{Y_T | s} \right] (X_s; s) ds + \sigma(s) dB_s; \quad X_0 = Y_T;$$

and the ODE

$$dX_s = \frac{1}{2} \left[r \log p_{Y_T | s} \right] (X_s; s) ds; \quad X_0 = Y_T;$$

both satisfy that $X_s \sim Y_{T-s}$, where $p_{Y_T | s}$ is the density of Y_{T-s} .

We need an approximation to the *score* $r \log p_{Y_T | s}$.

We can reverse the diffusion (proven via the Fokker-Planck equation):

Reverse-time generative SDE/ODE [Anderson, 1982, Song et al., 2020]

The solutions to the SDE

$$dX_s = \left[r \log p_{Y_T | s} - \frac{1}{2} \sigma^2(s) \right] ds + \sigma(s) dB_s; \quad X_0 = Y_T;$$

and the ODE

$$dX_s = \frac{1}{2} \left[r \log p_{Y_T | s} - \sigma^2(s) \right] ds; \quad X_0 = Y_T;$$

both satisfy that $X_s = Y_{T-s}$, where $p_{Y_T | s}$ is the density of Y_{T-s} .

We need an approximation to the *score* $r \log p_{Y_T | s}$.

Using the noise prediction network Φ , we obtain that

$$r \log p_{Y_t | Y_0}(Y_t | Y_0) = \frac{Y_t - Y_0}{\frac{2}{t}} = \frac{\Phi(Y_t; t)}{t}.$$

We can reverse the diffusion (proven via the Fokker-Planck equation):

Reverse-time generative SDE/ODE [Anderson, 1982, Song et al., 2020]

The solutions to the SDE

$$dX_s = \left[r \log p_{Y_T | s}(X_s; s) \right] ds + \sigma(s) dB_s; \quad X_0 = Y_T;$$

and the ODE

$$dX_s = \frac{1}{2} \left[r \log p_{Y_T | s}(X_s; s) \right] ds; \quad X_0 = Y_T;$$

both satisfy that $X_s \sim Y_{T-s}$, where $p_{Y_T | s}$ is the density of Y_{T-s} .

We need an approximation to the *score* $r \log p_{Y_T | s}$.

Using the noise prediction network Φ , we obtain that

$$r \log p_{Y_t | Y_0}(Y_t | Y_0) = \frac{Y_t - Y_0}{t} \frac{\Phi(Y_t; t)}{t}.$$

Sampling:

1. Sample $X_0 \sim N(0; 1)$.
2. Plug-in the approximate score and simulate the SDE (using Euler-Maruyama) or the ODE (analogous to time-continuous normalizing flows) to obtain samples X_T .

Variational Lower Bound

Up to a constant and a time-dependent weighting, the (negative) **noise prediction objective** also provides a **lower bound on the log-likelihood** $\log p_{X_T}(Y_0)$ of our model X (with score replaced by the NN approximation).

Variational Lower Bound

Up to a constant and a time-dependent weighting, the (negative) **noise prediction objective** also provides a **lower bound on the log-likelihood** $\log p_{X_T}(Y_0)$ of our model X (with score replaced by the NN approximation).

Proof idea with short-hands $f(t) := f(T - t)$, $D = \frac{1}{2} \nabla \cdot \nabla$, and $X = X$:

1. Fokker-Planck for p_X :

$$\partial_t p_X = \text{div} \cdot \text{div} D p_X - p_X$$

Variational Lower Bound

Up to a constant and a time-dependent weighting, the (negative) **noise prediction objective** also provides a **lower bound on the log-likelihood** $-\log p_{X_T}(Y_0)$ of our model X (with score replaced by the NN approximation).

Proof idea with short-hands $f(t) := f(T - t)$, $D = \frac{1}{2} \sigma^2$, and $X = X$:

1. Fokker-Planck for p_X :

$$\partial_t p_X = \text{div}(\text{div} D p_X) - p_X$$

2. Kolmogorov backwards equation for p_X :

$$\partial_t p_X = \text{tr}(D r^2 p_X) + r p_X + \text{div}(\cdot) p_X:$$

Variational Lower Bound

Up to a constant and a time-dependent weighting, the (negative) **noise prediction objective** also provides a **lower bound on the log-likelihood** $\log p_{X_T}(Y_0)$ of our model X (with score replaced by the NN approximation).

Proof idea with short-hands $f(t) := f(T - t)$, $D = \frac{1}{2} \sigma^2$, and $X = X$:

1. Fokker-Planck for p_X :

$$\partial_t p_X = \text{div}(\text{div} D p_X) - p_X$$

2. Kolmogorov backwards equation for p_X :

$$\partial_t p_X = \text{tr}(D r^2 p_X) + r p_X + \text{div}(\cdot) p_X$$

3. HJB equation for $V := -\log p_X$ (Hopf-Cole transformation):

$$\partial_t V = \text{tr}(D r^2 V) + r V - \text{div}(\cdot) + \frac{1}{2} \sigma^2 r V^2; \quad V(\cdot; T) = -\log p_{X_0}$$

Variational Lower Bound

Up to a constant and a time-dependent weighting, the (negative) **noise prediction objective** also provides a **lower bound on the log-likelihood** $\log p_{X_T}(Y_0)$ of our model X (with score replaced by the NN approximation).

Proof idea with short-hands $f(t) := f(T - t)$, $D = \frac{1}{2} \nabla^2$, and $X = X$:

1. Fokker-Planck for p_X :

$$\partial_t p_X = \text{div}(\text{div} D p_X) - p_X$$

2. Kolmogorov backwards equation for p_X :

$$\partial_t p_X = \text{tr}(D r^2 p_X) + r p_X + \text{div}(\cdot) p_X$$

3. HJB equation for $V := -\log p_X$ (Hopf-Cole transformation):

$$\partial_t V = \text{tr}(D r^2 V) + r V - \text{div}(\cdot) + \frac{1}{2} r^2 V^2; \quad V(\cdot; T) = -\log p_{X_0}$$

4. Reparametrize and use verification theorem from optimal control:

$$\log p_{X_T}(Y_0) = \int_0^T \text{div}(\Phi) + \frac{1}{2} k \Phi^2(Y_s; s) ds + \log p_{X_0}(Y_T)$$

Variational Lower Bound

Up to a constant and a time-dependent weighting, the (negative) **noise prediction objective** also provides a **lower bound on the log-likelihood** $\log p_{X_T}(Y_0)$ of our model X (with score replaced by the NN approximation).

Proof idea with short-hands $f(t) := f(T - t)$, $D = \frac{1}{2} \nabla^2$, and $X = X$:

1. Fokker-Planck for p_X :

$$\partial_t p_X = \text{div}(-\text{div} D p_X - r p_X)$$

2. Kolmogorov backwards equation for p_X :

$$\partial_t p_X = -\text{tr} D r^2 p_X + r p_X + \text{div}(\Phi) p_X$$

3. HJB equation for $V := -\log p_X$ (Hopf-Cole transformation):

$$\partial_t V = -\text{tr} D r^2 V + r V - \text{div}(\Phi) + \frac{1}{2} \nabla^2 r V^2; \quad V(\cdot; T) = -\log p_{X_0}$$

4. Reparametrize and use verification theorem from optimal control:

$$\log p_{X_T}(Y_0) = \int_0^T \text{div}(\Phi) ds - \frac{1}{2} \int_0^T k \Phi^2 ds + \log p_{X_0}(Y_T)$$

5. Employ Stokes' theorem to rewrite the divergence.

Anderson, B. D. (1982).

Reverse-time diffusion equation models.

Stochastic Processes and their Applications, 12(3):313–326.

Ho, J., Jain, A., and Abbeel, P. (2020).

Denoising diffusion probabilistic models.

Advances in Neural Information Processing Systems, 33:6840–6851.

Kingma, D., Salimans, T., Poole, B., and Ho, J. (2021).

Variational diffusion models.

Advances in Neural Information Processing Systems, 34:21696–21707.

Nichol, A. Q. and Dhariwal, P. (2021).

Improved denoising diffusion probabilistic models.

In *International Conference on Machine Learning*, pages 8162–8171. PMLR.

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. (2021).

High-resolution image synthesis with latent diffusion models.

Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. (2015).

Deep unsupervised learning using nonequilibrium thermodynamics.

In *International Conference on Machine Learning*, pages 2256–2265. PMLR.

Song, Y. and Ermon, S. (2019).

Generative modeling by estimating gradients of the data distribution.

Advances in Neural Information Processing Systems, 32.

Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. (2020).

Score-based generative modeling through stochastic differential equations.

In *International Conference on Learning Representations*.