

Deep Generative Models: Score Matching with Diffusion Models

Fall Semester 2025

René Vidal

Director of the Center for Innovation in Data Engineering and Science (IDEAS)

Rachleff University Professor, University of Pennsylvania

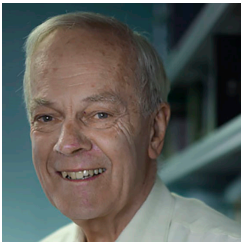
Amazon Scholar & Chief Scientist at NORCE



Diffusion Models

- **Derivation of Diffusion Models + Image Editing Applications (Last Lecture)**
 - Markov Hierarchical Variational Auto Encoders (MHVAE)
 - Diffusion Models are VAEs with Linear Gaussian Autoregressive latent space
 - ELBO for Diffusion Models is a particular case of ELBO for VAEs with extra structure
 - Implementation Details
 - Latent Diffusion Models (Stable Diffusion) + Controllable generation
- **A Different Viewpoint of Diffusion Models (Today's Lecture)**
 - Denoising Score Matching

Connection with DDPM and Score-based Models



Stochastic Processes and their Applications 12 (1982) 313–326
North-Holland Publishing Company

313

REVERSE-TIME DIFFUSION EQUATION MODELS*

Brian D.O. ANDERSON**
Department of Electrical Engineering, The University of Newcastle, N.S.W. 2308, Australia

Received 13 November 1979
Revised 12 May 1980

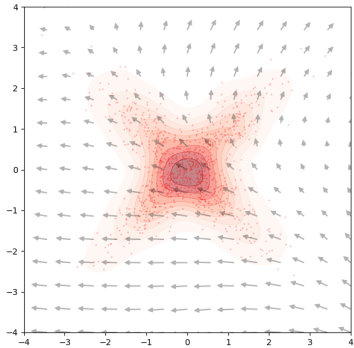
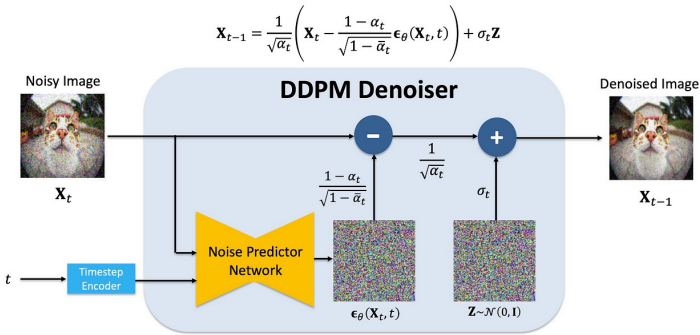
Reverse-time stochastic diffusion equation models are defined and it is shown how most processes defined via a forward-time or conventional diffusion equation model have an associated reverse-time model.

Stochastic equations	diffusion equations
reverse-time equations	Kolmogorov equations
Fokker-Planck equations	

SDE Perspective
(1982)

DDPM
(2020)

Score-based models
(2020)



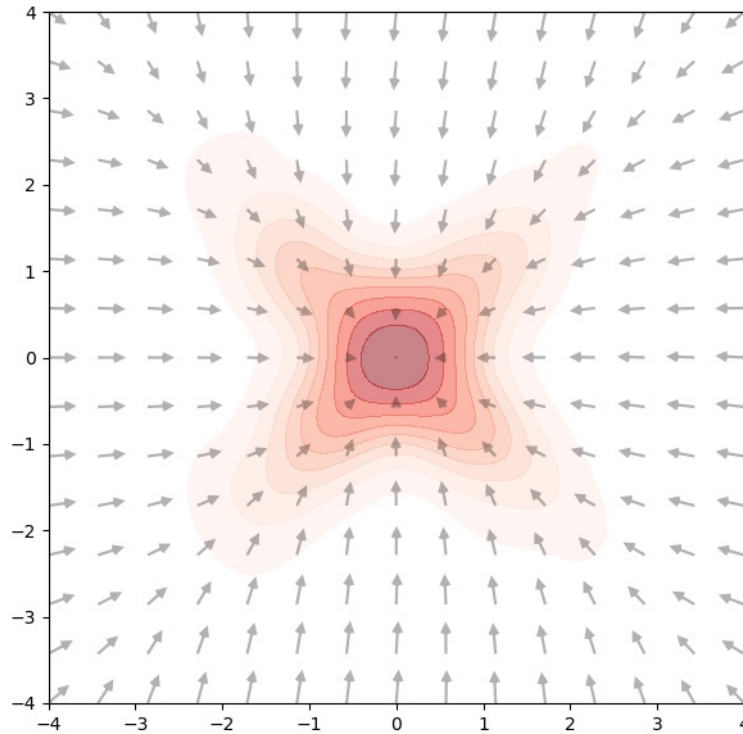
Generative Modelling via Score Functions

- **Goal:** Draw samples from $p(x)$
 - Rich history of sampling methods (MCMC, particle filter, etc.) – we won't cover these
- **So far:** Learning $p(x)$ by assuming latent variable model
 - Recall PPCA, GMM, VAE, HMM, LDS, DDPM
 - Sampling is easy, since one can sample latent variables z and then x from $p(x|z)$
- **Today:** Instead of learning $p(x)$ directly, we can turn to learn $\nabla_x \log p(x)$
 - This quantity is known as the *score function*
- **Challenges**
 - **Sampling:** Even if we had access to score function, how do we draw samples from $p(x)$?
 - **Score Estimation:** How can we estimate score function from data?
 - **Diffusion:** What does this have to do with diffusion models?

Sampling: Intuition for Score Function

- The score function $\nabla_x \log p(x)$ is a vector field that points to the direction of steepest increase in log-likelihood at any given point in data space

The arrows indicate the score function vector field – they point towards the mode



Data distribution on \mathbb{R}^2
Darker colors indicate higher probability density

Sampling: Naïve Idea

- **Simple sampling strategy**: start at any point in data space and take some steps in the direction of the score function (gradient ascent)

$$x_{t+1} = x_t + \eta \nabla_x \log p(x) \Big|_{x=x_t}$$

- If $\eta \rightarrow 0$, this process is described by the Ordinary Differential Equation (ODE)

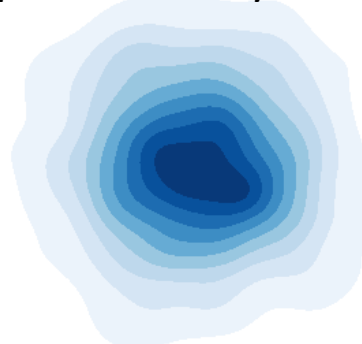
$$dx = \nabla_x \log p(x) dt$$

- However, the trajectories of this ODE converge to samples from only the modes of the distribution: i.e., **the ODE does not explore lower probability regions**

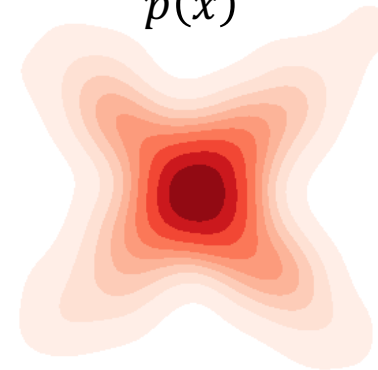
Different Samples



Empirical Density of Samples



$p(x)$



Sampling: Insights from Physics

- French physicist **Paul Langevin** proved exploration can be done by adding noise

$$dx = \nabla_x \log p(x)dt + \sqrt{2}B_t$$

- This gives a Stochastic Differential Equation (SDE), where the **Brownian motion** B_t is Gaussian noise with infinitesimally small variance:

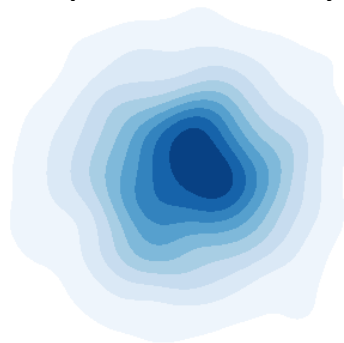
$$dB_t = \sqrt{dt} z, z \sim N(0, I)$$

- Langevin proved that, as $t \rightarrow \infty$, the solutions to the SDE visit each x with probability $p(x)$. That is, this process draws samples from $p(x)$!

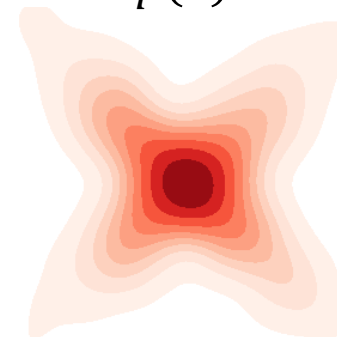
Different Samples



Empirical Density



$p(x)$



Langevin Dynamics Convergence Theorem

- **Theorem:** Consider the Langevin dynamics

$$dx = \nabla_x \log p(x) dt + \sqrt{2} B_t$$

where B_t is Brownian motion

$$dB_t = \sqrt{dt} z, z \sim N(0, I)$$

The process $(X_t)_{t \geq 0}$ defines a **Markov diffusion** that admits $p(x)$ as a **unique invariant distribution**. Moreover, if p_t denotes the distribution of X_t , then

$$p_t \xrightarrow[t \rightarrow \infty]{} p,$$

i.e., the distribution of X_t converges to p as $t \rightarrow \infty$.

Sampling as a Reverse Process

- Langevin dynamics **resemble a reverse process**: start with samples from a simple distribution such as $N(0, I)$ and **solve Langevin equation to get sample from $p(x)$**

$$dx_t = \nabla_x \log p(x) \Big|_{x=x_t} dt + \sqrt{2} B_t$$

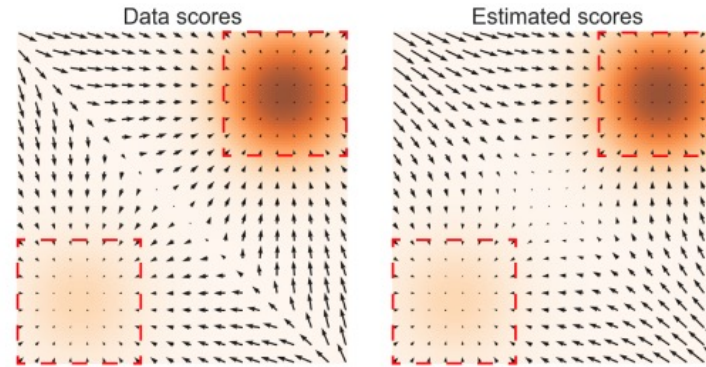
- However, we don't know the distribution $p(x)$, hence **we don't know the score function** $\nabla_x \log p(x)$. All we have are samples from $p(x)$.
- How can we estimate the score function $\nabla_x \log p(x)$ from samples $p(x)$? Can we train a deep network that predicts the score by minimizing a reconstruction loss?

$$\operatorname{argmin}_{\theta} \frac{1}{2} E_{x \sim p(x)} \left[\|s_{\theta}(x) - \nabla_x \log p(x)\|_2^2 \right]$$

Challenges with Learning Score Function

$$\operatorname{argmin}_{\theta} \frac{1}{2} E_{x \sim p(x)} \left[\|s_{\theta}(x) - \nabla_x \log p(x)\|_2^2 \right]$$

- **Challenge 1:** Inaccurate score estimation in low data density regions



- **Solution:** Perturb data with various levels of Gaussian noise, which “fills” in low density regions of the space
- **Challenge 2:** How do we optimize this because we don’t have ground truth score data to train on?
 - **Solution:** Clever manipulation of score function that we will see later

Perturbing Data: Forward Process

- We can arrive at the forward process mathematically by simply **inverting the Langevin equation**
- Instead of starting at $N(0, I)$ and drawing samples from $p(x)$, what if we just start from $p(x)$, and follow score of standard Gaussian, and end at $N(0, I)$?

$$\begin{aligned} dx &= \nabla_x \log N(0, I) dt + \sqrt{2} dB_t \\ &= -x dt + \sqrt{2} dB_t \end{aligned}$$

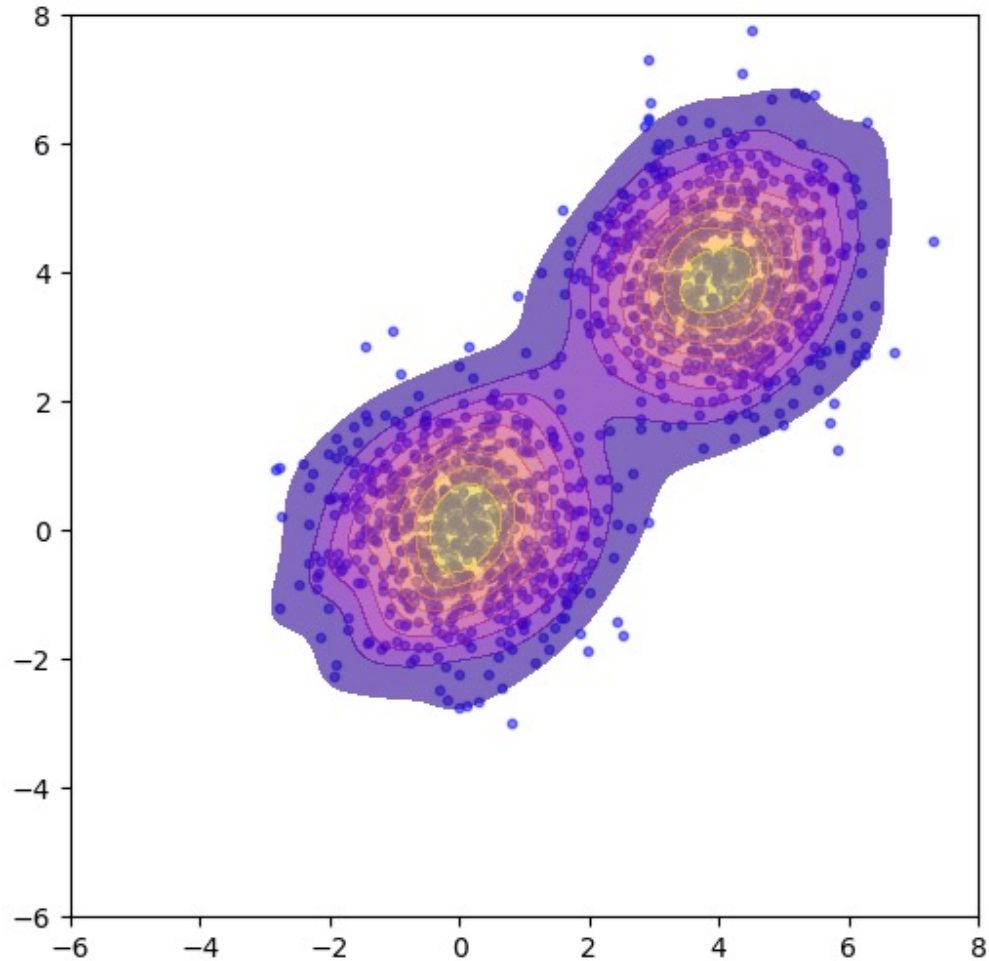
- Recall that $dB_t = \sqrt{dt} z$, $z \sim N(0, I)$. After discretizing we get:

$$x_{t+dt} = (1 - dt)x_t + \sqrt{2dt} z$$

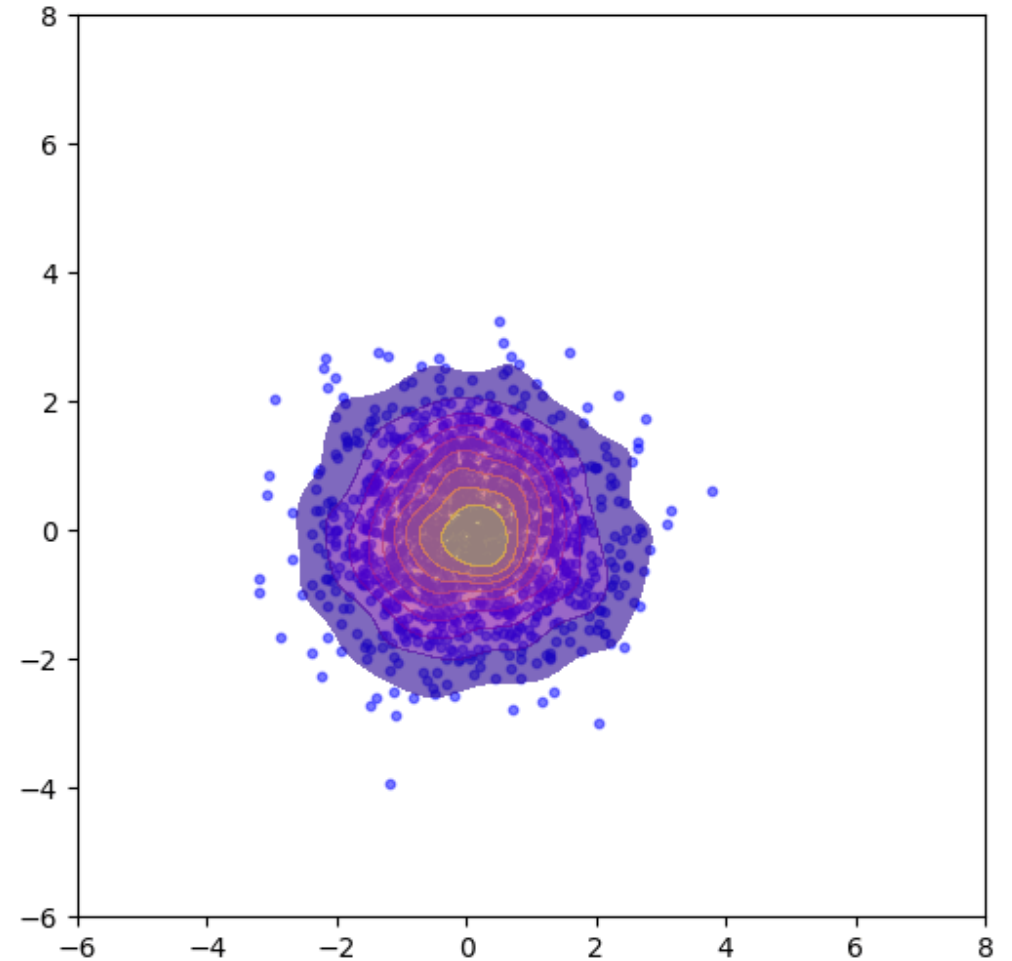
- Looks like the forward process we already know!

Back and forth from GMM to Normal distribution

Forward process



Reverse process



Summary so Far

- We have shown that with the score function given, we can draw samples from the underlying data distribution using the Langevin equation
- We described the forward process to get a learning objective (for all t)

$$\operatorname{argmin}_{\theta} \frac{1}{2} E_{x \sim p(x)} \left[\left\| s_{\theta}(x, t) - \nabla_{x_t} \log p(x_t) \right\|_2^2 \right]$$

- Next: How to actually learn this without ground truth access to the score function?

Rewriting the Score Function

- From the forward process, we generated data

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(\epsilon; 0, I)$$

- Let's cleverly rewrite the score function given this

$$\begin{aligned} \nabla_{x_t} \log p(x_t) &= \frac{\nabla_{x_t} p(x_t)}{p(x_t)} \\ &= \frac{1}{p(x_t)} \int \nabla_{x_t} p(x_t | x_0) p(x_0) dx_0 \\ &= \frac{1}{p(x_t)} \int p(x_t | x_0) \nabla_{x_t} \log p(x_t | x_0) p(x_0) dx_0 \\ &= \int \nabla_{x_t} \log p(x_t | x_0) p(x_0 | x_t) dx_0 \\ &= \int \frac{\sqrt{\bar{\alpha}_t} x_0 - x_t}{1 - \bar{\alpha}_t} p(x_0 | x_t) dx_0 \\ &= \frac{\sqrt{\bar{\alpha}_t} E[x_0 | x_t] - x_t}{1 - \bar{\alpha}_t} \end{aligned}$$

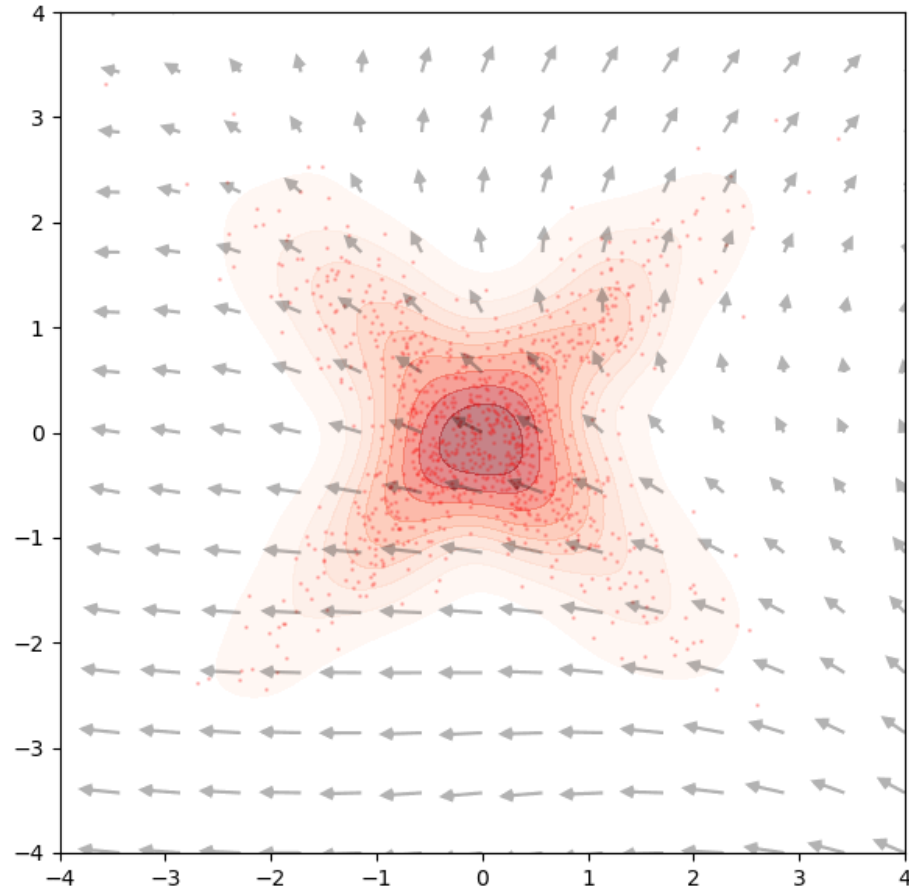
Denoising Score Matching

- Plugging into our objective before and simplifying, we have

$$\begin{aligned} & \operatorname{argmin}_{\theta} \frac{1}{2} E_{x_t \sim p(x_t)} \left[\left\| s_{\theta}(x, t) - \frac{\sqrt{\bar{\alpha}_t} E[x_0 | x_t] - x_t}{1 - \bar{\alpha}_t} \right\|_2^2 \right] \\ & \operatorname{argmin}_{\theta} \frac{1}{2} E_{x_t \sim p(x_t), x_0 \sim p(x_0 | x_t)} \left[\left\| s_{\theta}(x, t) - \frac{\sqrt{\bar{\alpha}_t} x_0 - x_t}{1 - \bar{\alpha}_t} \right\|_2^2 \right] \\ & \operatorname{argmin}_{\theta} \frac{1}{2(1 - \bar{\alpha}_t)} E_{x_t, x_0 \sim p(x_t, x_0)} \left[\left\| \epsilon_{\theta}(x, t) - \epsilon \right\|_2^2 \right] \end{aligned}$$

- We have a way to estimate the score without ever needing ground truth score, which is remarkable
- Moreover, this is exactly our denoising objective from diffusion models!
 - Maximizing our ELBO is equivalent to learning the score function at x_t !

Learning Score Function Example



Why is this viewpoint useful?

- It unifies diffusion model forward process as following Langevin equation to transform data distribution into Gaussian distribution
 - Different noise schedules just correspond to how fast this transition is
- Connects learning denoisers with estimating the score function, which we can then use to sample from the distribution
- Allows us to naturally think of continuous-time diffusion models
- Our rewriting of the score function is called Tweedie's formula
 - Can use it to solve $E[x_0|x_t]$ given the score function – one step denoising
 - Precisely what is used in DDIM to predict the x_0
- Conditional Sampling becomes straightforward (see next slide)

Conditional Diffusion Models - Guidance

- **Conditional Sampling:** Sample from $p(x | y)$ where y is another image, text, etc.
- The score function for this distribution can be obtained via Bayes rule
$$\nabla_{x_t} \log p(x_t | y) = \nabla_{x_t} \log p(y | x_t) + \nabla_{x_t} \log p(x_t)$$
- The second term is already modelled by a unconditional diffusion model!
- For the first term, we could use a classifier. This is known as classifier guidance, where we scale the effect of classifier by some factor γ in front of the first term



Samples from an unconditional diffusion model with classifier guidance, for guidance scales 1.0 (left) and 10.0 (right), taken from Dhariwal & Nichol (2021).

Classifier Free Guidance

- Classifier Guidance required access to a classifier that can handle noisy inputs
- Classifier-free Guidance allows us to do guidance without training any auxiliary classifiers

$$\begin{aligned}\nabla_{x_t} \log p_\gamma(x_t|y) &= \gamma \nabla_{x_t} \log p(y|x_t) + \nabla_{x_t} \log p(x_t) \\ &= \gamma (\nabla_{x_t} \log p(x_t|y) - \nabla_{x_t} \log p(x_t)) + \nabla_{x_t} \log p(x_t) \\ &= (1 - \gamma) \nabla_{x_t} \log p(x_t) + \gamma \nabla_{x_t} \log p(x_t|y)\end{aligned}$$

- $\gamma > 1$ is found to be the regime that works well
- Train a diffusion model with conditioning dropout: some percentage of the time (10-20%), just remove the conditioning information, so the model learns to be both a conditional and unconditional model

Classifier Free Guidance Results

Seems to work better than classifier guidance because the “classifier” is constructed from the generative model itself



Two sets of samples from OpenAI's GLIDE model, for the prompt 'A stained glass window of a panda eating bamboo.', taken from their paper. Guidance scale 1 (no guidance) on the left, guidance scale 3 on the right.



Two sets of samples from OpenAI's GLIDE model, for the prompt "A cozy living room with a painting of a corgi on the wall above a couch and a round coffee table in front of a couch and a vase of flowers on a coffee table.", taken from their paper. Guidance scale 1 (no guidance) on the left, guidance scale 3 on the right.

Conclusion

- We saw how diffusion models implicitly learn the score function and perform reverse process sampling using the Langevin equation
 - Connects to many core topics in stochastic differential equations and clean mathematical derivations
- Used the score function viewpoint to easily derive guidance (classifier and classifier-free) for conditional sampling from diffusion models
- Three great resources for more information
 - ICLR Blogpost – Building Diffusion Models Theory From Scratch (<https://iclr-blogposts.github.io/2024/blog/diffusion-theory-from-scratch/>)
 - Sander Dieleman Blog – Guidance: A Cheat Code for Diffusion Models (<https://sander.ai/2022/05/26/guidance.html>)
 - See references in the blog posts

Deep Generative Models: Diffusion Models – SDE perspective

Fall Semester 2025

René Vidal

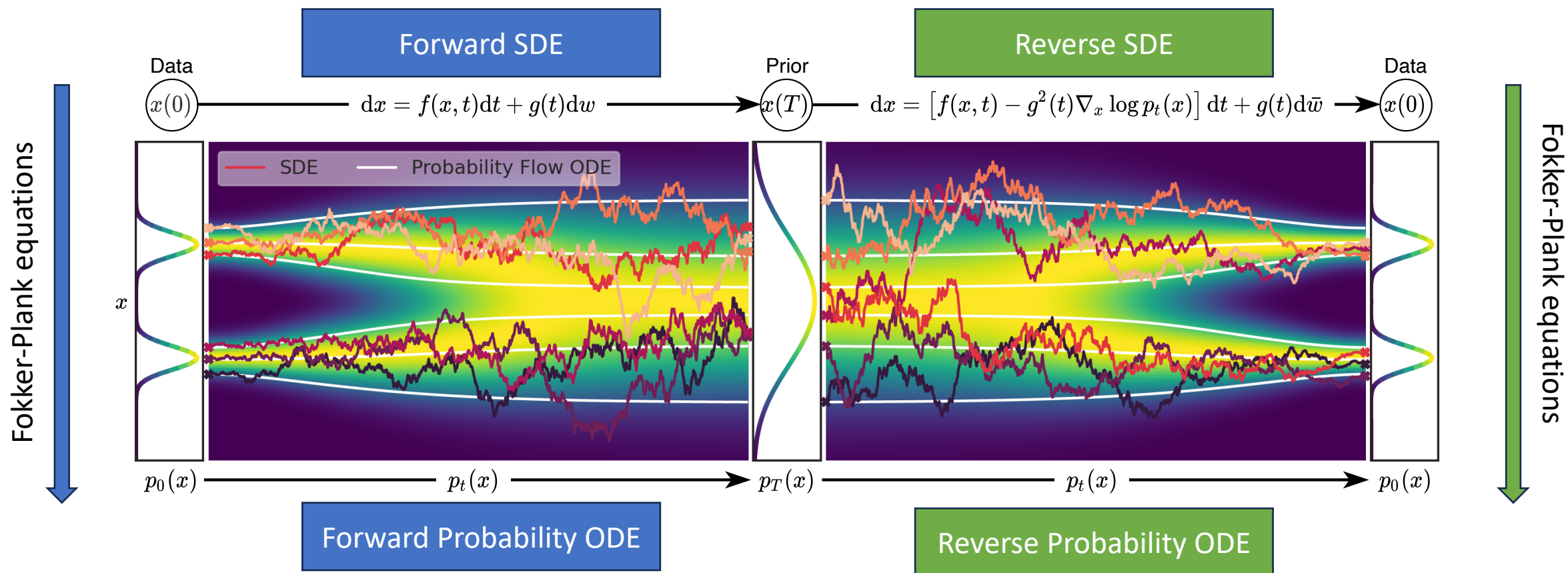
Director of the Center for Innovation in Data Engineering and Science (IDEAS)

Rachleff University Professor, University of Pennsylvania

Amazon Scholar & Chief Scientist at NORCE



Stochastic Differential Equation (SDE) view



Stochastic Differential Equation

- SDE is a differential equation in which one or more terms are stochastic processes.

$$dx_t = \mu(x_t, t)dt + \sigma(x_t, t)dw_t$$

- $\mu(x_t, t)$ is called as the **drift coefficient** and $\sigma(x_t, t)$ is called the **diffusion coefficient**.
- w_t is called the Wiener process and has the following properties:
 1. Initial zero condition: $w_0 = 0$ almost surely.
 2. Independent increments: For every $t, u \geq 0$, the increments $w_{t+u} - w_t$ are independent of w_s for $s < t$.
 3. Gaussian increments: For every $t, u \geq 0$, the increments $w_{t+u} - w_t \sim \mathcal{N}(0, uI)$.
 4. Continuity paths: w_t is continuous in t almost surely.
- How to obtain x_t given x_0 ?

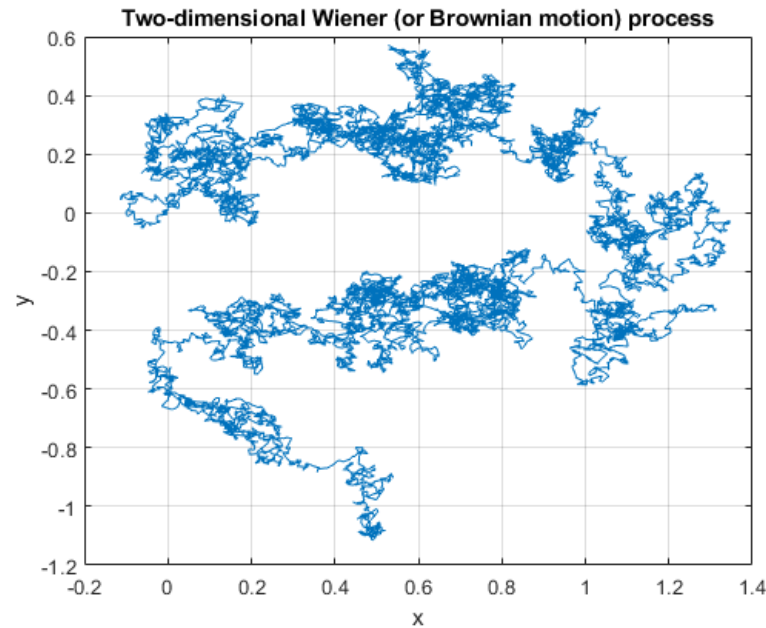
$$x_t - x_0 = \underbrace{\int_0^t \mu(x_t, t)dt}_{\text{Lebesgue integral}} + \underbrace{\int_0^t \sigma(x_t, t)dw_t}_{\text{It\^o integral}}$$

μ, σ should be Lipschitz continuous in (x, t) for the existence and uniqueness of SDE solution.

Example of SDEs

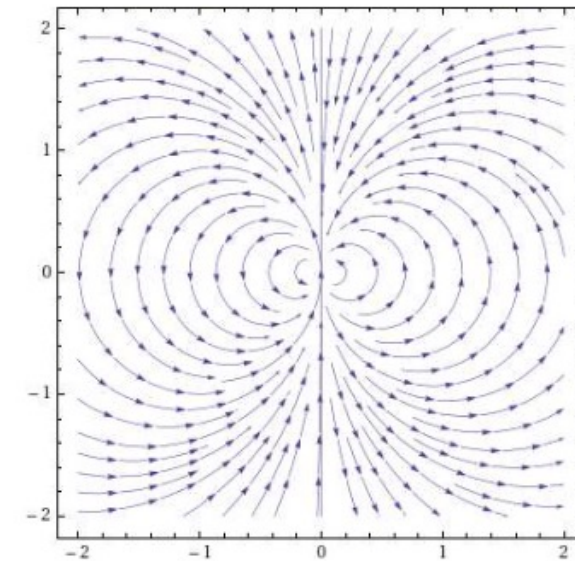
- Two distinct cases for the **drift** and **diffusion** coefficients in the process

$$dx_t = \mu(x_t, t)dt + \sigma(x_t, t)dw_t$$



$$dx_t = dw_t$$

SDE with **zero drift**
the process is purely stochastic.



$$dx_t = \mu(x_t, t)dt$$

SDE with **zero diffusion**
the process is purely deterministic.

Reversal of the diffusion process

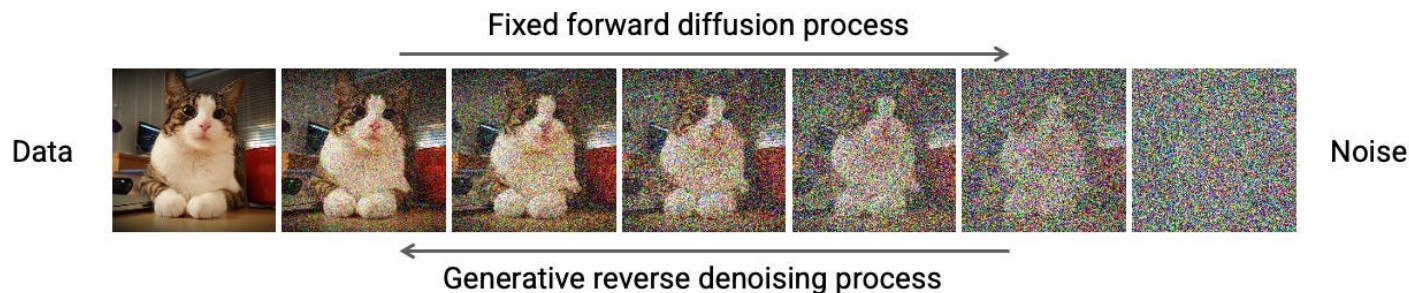
- [Brian Anderson](#) showed this in 1982 which later became a building block in [score-matching diffusion models](#) introduced in 2020.
- The forward SDE

$$dx_t = \mu(x_t, t)dt + \sigma(x_t, t)dw_t, \quad x_0 \sim p_0$$

has the time-reversal equation:

$$dx'_t = \left(\frac{\nabla_x(p(x'_{1-t}, 1-t) \cdot \sigma(x'_{1-t}, 1-t)^2)}{p(x'_{1-t}, 1-t)} - \mu(x'_{1-t}, 1-t) \right) dt + \sigma(x'_{1-t}, 1-t)d\bar{w}_t$$

with $x'_0 \sim p_1$, and (w_t, \bar{w}_t) being independent Wiener processes.



SDE \longrightarrow DDPM

$$\epsilon_i \sim \mathcal{N}(0, I)$$

- Set $\mu(x, t) = -\frac{1}{2}\sqrt{1 - \alpha_t} \cdot x_t$, and $\sigma(x, t) = N\sqrt{1 - \alpha_t}$ for some sufficiently large integer N .
- We have the forward process: $dx_t = -\frac{1}{2}\sqrt{1 - \alpha_t} \cdot x_t dt + N\sqrt{1 - \alpha_t} \cdot dw_t$.
- Now we discretize the time with N scales: $t = \frac{i}{N}$, and $\Delta t = \frac{1}{N}$.
- We can write $dx_t \approx x_{i+1} - x_i$,
$$dx_t \approx x_{i+1} - x_i = -\frac{1}{2}\sqrt{1 - \alpha_i} \cdot x_i \Delta t + \sqrt{1 - \alpha_t} \cdot \epsilon_i,$$
$$x_{i+1} = \left(1 - \frac{1}{2}\sqrt{1 - \alpha_i}\right) x_i + \sqrt{1 - \alpha_t} \cdot \epsilon_i$$
- Let $N \rightarrow \infty$, then we have $\left(1 - \frac{1}{2}\sqrt{1 - \alpha_i}\right) \approx \sqrt{1 - 1 + \alpha_i} = \sqrt{\alpha_i}$.
- Finally, we have $x_{i+1} \approx \sqrt{\alpha_i} x_i + \sqrt{1 - \alpha_i} \epsilon_i$, which is the same as DDPM.

$$\epsilon_t \sim \mathcal{N}(0, I)$$

DDPM \longrightarrow SDE

- The forward process in DDPM is given by

$$x_{t+\Delta t} = \sqrt{\alpha_t} x_t + \sqrt{1 - \alpha_t} \epsilon_t,$$

where t goes from 0 to 1, $\Delta t = \frac{1}{N}$, and N is the number of noise scales.

- Define $\tilde{a}(t) = (1 - \alpha_t)/\Delta t$.
- In the limit as $N \rightarrow \infty$, we have

$$x_{t+\Delta t} = \sqrt{1 - \tilde{a}(t)\Delta t} \cdot x_t + \sqrt{\tilde{a}(t) \cdot \Delta t} \cdot \epsilon_t,$$

$$x_{t+\Delta t} \approx \left(1 - \frac{\tilde{a}(t)\Delta t}{2}\right) x_t + \sqrt{\tilde{a}(t)} \cdot \sqrt{\Delta t} \cdot \epsilon_t,$$

$$x_{t+\Delta t} - x_t \approx -\frac{\tilde{a}(t)}{2} x_t \cdot \Delta t + \sqrt{\tilde{a}(t)} \cdot \sqrt{\Delta t} \cdot \epsilon_t,$$

$$dx_t \approx -\frac{\tilde{a}(t)}{2} x_t dt + \sqrt{\tilde{a}(t)} \cdot \sqrt{dt} \cdot \epsilon_t.$$

- Recall that $\sqrt{dt} \cdot \epsilon_t$ is a Weiner process.

SDE \longrightarrow Score-based models

- Let $\mu(x, t) = 0$, and $\sigma(x, t) = \sigma(t)$, then we have:

Forward: $dx_t = \sigma(t)dw_t; x_0 \sim p_0,$

Reverse: $dx'_t = \sigma(1-t)^2 \underbrace{\nabla_x \ln(p(x'_{1-t}, 1-t))}_{\text{Stein's score function}} dt + \sigma(1-t)d\bar{w}_t; x'_0 \sim p_1.$

- In the reverse process, we are forcing the samples to drift to **high**-likelihood regions just like in score-based diffusion model.
- After discretizing the reverse process, we obtain:

$$x'_{t+\Delta t} = x'_t + \sigma(1-t)^2 \Delta t \cdot \nabla_x \ln(p(x'_{1-t}, 1-t)) + \sigma(1-t) \sqrt{\Delta t} \epsilon_t$$

