

# Deep Generative Models

## Lecture 11

Roman Isachenko



AI Masters

2025, Spring

# Recap of previous lecture

## DDPM objective

$$\mathbb{E}_{\pi(\mathbf{x}_0)} \mathbb{E}_{t \sim U\{1, T\}} \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \left[ \frac{(1 - \alpha_t)^2}{2\tilde{\beta}_t \alpha_t} \left\| \mathbf{s}_{\theta, t}(\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0) \right\|_2^2 \right]$$

In practice the coefficient is omitted.

## NCSN objective

$$\mathbb{E}_{\pi(\mathbf{x}_0)} \mathbb{E}_{t \sim U\{1, T\}} \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \left\| \mathbf{s}_{\theta, \sigma_t}(\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0) \right\|_2^2$$

**Note:** The objective of DDPM and NCSN is almost identical. But the difference in sampling scheme:

- ▶ NCSN uses annealed Langevin dynamics;
- ▶ DDPM uses ancestral sampling.

# Recap of previous lecture

## Unconditional generation

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{1-\beta_t}} \cdot \mathbf{x}_t + \frac{\beta_t}{\sqrt{1-\beta_t}} \cdot \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\boldsymbol{\theta}) + \sigma_t \cdot \epsilon$$

## Conditional generation

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{1-\beta_t}} \cdot \mathbf{x}_t + \frac{\beta_t}{\sqrt{1-\beta_t}} \cdot \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}, \boldsymbol{\theta}) + \sigma_t \cdot \epsilon$$

## Conditional distribution

$$\begin{aligned}\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}, \boldsymbol{\theta}) &= \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\boldsymbol{\theta}) \\ &= \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) - \frac{\epsilon_{\theta,t}(\mathbf{x}_t)}{\sqrt{1-\bar{\alpha}_t}}\end{aligned}$$

Here  $p(\mathbf{y}|\mathbf{x}_t)$  – classifier on noisy samples (we have to learn it separately).

# Recap of previous lecture

## Classifier-corrected noise prediction

$$\epsilon_{\theta,t}(\mathbf{x}_t, \mathbf{y}) = \epsilon_{\theta,t}(\mathbf{x}_t) - \sqrt{1 - \bar{\alpha}_t} \cdot \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)$$

## Guidance scale

$$\epsilon_{\theta,t}(\mathbf{x}_t, \mathbf{y}) = \epsilon_{\theta,t}(\mathbf{x}_t) - \gamma \cdot \sqrt{1 - \bar{\alpha}_t} \cdot \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)$$

- ▶ Train DDPM as usual.
- ▶ Train the additional classifier  $p(\mathbf{y}|\mathbf{x}_t)$  on the noisy samples  $\mathbf{x}_t$ .

## Guided sampling

$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \cdot \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{\alpha_t(1 - \bar{\alpha}_t)}} \cdot \epsilon_{\theta,t}(\mathbf{x}_t, \mathbf{y}) + \sigma_t \cdot \epsilon$$

**Note:** Guidance scale  $\gamma$  tries to sharpen the distribution  $p(\mathbf{y}|\mathbf{x}_t)$  (in this case  $Z$  should not depend on  $\mathbf{x}_t$ ).

## Recap of previous lecture

- ▶ Previous method requires training the additional classifier model  $p(\mathbf{y}|\mathbf{x}_t)$  on the noisy data.
- ▶ Let try to avoid this requirement.

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}, \boldsymbol{\theta}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\boldsymbol{\theta})$$

$$\begin{aligned}\nabla_{\mathbf{x}_t}^\gamma \log p(\mathbf{x}_t|\mathbf{y}, \boldsymbol{\theta}) &= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\boldsymbol{\theta}) + \gamma \cdot \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) = \\ &= (1 - \gamma) \cdot \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\boldsymbol{\theta}) + \gamma \cdot \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}, \boldsymbol{\theta})\end{aligned}$$

## Classifier-free-corrected noise prediction

$$\hat{\epsilon}_{\boldsymbol{\theta},t}(\mathbf{x}_t, \mathbf{y}) = \gamma \cdot \epsilon_{\boldsymbol{\theta},t}(\mathbf{x}_t, \mathbf{y}) + (1 - \gamma) \cdot \epsilon_{\boldsymbol{\theta},t}(\mathbf{x}_t)$$

- ▶ Train the single model  $\epsilon_{\boldsymbol{\theta},t}(\mathbf{x}_t, \mathbf{y})$  on **supervised** data alternating with real conditioning  $\mathbf{y}$  and empty conditioning  $\mathbf{y} = \emptyset$ .
- ▶ Apply the model twice during inference.

# Outline

1. Continuity equation for NF log-likelihood
2. FFJORD (Hutchinson's trace estimator)
3. SDE basics

# Outline

1. Continuity equation for NF log-likelihood
2. FFJORD (Hutchinson's trace estimator)
3. SDE basics

# Continuous-in-time NF

## Theorem (continuity equation)

If  $\mathbf{f}$  is uniformly Lipschitz continuous in  $\mathbf{x}$  and continuous in  $t$ , then

$$\frac{d \log p_t(\mathbf{x}(t))}{dt} = -\text{tr} \left( \frac{\partial \mathbf{f}(\mathbf{x}(t), t)}{\partial \mathbf{x}(t)} \right)$$

It means that if we have the value  $\mathbf{x}_0 = \mathbf{x}(0)$  then the solution of the continuity equation will give us the density  $p_1(\mathbf{x}(1))$ .

## Solution of continuity equation

$$\log p_1(\mathbf{x}(1)) = \log p_0(\mathbf{x}(0)) - \int_0^1 \text{tr} \left( \frac{\partial \mathbf{f}(\mathbf{x}(t), t)}{\partial \mathbf{x}(t)} \right) dt.$$

**Note:** This solution will give us the density along the trajectory (not the total probability path).



# Continuous-in-time NF

## Forward transform + log-density

$$\mathbf{x}(1) = \mathbf{x}(0) + \int_0^1 \mathbf{f}_\theta(\mathbf{x}(t), t) dt$$

$$\log p_1(\mathbf{x}(1)|\theta) = \log p_0(\mathbf{x}(0)) - \int_0^1 \text{tr} \left( \frac{\partial \mathbf{f}_\theta(\mathbf{x}(t), t)}{\partial \mathbf{x}(t)} \right) dt$$

- ▶ **Discrete-in-time NF**: evaluation of determinant of the Jacobian costs  $O(m^3)$  (we need invertible  $\mathbf{f}$ ).
- ▶ **Continuous-in-time NF**: getting the trace of the Jacobian costs  $O(m^2)$  (we need smooth  $\mathbf{f}$ ).

## Why $O(m^2)$ ?

$\text{tr} \left( \frac{\partial \mathbf{f}_\theta(\mathbf{x}(t))}{\partial \mathbf{x}(t)} \right)$  costs  $O(m^2)$  ( $m$  evaluations of  $\mathbf{f}$ ), since we have to compute a derivative for each diagonal element. It is possible to reduce cost from  $O(m^2)$  to  $O(m)$ !

# Outline

1. Continuity equation for NF log-likelihood
2. FFJORD (Hutchinson's trace estimator)
3. SDE basics

# Continuous-in-time NF

## Hutchinson's trace estimator

If  $\epsilon \in \mathbb{R}^m$  is a random variable with  $\mathbb{E}[\epsilon] = 0$  and  $\text{Cov}(\epsilon) = \mathbf{I}$ , then

$$\begin{aligned}\text{tr}(\mathbf{A}) &= \text{tr}(\mathbf{A} \cdot \mathbf{I}) = \text{tr}\left(\mathbf{A} \cdot \mathbb{E}_{p(\epsilon)} \left[ \epsilon \epsilon^T \right]\right) = \\ &= \mathbb{E}_{p(\epsilon)} \left[ \text{tr}\left(\mathbf{A} \epsilon \epsilon^T\right) \right] = \mathbb{E}_{p(\epsilon)} \left[ \epsilon^T \mathbf{A} \epsilon \right]\end{aligned}$$

Jacobian vector products  $\mathbf{v}^T \frac{\partial \mathbf{f}}{\partial \mathbf{x}}$  can be computed for approximately the same cost as evaluating  $\mathbf{f}$  (`torch.func.jvp`).

## FFJORD density estimation

$$\begin{aligned}\log p_1(\mathbf{x}(1)) &= \log p_0(\mathbf{x}(0)) - \int_0^1 \text{tr} \left( \frac{\partial \mathbf{f}(\mathbf{x}(t), t)}{\partial \mathbf{x}(t)} \right) dt = \\ &= \log p_0(\mathbf{x}(0)) - \mathbb{E}_{p(\epsilon)} \int_0^1 \left[ \epsilon^T \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \epsilon \right] dt.\end{aligned}$$

# Outline

1. Continuity equation for NF log-likelihood
2. FFJORD (Hutchinson's trace estimator)
3. SDE basics

# Stochastic differential equation (SDE)

Let define stochastic process  $\mathbf{x}(t)$  with initial condition  $\mathbf{x}(0) \sim p_0(\mathbf{x}) = \pi(\mathbf{x})$ :

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

- ▶  $\mathbf{f}(\mathbf{x}, t) : \mathbb{R}^m \times [0, 1] \rightarrow \mathbb{R}^m$  is the **drift** function of  $\mathbf{x}(t)$ .
- ▶  $g(t) : \mathbb{R} \rightarrow \mathbb{R}$  is the **diffusion** function of  $\mathbf{x}(t)$ .
- ▶  $\mathbf{w}(t)$  is the standard Wiener process (Brownian motion):
  1.  $\mathbf{w}(0) = 0$  (almost surely);
  2.  $\mathbf{w}(t)$  has independent increments;
  3.  $\mathbf{w}(t) - \mathbf{w}(s) \sim \mathcal{N}(0, (t-s)\mathbf{I})$ , for  $t > s$ .
- ▶  $d\mathbf{w} = \mathbf{w}(t+dt) - \mathbf{w}(t) = \mathcal{N}(0, \mathbf{I} \cdot dt) = \epsilon \cdot \sqrt{dt}$ , where  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ .
- ▶ If  $g(t) = 0$  we get standard ODE.

# Stochastic differential equation (SDE)

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

- ▶ In contrast to ODE, initial condition  $\mathbf{x}(0)$  does not uniquely determine the process trajectory.
- ▶ We have two sources of randomness: initial distribution  $p_0(\mathbf{x})$  and Wiener process  $\mathbf{w}(t)$ .

## Discretization of SDE (Euler method) - SDESolve

$$\mathbf{x}(t + dt) = \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), t) \cdot dt + g(t) \cdot \epsilon \cdot \sqrt{dt}$$

If  $dt = 1$ , then

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \mathbf{f}(\mathbf{x}_t, t) + g(t) \cdot \epsilon$$

- ▶ At each moment  $t$  we have the density  $p_t(\mathbf{x}) = p(\mathbf{x}, t)$ .
- ▶  $p : \mathbb{R}^m \times [0, 1] \rightarrow \mathbb{R}_+$  is a **probability path** between  $p_0(\mathbf{x})$  and  $p_1(\mathbf{x})$ .
- ▶ How to get the distribution path  $p_t(\mathbf{x})$  for  $\mathbf{x}(t)$ ?

# Stochastic differential equation (SDE)

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w}, \quad d\mathbf{w} = \boldsymbol{\epsilon} \cdot \sqrt{dt}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}).$$

## Theorem (Kolmogorov-Fokker-Planck)

Evolution of the distribution  $p_t(\mathbf{x})$  is given by the following equation:

$$\frac{\partial p_t(\mathbf{x})}{\partial t} = -\text{div}(\mathbf{f}(\mathbf{x}, t)p_t(\mathbf{x})) + \frac{1}{2}g^2(t)\Delta_{\mathbf{x}}p_t(\mathbf{x})$$

Here

$$\text{div}(\mathbf{v}) = \sum_{i=1}^m \frac{\partial v_i(\mathbf{x})}{\partial x_i} = \text{tr} \left( \frac{\partial \mathbf{v}(\mathbf{x})}{\partial \mathbf{x}} \right)$$

$$\Delta_{\mathbf{x}}p_t(\mathbf{x}) = \sum_{i=1}^m \frac{\partial^2 p_t(\mathbf{x})}{\partial x_i^2} = \text{tr} \left( \frac{\partial^2 p_t(\mathbf{x})}{\partial \mathbf{x}^2} \right)$$

$$\frac{\partial p_t(\mathbf{x})}{\partial t} = \text{tr} \left( -\frac{\partial}{\partial \mathbf{x}} [\mathbf{f}(\mathbf{x}, t)p_t(\mathbf{x})] + \frac{1}{2}g^2(t) \frac{\partial^2 p_t(\mathbf{x})}{\partial \mathbf{x}^2} \right)$$

# Stochastic differential equation (SDE)

## Theorem (Kolmogorov-Fokker-Planck)

$$\frac{\partial p_t(\mathbf{x})}{\partial t} = \text{tr} \left( -\frac{\partial}{\partial \mathbf{x}} [\mathbf{f}(\mathbf{x}, t) p_t(\mathbf{x})] + \frac{1}{2} g^2(t) \frac{\partial^2 p_t(\mathbf{x})}{\partial \mathbf{x}^2} \right)$$

- ▶ KFP theorem does not define the SDE uniquely in general case.
- ▶ This is the generalization of continuity equation that we used in continuous-in-time NF:

$$\frac{d \log p_t(\mathbf{x}(t))}{dt} = -\text{tr} \left( \frac{\partial \mathbf{f}(\mathbf{x}, t)}{\partial \mathbf{x}} \right).$$

## Langevin SDE (special case)

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

$$d\mathbf{x} = \frac{1}{2} \frac{\partial}{\partial \mathbf{x}} \log p_t(\mathbf{x}) dt + \mathbf{1} \cdot d\mathbf{w}$$

Let apply KFP theorem to this SDE.



## Langevin SDE (special case)

$$d\mathbf{x} = \frac{1}{2} \frac{\partial}{\partial \mathbf{x}} \log p_t(\mathbf{x}) dt + 1 \cdot d\mathbf{w}$$

$$\begin{aligned} \frac{\partial p_t(\mathbf{x})}{\partial t} &= \text{tr} \left( -\frac{\partial}{\partial \mathbf{x}} \left[ p_t(\mathbf{x}) \frac{1}{2} \frac{\partial}{\partial \mathbf{x}} \log p_t(\mathbf{x}) \right] + \frac{1}{2} \frac{\partial^2 p_t(\mathbf{x})}{\partial \mathbf{x}^2} \right) = \\ &= \text{tr} \left( -\frac{\partial}{\partial \mathbf{x}} \left[ \frac{1}{2} \frac{\partial}{\partial \mathbf{x}} p_t(\mathbf{x}) \right] + \frac{1}{2} \frac{\partial^2 p_t(\mathbf{x})}{\partial \mathbf{x}^2} \right) = 0 \end{aligned}$$

The density  $p_t(\mathbf{x}) = \text{const}(t)$ !

If  $\mathbf{x}(0) \sim p_0(\mathbf{x})$ , then  $\mathbf{x}(t) \sim p_0(\mathbf{x})$ .

## Discretized Langevin SDE

$$\mathbf{x}_{t+1} - \mathbf{x}_t = \frac{\eta}{2} \cdot \frac{\partial}{\partial \mathbf{x}} \log p_t(\mathbf{x}) + \sqrt{\eta} \cdot \epsilon, \quad \eta \approx dt.$$

## Langevin dynamic

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\eta}{2} \cdot \nabla_{\mathbf{x}} \log p(\mathbf{x}|\theta) + \sqrt{\eta} \cdot \epsilon, \quad \eta \approx dt.$$

# Summary

- ▶ Continuity equation allows to calculate  $\log p(\mathbf{x}, t)$  at arbitrary moment  $t$ .
- ▶ FFJORD model makes such kind of NF scalable.
- ▶ Adjoint method are the continuous analog of backpropagation in the discrete time. Pontryagin theorem gives the way to compute the adjoint functions.
- ▶ Using numerical solvers it is possible to make forward and backward passes for the continuous-in-time NF.
- ▶ SDE defines a stochastic process with drift and diffusion terms. ODEs are the special case of SDEs.