Introduction to Deep Generative Modeling

Lecture #15

HY-673 – Computer Science Dep, University of Crete <u>Professors:</u> Yannis Pantazis, Yannis Stylianou <u>Teaching Assistant:</u> Thomas Marchioro

Taxonomy of GMs



Lecture #15 HY-673

Recap: Forward Diffusion Process

Lecture #15 HY-673

Data

The forward diffusion process:



Recap: Reverse Denoising Process

Lecture #15 HY-673

The formal definition of the reverse process in T steps:

Reverse Denoising Process (generative)

Data

Noise x_0 x_1 x_4 x_2 x_3 x_T • • • $p(x_T) = \mathcal{N}(x_T; 0, I_d)$ $p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^{T} p_{\theta}(x_{t-1}|x_t).$ $p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \sigma_t^2 I_d)$ $\approx q(x_{t-1}|x_t)$ (true posterior; intractable) Trainable network (U-net, Denoising Autoencoder)

Recap: Training and Sampling

Minimize a simplification of negative ELBO:

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{x_0 \sim p_d(x_0), \epsilon \sim \mathcal{N}(0, I_d), t \sim \mathcal{U}(1, \mathrm{T})} \left[\left| \left| \epsilon - \epsilon_\theta \left(\underbrace{\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon}_{x_t}, t \right) \right| \right|^2 \right]$$

Lecture #15

HY-673

Algorithm 1 Training	Algorithm 2 Sampling
1: repeat	1: $x_{\mathrm{T}} \sim \mathcal{N}(0, I_d)$
2: $x_0 \sim p_d(x_0)$ 2. $t \in \text{Uniform}(1, T)$	2: for $t = T,, 1$ do
3: $t \sim \text{Omorm}(1, \dots, 1)$ 4: $\epsilon \sim \mathcal{N}(0, I_d)$	3: $z \sim \mathcal{N}(0, I_d)$
5: Take gradient descent step on	4: $x_{t-1} = \frac{1}{\sqrt{1-\beta_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_{\theta}(x_t, t) \right) + \sigma_t z$
	5: end for
$\nabla_{\theta} \epsilon - \epsilon_{\theta} (\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t) ^2$	6: return x_0

6: until converged



Forward Diffusion Process as Stochastic Differential Equation

Consider the limit of many small steps:

Data



Lecture #15

HY-673

Stochastic Differential Equation (SDE) describing the diffusion process in the infinitesimal limit

https://arxiv.org/abs/2011.13456

Crash Course in Differential Equations

Ordinary Differential Equation (ODE):

 $\frac{dx}{dt} = f(x,t)$ or dx = f(x,t)dt



Lecture #15 HY-673

Crash Course in Differential Equations

Lecture #15 HY-673

Ordinary Differential Equation (ODE):

 $\frac{dx_t}{dt} = f(x_t, t) \text{ or } dx_t = f(x_t, t)dt$



Stochastic Differential Equation (SDE):

$$\underbrace{\frac{dx_t}{dt}}_{\text{drift coefficient diffusion coefficient}} = \underbrace{f(x_t, t)}_{\text{drift coefficient diffusion coefficient}} W'_t \quad \Leftarrow$$

$$\left(dx_t = f(x_t, t)dt + \sigma(x_t, t)dW_t\right)$$

Wiener Process Derivative (i.e., Gaussian White Noise)



 $x(t + \Delta t) \approx x(t) + f(x(t), t)\Delta t + \sigma(x(t), t)\sqrt{\Delta t}z_t$

Crash Course in Differential Equations

Lecture #15 HY-673

Ordinary Differential Equation (ODE):

 $\frac{dx_t}{dt} = f(x_t, t) \text{ or } dx_t = f(x_t, t)dt$



Stochastic Differential Equation (SDE):

$$\frac{dx_t}{dt} = \underbrace{f(x_t, t)}_{\text{drift coefficient diffusion coefficient}} + \underbrace{\sigma(x_t, t)}_{\text{drift coefficient}} W_t' \quad \Leftarrow$$

$$\left(dx_t = f(x_t, t)dt + \sigma(x_t, t)dW_t\right)$$

$$\left(dx_t = f(x_t, t)dt + \sigma(x_t, t)dW_t\right)$$



Wiener Process Derivative (i.e., Gaussian White Noise)



 $x(t + \Delta t) \approx x(t) + f(x(t), t)\Delta t + \sigma(x(t), t)\sqrt{\Delta t}z_t$

Forward Diffusion Process as Stochastic Differential Equation

Consider the limit of many small steps:

Data



Lecture #15

HY-673

Stochastic Differential Equation (SDE) describing the diffusion process in the infinitesimal limit

https://arxiv.org/abs/2011.13456

Forward Diffusion Process as Stochastic Differential Equation

Forward diffusion process (fixed) $i \to i$ $i \to i$ i

Forward Diffusion SDE:

https://arxiv.org/abs/2011.13456

$$dx_t = \underbrace{-\frac{1}{2}\beta(t)x_tdt}_{-\frac{1}{2}} + \underbrace{\sqrt{\beta(t)}dW_t}_{-\frac{1}{2}}.$$

drift term (pulls towards mode) diffusion term (injects noise) Lecture #15 HY-673

Forward Diffusion Process as Stochastic Differential Equation

Lecture #15 HY-673



Special case of more general SDEs used in generative diffusion models:

 $dx_t = f(t)x_t dt + g(t)dW_t.$

https://arxiv.org/abs/2011.13456

The Generative Reverse Stochastic Differential Equation

Lecture #15 HY-673



Forward Diffusion SDE:

$$dx_t = -\frac{1}{2}\beta(t)x_tdt + \sqrt{\beta(t)}dW_t$$

But what about the reverse process, necessary for generation?



The Generative Reverse Stochastic Differential Equation

Lecture #15 HY-673



Forward Diffusion SDE: Reverse Generative Diffusion SDE:



https://arxiv.org/abs/2011.13456 https://www.sciencedirect.com/science/article/pii/0304414982900515



Simulate reverse diffusion process: Data generation from random noise!

The Generative Reverse Stochastic Differential Equation

Lecture #15 HY-673



Score Matching

Lecture #15 HY-673



• Naive idea: learn a model for the score function by direct regression.



 \implies But $\nabla_{\mathbf{x}} \log q_t(x_t)$ (score of the marginal diffused density $q_t(x_t)$) is not tractable!

https://ieeexplore.ieee.org/document/ 6795935 https://arxiv.org/abs/1907.05600 https://arxiv.org/abs/2011.13456

"Variance Preserving" SDE: $dx_t = -\frac{1}{2}\beta(t)x_t dt + \sqrt{\beta(t)}dW_t$ $q_t(x_t|x_0) = \mathcal{N}(x_t;\gamma_t x_0, \sigma_t^2 I)$ $\gamma_t = \exp^{-\frac{1}{2}\int_0^t \beta(s)ds}$ $\sigma_t^2 = 1 - \exp^{-\frac{1}{2}\int_0^t \beta(s)ds}$

https://arxiv.org/abs/2011.13456

Lecture #15

HY-673

- Instead, diffuse individual data points x_0 . Conditional $q_t(x_t|x_0)$ is tractable!
- Denoising Score Matching:

$$\min_{\theta} \underbrace{\mathbb{E}_{t \sim \mathcal{U}(0,T)}}_{\text{diffusion time } t} \underbrace{\mathbb{E}_{x_0 \sim q_0(x_0)}}_{\text{data } x_0} \underbrace{\mathbb{E}_{x_t \sim q_t(x_t|x_0)}}_{\text{diffused data sample } x_t|x_0} || \underbrace{s_{\theta}(x_t, t)}_{\text{neural network}} - \underbrace{\nabla x \log q_t(x_t|x_0)}_{\text{score of diffused data sample}} ||_2^2.$$

$$\sup_{\theta \in \mathcal{U}(x_t, t)}_{\text{data } x_0} \underbrace{\mathbb{E}_{x_t \sim q_t(x_t|x_0)}}_{\text{diffused data sample } x_t|x_0} || \underbrace{s_{\theta}(x_t, t)}_{\text{neural network}} - \underbrace{\nabla x \log q_t(x_t|x_0)}_{\text{data sample}} ||_2^2.$$

$$\lim_{\theta \in \mathcal{U}(x_t, t)}_{\text{data } x_0} \underbrace{\mathbb{E}_{x_t \sim q_t(x_t|x_0)}}_{\text{diffused data sample } x_t|x_0} || \underbrace{s_{\theta}(x_t, t)}_{\text{neural network}} - \underbrace{\nabla x \log q_t(x_t|x_0)}_{\text{data sample}} ||_2^2.$$

 \Rightarrow After expectations, $s_{\theta}(x_t, t) \approx \nabla_x \log q_t(x_t)!$

 $\begin{array}{c} \text{Implementation 1: Noise Prediction} \\ \text{``Variance Preserving'' SDE:} \\ dx_t = -\frac{1}{2}\beta(t)x_t dt + \sqrt{\beta(t)}dW_t \\ dx_t = -\frac{1}{2}\beta(t)x_t dt + \sqrt{\beta(t)}dW_t \\ q_t(x_t|x_0) = \mathcal{N}(x_t;\gamma_t x_0,\sigma_t^2 I) \\ \gamma_t = \exp^{-\frac{1}{2}\int_0^t \beta(s)ds} \\ \sigma_t^2 = 1 - \exp^{-\frac{1}{2}\int_0^t \beta(s)ds} \end{array}$

Lecture #15

HY-673

• Denoising Score Matching: $\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{x_0 \sim q_0(x_0)} \mathbb{E}_{x_t \sim q_t(x_t|x_0)} ||s_{\theta}(x_t,t) - \nabla x \log q_t(x_t|x_0)||_2^2.$

• Re-parametrized Sampling: $x_t = \gamma_t x_0 + \sigma_t \epsilon, \ \epsilon \sim \mathcal{N}(0, I).$

• Score Function: $\nabla_x \log q_t(x_t | x_0) = -\nabla_x \frac{(x_t - \gamma_t x_0)^2}{2\sigma_t^2} = -\frac{x_t - \gamma_t x_0}{\sigma_t^2} = -\frac{\gamma_t x_0 + \sigma_t \epsilon - \gamma_t x_0}{\sigma_t^2} = -\frac{\epsilon}{\sigma_t}.$ • Neural Network Model: $\sigma_{\theta}(x_t, t) := -\frac{\epsilon_{\theta}(x_t, t)}{\sigma_t}.$

 $\begin{array}{c} \text{Implementation 1: Noise Prediction} \\ \text{Forward diffusion process (fixed)} \\ \text{Forward$

• Denoising Score Matching:

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{x_0 \sim q_0(x_0)} \mathbb{E}_{x_t \sim q_t(x_t|x_0)} ||s_{\theta}(x_t,t) - \nabla x \log q_t(x_t|x_0)||_2^2$$

$$\implies \min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{x_0 \sim q_0(x_0)} \mathbb{E}_{\epsilon \sim \mathcal{N}(0,I)} \frac{1}{\sigma_t^2} ||\epsilon - \epsilon_{\theta}(x_t,t)||_2^2.$$

https://ieeexplore.ieee.org/document/ 6795935 https://arxiv.org/abs/1907.05600 https://arxiv.org/abs/2011.13456

Lecture #15

HY-673

Lecture #15 HY-673



"Variance Preserving" SDE: $dx_t = -\frac{1}{2}\beta(t)x_t dt + \sqrt{\beta(t)}dW_t$ $q_t(x_t|x_0) = \mathcal{N}(x_t;\gamma_t x_0,\sigma_t^2 I)$ $\gamma_t = \exp^{-\frac{1}{2}\int_0^t \beta(s)ds}$ $\sigma_t^2 = 1 - \exp^{-\frac{1}{2}\int_0^t \beta(s)ds}$

• Denoising Score Matching objective with loss weighting $\lambda(t)$: $\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}(0,T)} \mathbb{E}_{x_0 \sim q_0(x_0)} \mathbb{E}_{\epsilon \sim \mathcal{N}(0,I)} \frac{\lambda(t)}{\sigma_t^2} ||\epsilon - \epsilon_{\theta}(x_t,t)||_2^2$.

Different loss weighting trade off between model with good quality vs. high log-likelihood:

- Perceptual quality: $\lambda(t) = \sigma_t^2$
- Maximum log-likelihood: $\lambda(t)\beta(t)$ (negative ELBO)

Same objectives as derived in Denoising DPMs!

https://ieeexplore.ieee.org/document/ 6795935 https://arxiv.org/abs/1907.05600 https://arxiv.org/abs/2011.13456

Lecture #15 HY-673



Probability Flow ODE

Lecture #15 HY-673



• Consider reverse generative diffusion SDE:

• Equivalent to "**Probability Flow ODE**" in distribution:
(i.e., solving this ODE results in the same
$$q_t(x_t)$$
 when
initializing $q_T(x_T) \approx \mathcal{N}(x_T; 0, I)$)

$$dx_t = -\frac{1}{2}\beta(t) \left[x_t + 2\nabla_x \log q_t(x_t) \right] dt + \sqrt{\beta(t)} d\bar{W}_t.$$
$$dx_t = -\frac{1}{2}\beta(t) \left[x_t + \nabla_x \log q_t(x_t) \right] dt.$$

https://arxiv.org/abs/2011.13456

Probability Flow ODE

Lecture #15 HY-673



$$dx_t = -\frac{1}{2}\beta(t)\Big[x_t + s_\theta(x_t)\Big]dt.$$

https://arxiv.org/abs/2011.13456

Synthesis with SDE vs. ODE



• Generative Reverse Diffusion SDE (stochastic):

$$dx_t = -\frac{1}{2}\beta(t)\Big[x_t + 2s_\theta(x_t)\Big]dt + \sqrt{\beta(t)}d\bar{W}_t.$$

• Generative Probability Flow ODE (deterministic):

$$dx_t = -\frac{1}{2}\beta(t)\Big[x_t + s_\theta(x_t)\Big]dt.$$

https://arxiv.org/abs/2011.13456

Lecture #15

HY-673

Sampling from "Continuous-Time" Diffusion Models



• Generative (Reverse) Diffusion SDE:

$$dx_t = -\frac{1}{2}\beta(t) \Big[x_t + 2s_\theta(\mathbf{x}_t) \Big] dt + \sqrt{\beta(t)} d\bar{W}_t.$$

1. Euler-Maruyama:

$$x_{t-1} = x_t + \frac{1}{2}\beta(t)\left[x_t + 2s_\theta(x_t)\right]\Delta t + \sqrt{\beta(t)\Delta t}\bar{z}_t$$

2. Ancestral Sampling is also a generative SDE sampler!



- Probability Flow ODE: $dx_t = -\frac{1}{2}\beta(t)\Big[x_t + s_{\theta}(x_t)\Big]dt.$
- 1. Euler's Method: $x_{t-1} = x_t + \frac{1}{2}\beta(t) \Big[x_t + s_{\theta}(x_t) \Big] \Delta t.$
- 2. In practice: Higher-Order ODE solvers (Runge-Kutta, linear multistep methods, exponential integrators). <u>https://arxiv.org/abs/2011.13456</u>

Sampling from "Continuous-Time" Diffusion Models

Lecture #15 HY-673

How to solve the generative SDE or ODE in practice?

- Runge-Kutta adaptive step-size ODE solver [1]
- Higher-Order adaptive step-size SDE solver [2]
- Reparametrized, smoother ODE [3]
- Higher-Order ODE soolver with linear multistepping [4]
- Exponential ODE Integrators [5,6]
- Higher-Order ODE solver with Heun's Method [7]

 \mathbf{x}_t) $\Delta t + \sqrt{\beta(t)\Delta t}\mathcal{N}(\mathbf{0}, \mathbf{I})$

2. In practice: Higher-Order ODE solvers (Runge-Kutta, linear multistep methods, exponential integrators).

[2] https://arxiv.org/abs/2105.14080
[3] https://arxiv.org/abs/2010.02502
[4] https://arxiv.org/abs/2202.09778
[5] https://arxiv.org/abs/2204.13902
[6] https://arxiv.org/abs/2206.00927
[7] https://arxiv.org/abs/2206.00364

[1] https://arxiv.org/abs/2011.13456

Sampling from "Continuous-Time" Diffusion Models

Lecture #15 HY-673



• Generative (Reverse) Diffusion SDE:

$$dx_t = \underbrace{-\frac{1}{2}\beta(t)[x_t + s_{\theta}(x_t)]dt}_{\text{Probability Flow ODE}} \underbrace{-\frac{1}{2} + \beta(t)[s_{\theta}(x_t)]dt}_{\text{Langevin Diffusion}} \underbrace{-\frac{1}{2} + \beta(t)[s_{\theta}(x_t)]dt}_{\text{L$$

- 1. **Pros:** Continuous noise injection can help to compensate errors during diffusion process (Langevin sampling acttively pushes towards correct distribution).
- 2. **Cons:** Often slower, beacuse the stochastic terms themselves require fine discretization during solve.



• Probability Flow ODE:

$$dx_t = -\frac{1}{2}\beta(t)\Big[x_t + s_\theta(x_t)\Big]dt.$$

- 1. **Pros:** Can leverage fast ODE solvers. Best when targeting very fast sampling.
- 2. **Cons:** No "stochastic" error correction, often slightly lower performance than stochastic sampling.

https://arxiv.org/abs/2206.00364

Why use Differential Equation Framework?

Lecture #15 HY-673



Advantages of the Differential Equation framework for Diffusion Models:

- Can leverage broad existing literature on advanced and fast SDE and ODE solvers.
- Allows us to construct deterministic **Probability Flow ODE**.
 - Deterministic Data Encodings.
 - Log-likelihood Estimation.
- Clean mathematical framework based on Diffusion Processes and Score Matching; connections to Neural ODEs, Continuous Normalizing Flows and Energy-based Models.

Diffusion Models as Energy-based Models



- Sample EBM via Langevin dynamics: $x_{t-1} = x_t \eta \nabla_x E_\theta(x_t, t) + \sqrt{2\eta} z_t.$
- Requires only gradient of energy $-\nabla_x E(x,t)$, not $E_{\theta}(x,t)$ itself, nor $Z_{\theta,t}$!

In diffusion models, we learn "energy gradients" for all diffused distributions directly:

$$\nabla_x \log q_t(x) \approx s_\theta(x,t) =: \nabla_x \log p_\theta(x,t) = -\nabla_x E_\theta(x,t) - \underbrace{\nabla_x \log Z_{\theta,t}}_{t} = -\nabla_x E_\theta(x,t).$$

=0 Diffusion Models model energy gradient directly, along entire diffusion process, and avoid modeling partition function. Different noise levels along diffusion are analogous to annealed sampling in EBMs.

Introduction to Deep Generative Modeling

Lecture #15

HY-673 – Computer Science Dep, University of Crete <u>Professors:</u> Yannis Pantazis, Yannis Stylianou <u>Teaching Assistant:</u> Thomas Marchioro