

Generative modeling

: use the dataset to learn a model $p_\theta(x)$ for generating new samples from $p_d(x)$

(Stein) Score matching

: learn an unnormalized density $\tilde{p}_d(x)$ s.t. $p_d(x) = \frac{\tilde{p}_d(x)}{\int \tilde{p}_d(x) dx} \rightarrow Z$

$$\Rightarrow \text{score}(s(x)) : \nabla_x \log p_d(x) = \nabla_x \log \tilde{p}_d(x) - \cancel{\nabla_x \tilde{Z}}$$

→ no need to take care of the intractable partition function Z

$$\Rightarrow \text{score function } (S_\theta(x)) : \nabla_x \log p_\theta(x)$$

$$\rightarrow \min_{\theta} L(\theta) = \frac{1}{2} \mathbb{E}_{p_d(x)} [\|S_\theta(x) - S(x)\|_2^2]$$

(f-divergence between $p_d(x)$ and $p_\theta(x)$)

$$\rightarrow \min_{\theta} J(\theta) = \mathbb{E}_{p_d(x)} \left[\underbrace{\text{tr}(\nabla_x S_\theta(x))}_{\text{computational burden}} + \frac{1}{2} \|S_\theta(x)\|_2^2 \right]$$

\Rightarrow Sliced score matching: faster computation

$$\rightarrow \min_{\theta} L(\theta, p_v) = \frac{1}{2} \mathbb{E}_{p_v} \mathbb{E}_{p_d(x)} [(v^T S_\theta(x) - v^T S(x))^2]$$

$$\left. \begin{array}{l} \text{(projecting score (function) to random direction } v \sim p_v) \\ \text{s.t. } \mathbb{E}_{p_v}[vv^T] \succ 0, \mathbb{E}_{p_v}[\|v\|_2^2] < \infty \end{array} \right\}$$

$$\rightarrow \min_{\theta} J(\theta, p_v) = \mathbb{E}_{p_v} \mathbb{E}_{p_d(x)} [v^T \nabla_x S_\theta(x) v + \frac{1}{2} (v^T S_\theta(x))^2]$$

$$= \mathbb{E}_{p_d(x)} \left[\underbrace{\mathbb{E}_{p_v} [v^T \nabla_x S_\theta(x) v]}_{\text{less computation}} + \frac{1}{2} \|S_\theta(x)\|_2^2 \right]$$

\Rightarrow Denoised score-matching: learn denoising direction

$$\rightarrow q_\sigma(\tilde{x}|x) = \mathcal{N}(\tilde{x}|x, \sigma^2 I) \text{ s.t. } \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x) = - \frac{\tilde{x} - x}{\sigma^2}$$

$$\rightarrow \min_{\theta} \ell(\theta, \sigma) = \frac{1}{2} \mathbb{E}_{p_d(x)} \mathbb{E}_{q(\tilde{x}|x)} [\|S_\theta(\tilde{x}, \sigma) + \frac{\tilde{x} - x}{\sigma^2}\|_2^2]$$

Score-based generative modeling

• sample data via Langevin dynamics using the estimated score function

⇒ Langevin dynamics

$$\rightarrow \tilde{x}_t = \tilde{x}_{t-1} + \frac{\epsilon}{2} S_\theta(\tilde{x}_{t-1}) + \sqrt{\epsilon} z_t \text{ where } z_t \sim N(0, I)$$

→ the distribution of \tilde{x}_T equals $p_\theta(x)$ when $\epsilon \rightarrow 0, T \rightarrow \infty$

⇒ challenges

- inaccurate $S_\theta(x)$ due to manifold hypothesis
- slow mixing of Langevin dynamics due to low data density region

⇒ Noise Conditional Score Networks (NCSN)

① perturb data with random Gaussian noise for better score-matching

(use multiple noise level and train a single conditioned score function $S_\theta(x, \sigma)$)

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

$$\rightarrow \min_{\theta} L(\theta, \{\sigma_i\}_{i=1}^L) = \frac{1}{L} \sum_{i=1}^L \lambda(\sigma_i) \ell(\theta, \sigma_i) \xrightarrow{\text{denoising score-matching loss.}}$$

$$\text{where } \ell(\theta, \sigma_i) = \frac{1}{2} \mathbb{E}_{p_\theta(x)} \mathbb{E}_{q(\tilde{x}|x)} \left[\|S_\theta(\tilde{x}, \sigma_i) + \frac{\tilde{x} - x}{\sigma_i^2}\|_2^2 \right] \text{ and } \lambda(\sigma_i) = \sigma_i^2$$

② initially use scores with large noise and gradually scale down the noise

$$\rightarrow \tilde{x}_t = \tilde{x}_{t-1} + \frac{\alpha_i}{2} S_\theta(\tilde{x}_{t-1}, \sigma_i) + \sqrt{\alpha_i} z_t \text{ where } \alpha_i = \epsilon \cdot \frac{\sigma_i^2}{\sigma_L^2}, z_t \sim N(0, I)$$

Denoising Diffusion Probabilistic Models (DDPM)

: sample data via a parameterized Markov Chain trained to reverse the diffusion process

\Rightarrow diffusion process

\rightarrow given x_0 , forward pass gradually adds noise until signal is destroyed

$$\Rightarrow q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}) \text{ where } q(x_t|x_{t-1}) = N(x_t | \sqrt{1-\beta_t} x_{t-1} + \beta_t I)$$

\rightarrow allow sampling x_t at an arbitrary time-step t in closed-form

$$\Rightarrow q(x_t|x_0) = N(x_t | \sqrt{\bar{\alpha}_t} x_0, (1-\bar{\alpha}_t) I) \text{ where } \alpha_t = 1-\beta_t, \bar{\alpha}_t = \prod_{s=1}^t \alpha_s$$

\rightarrow given x_0 , backward pass can be estimated in closed-form

$$\begin{aligned} \Rightarrow q(x_{t-1}|x_t, x_0) &= \frac{q(x_t|x_0)}{q(x_t|x_0)} q(x_t|x_{t-1}) \\ &= \frac{N(x_t | \sqrt{\bar{\alpha}_{t-1}} x_0, (1-\bar{\alpha}_{t-1}) I)}{N(x_t | \sqrt{\bar{\alpha}_t} x_0, (1-\bar{\alpha}_t) I)} N(x_t | \sqrt{1-\beta_t} x_{t-1} + \beta_t I) \\ &= N\left(x_{t-1} \Big| \underbrace{\frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1-\bar{\alpha}_t} x_0 + \frac{\sqrt{\bar{\alpha}_t} (1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_t}_{\tilde{\mu}_t(x_t, x_0)}, \underbrace{\frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t I}_{\tilde{\Sigma}_t(x_t, x_0)}\right) \end{aligned}$$

⇒ reverse process

$$\rightarrow p(x_{0:T}) = p(x_T) \cdot \prod_{t=1}^T p_\theta(x_{t-1}|x_t) \text{ where } p_\theta(x_{t-1}|x_t) = N(x_{t-1} | \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

→ $\mu_\theta(x_t, t)$ and $\Sigma_\theta(x_t, t)$ can be formulated according to $g(x_{t-1}|x_t, x_0)$

$$\Rightarrow \mu_\theta(x_t, t) = \tilde{\mu}_t \left(x_t, \underbrace{\frac{1}{\sqrt{\alpha_t}} (x_t - \sqrt{1-\alpha_t} \varepsilon_\theta(x_t, t))}_{\propto} \right) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \varepsilon_\theta(x_t, t) \right)$$

$$\Rightarrow \Sigma_\theta(x_t, t) = \sum_t (x_t, x_0) = \sigma_t^2 I$$

→ Training with variational inference

$$\begin{aligned} \Rightarrow \min_{\theta} \mathbb{E}_{p_d(x_0)} [-\log p_\theta(x_0)] &= \mathbb{E}_{p_d(x_0)} \left[\int -\log g(x_{1:T}|x_0) \cdot \frac{p(x_{0:T})}{g(x_{1:T}|x_0)} dx_{1:T} \right] \\ &\leq \mathbb{E}_{g(x_{1:T}|x_0)p_d(x_0)} \left[-\log \frac{p(x_{0:T})}{g(x_{1:T}|x_0)} \right] \\ &= \mathbb{E}_{g(x_{1:T}|x_0)p_d(x_0)} \left[-\log p(x_T) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1}|x_t)}{g(x_t|x_{t-1})} \right] \\ &= \mathbb{E}_{g(x_{1:T}|x_0)p_d(x_0)} \left[-\log p(x_T) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1}|x_t)}{g(x_{t-1}|x_t, x_0)} - \log \frac{p_\theta(x_t|x_t)}{g(x_t|x_0)} \right] \\ &= \mathbb{E}_{g(x_{1:T}|x_0)p_d(x_0)} \left[-\log p(x_T) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1}|x_t)}{g(x_{t-1}|x_t, x_0)} \cdot \frac{g(x_{t-1}|x_0)}{g(x_t|x_0)} - \log \frac{p_\theta(x_t|x_t)}{g(x_t|x_0)} \right] \\ &= \mathbb{E}_{g(x_{1:T}|x_0)p_d(x_0)} \left[\underbrace{\text{KL}(g(x_t|x_0) || p(x_t))}_{0} + \sum_{t \geq 1} \underbrace{\text{KL}(g(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))}_{L(\theta, t)} - \log p(x_0|x_0) \right] \end{aligned}$$

$$\Rightarrow \min_{\theta} L(\theta, t) = \mathbb{E}_{g(x_{1:T}|x_0)p_d(x_0)} \left[\frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t)\|_2^2 \right] + C \quad \text{predict } \mu_\theta(x_t, t) \text{ or }$$

$$= \mathbb{E}_{p \in P_d(x_0)} \left[\underbrace{\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1-\alpha_t)}}_{x_t} \|\varepsilon - \varepsilon_\theta(\underbrace{\sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon, t})\|_2^2 \right] \quad \text{predict } \varepsilon_\theta(x_t, t)$$

remove for simplicity ⇒ down-weight for small t

⇒ more focus on large t for large noise

⇒ inference

Step 1: $x_T \sim N(0, I)$

Step 2: for $t = T, \dots, 1$, $x_{t-1} \sim p_\theta(x_{t-1}|x_t)$

i) $z_t \sim N(0, I)$

ii) $x_{t-1} = \frac{1}{\sqrt{\alpha_t}} (x_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \varepsilon_\theta(x_t, t)) + \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t z_t \approx \text{Langevin dynamics}$

Denoising Diffusion Implicit Models (DDIM)

: use deterministic non-Markovian diffusion process

↳ enable consistency ↳ enable short generative process

↳ enable latent traverse ↳ enable non-Gaussian diffusion process (Categorical)

⇒ non-Markovian forward process

$$\rightarrow q(x_{1:T}|x_0) = q(x_T|x_0) \cdot \prod_{t=2}^T q(x_{t-1}|x_t, x_0)$$

$$\text{where } q(x_{t-1}|x_t, x_0) = N(x_{t-1} | \sqrt{\bar{\alpha}_{t-1}} x_0 + \sqrt{1-\bar{\alpha}_{t-1}-\sigma_t^2} \cdot \frac{x_t - \sqrt{\bar{\alpha}_t} x_0}{\sqrt{1-\bar{\alpha}_t}}, \sigma_t^2 I)$$

$$\text{s.t. } q(x_t|x_0) = N(x_t | \sqrt{\bar{\alpha}_t} x_0, \sqrt{1-\bar{\alpha}_t} I)$$

⇒ reverse process

: given x_t , predict x_0 , and estimate x_{t-1}

$$\rightarrow p_\theta(x_{t-1}|x_t) = \begin{cases} N(x_{t-1} | f_\theta(x_t, t), \sigma_t^2 I) & \text{if } t=1 \\ q(x_{t-1} | x_t, f_\theta(x_t, t)) & \text{otherwise} \end{cases}$$

where $f_\theta(x_t, t) = \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \sqrt{1-\bar{\alpha}_t} \cdot \mathcal{E}_\theta(x_t, t))$ is the predicted x_0 from $q(x_t|x_0)$

⇒ inference

Step 1: $x_T \sim N(0, I)$

Step 2: for $t=T, \dots, 1$, $x_{t-1} \sim p_\theta(x_{t-1}|x_t)$

$$\text{i)} z_t \sim N(0, I)$$

$$\text{ii)} x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \cdot f_\theta(x_t, t) + \sqrt{1-\bar{\alpha}_{t-1}-\sigma_t^2} \cdot \underbrace{\mathcal{E}_\theta(x_t, t)}_{\frac{x_t - \sqrt{\bar{\alpha}_t} f_\theta(x_t, t)}{\sqrt{1-\bar{\alpha}_t}}} + \sigma_t z_t$$

(if $\sigma_t^2=0$ for all t , it becomes an implicit model)

(a subset of time steps can be used if $\sum_{t \in S} KL(q(x_t|x_0) \| p_\theta(x_t|x_0))$ is added to the loss)

$$\rightarrow N(x_0 | f_\theta(x_t, t), \sigma_t^2 I)$$

Improved Denoising Diffusion Probabilistic Models

: achieve competitive log-likelihood with high sample quality

⇒ learn the reverse process variances and use hybrid objective

$$\rightarrow \Sigma_\theta(x_t, t) = \exp(v_\theta(x_t, t) \log \beta_t + (1 - v_\theta(x_t, t)) \log \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t)$$

$$\rightarrow L_{\text{hybrid}} = \underbrace{L_{\text{simple}}}_{\text{for } \mu_\theta(x_t, t)} + \underbrace{L_{\text{vrb}}}_{\text{for } \Sigma_\theta(x_t, t)}$$

⇒ improve the noise schedule

$$\rightarrow \bar{\alpha}_t = \frac{f(t)}{f(0)} \text{ where } f(t) = \cos\left(\frac{t/T+\epsilon}{1+\epsilon} \cdot \frac{\pi}{2}\right) \text{ s.t. } \beta_t = 1 - \frac{\bar{\alpha}_t}{\bar{\alpha}_{t-1}}$$

⇒ reduce gradient noise of L_{vrb}

$$\rightarrow L_{\text{vrb}} = \mathbb{E}_{p_t} \left[\frac{L_t}{p_t} \right] \text{ where } p_t \propto \sqrt{\mathbb{E}[L_t]} \text{ and } \sum p_t = 1$$

↳ maintain a history of the previous 10 values

Diffusion Models Beat GAN on Image Synthesis

⇒ architecture improvement

- ★ U-Net
 - stack of residual layers & downsampling convolutions
 - stack of residual layers & upsampling convolutions
 - skip connection connecting layers with the same sampling size
 - global attention layer at 16x16 resolution with a single head
 - projection of the timestep embedding into each residual block

(Δ) ⇒ increase depth versus width, holding model size relatively constant

(○) ⇒ increase the number of attention heads

(○) ⇒ use attention at 32x32, 16x16, and 8x8 resolutions rather than not only 16x16

(○) ⇒ use the BigGAN residual block for upsampling and downsampling

(✗) ⇒ rescale residual connections with $\frac{1}{\sqrt{2}}$

(○) ⇒ more heads or fewer channels per head

(○) ⇒ use adaptive group normalization to incorporate the timestep and class embedding into residual block

$$(\text{AdaGN}(h, y) = y_{\text{time}} \text{GN}(h) + y_{\text{class}} \text{ s.t. } y = [y_{\text{time}}, y_{\text{class}}])$$

⇒ classifier guided sampling

: train a classifier $p_\theta(y|x_{t-1})$ and use $\nabla_{x_{t-1}} \log p_\theta(y|x_{t-1})$ to further guide inference

→ conditional diffusion process

(Suppose $\hat{q}_f(y|x_0)$ is given. Then, \hat{q}_f behaves like q_f)

$$\textcircled{1} \quad \hat{q}_f(x_t|x_{t-1}, y) := q_f(x_t|x_{t-1}) \Rightarrow \hat{q}_f(x_t|x_{t-1}) = q_f(x_t|x_{t-1})$$

$$\textcircled{2} \quad \hat{q}_f(x_{1:T}|x_0, y) := \prod_{t=1}^T \hat{q}_f(x_t|x_{t-1}, y) \Rightarrow \hat{q}_f(x_{1:T}|x_0) = q_f(x_{1:T}|x_0)$$

$$\therefore \hat{q}_f(x_t) = q_f(x_t) \text{ and } \hat{q}_f(x_{t-1}|x_t) = q_f(x_{t-1}|x_t) \text{ and } \hat{q}_f(y|x_{t-1}, x_t) = \hat{q}_f(y|x_t)$$

→ conditional reverse process

$$\begin{aligned} \Rightarrow \hat{q}_f(x_{t-1}|x_t, y) &= \frac{q(x_{t-1}|x_t) \hat{q}_f(y|x_{t-1}, x_t)}{\hat{q}_f(y|x_t)} = \frac{q(x_{t-1}|x_t) \hat{q}_f(y|x_{t-1})}{\hat{q}_f(y|x_t)} \\ &\propto q(x_{t-1}|x_t) \hat{q}_f(y|x_{t-1}) \\ &\quad \underset{ss}{p_\theta(x_{t-1}|x_t)} \quad \underset{ss}{p_\theta(y|x_{t-1})} \end{aligned}$$

$$\Rightarrow \nabla_{x_{t-1}} \log \hat{q}_f(x_{t-1}|x_t, y) = \nabla_{x_{t-1}} \log p_\theta(x_{t-1}|x_t) + \nabla_{x_{t-1}} \log p_\theta(y|x_{t-1})$$

$$\rightarrow \log p_\theta(x_{t-1}|x_t) = \log N(\mu, \Sigma) = -\frac{1}{2} (x_{t-1} - \mu)^T \Sigma^{-1} (x_t - \mu) + C$$

$$\begin{aligned} \rightarrow \log p_\theta(y|x_{t-1}) &\approx \log p_\theta(y|x_{t-1}) \Big|_{x_{t-1}=\mu} + (x_{t-1} - \mu) \nabla_{x_{t-1}} \log p_\theta(y|x_{t-1}) \Big|_{x_{t-1}=\mu} \\ &= (x_{t-1} - \mu) g + C_1 \end{aligned}$$

$$\begin{aligned} \rightarrow \log p_\theta(x_{t-1}|x_t) p_\theta(y|x_{t-1}) &\approx -\frac{1}{2} (x_t - \mu - \underbrace{\Sigma g}_{\text{shifting mean}})^T \Sigma^{-1} (x_t - \mu - \Sigma g) + C_2 \\ &= \log N(\mu + \underbrace{\Sigma g}_{\text{shifting mean}}, \Sigma) + C_3 \end{aligned}$$

\Rightarrow inference

Step 1: $x_t \sim N(0, I)$

Step 2: for $t = T, \dots, 1$,

"option 1": $x_{t-1} \sim \hat{q}(x_{t-1} | x_t, y)$

i) estimate $p_\theta(x_{t-1} | x_t) = N(x_{t-1} | \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$

ii) $x_{t-1} \sim N(\mu_\theta(x_t, t) + \underbrace{s \nabla_{x_t} \log p_\phi(y | x_t)}_{\text{introduces trade-off b/t fidelity and diversity}}, \Sigma_\theta(x_t, t))$

\hookrightarrow needs to be set to larger than 1.

"option 2": DDIM sampling with calibrated noise prediction

i) $\hat{\varepsilon} = \varepsilon_\theta(x_t, t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log p_\phi(y | x_t) \implies$ score function for $p(x_t, y) = p(x_t) \cdot p(y | x_t)$

$$\left(\because \nabla_{x_t} \log p_\theta(x_t) + \nabla_{x_t} \log p_\phi(y | x_t) = -\frac{1}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_\theta(x_t, t) + \nabla_{x_t} \log p_\phi(y | x_t) \right)$$

$$\text{ii) } x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \cdot \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\varepsilon}}{\sqrt{\bar{\alpha}_t}} + \sqrt{1 - \bar{\alpha}_{t-1}} \cdot \hat{\varepsilon}$$