

DDPM

Denoising Diffusion Probabilistic Models (2020) [1k](#)

Diffusion process

Forward diffusion process: input image x_0 , add Gaussian noise through T steps

Reverse diffusion process, or in general, sampling process of a generative model: train a neural network to recover the image

Forward diffusion

given x_0 sampled from real distribution $q(x)$:

$$x_0 \sim q(x)$$

add Gaussian noise with variance β_t to x_{t-1} , producing a new latent variable x_t :

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \mu_t = \sqrt{1 - \beta_t}x_{t-1}, \Sigma_t = \beta_t I)$$

from x_0 to x_T :

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

reparameterization trick: define $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{s=0}^t \alpha_s$, $\epsilon_t \sim \mathcal{N}(0, I)$

$$\begin{aligned} x_t &= \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_{t-1} \\ &= \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1} \\ &= \dots \\ &= \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0 \end{aligned}$$

Thus we can sample x_t at any arbitrary timestep using

$$x_t \sim q(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$$

This will be our target later on to calculate our tractable objective loss L_t

The variance parameter β_t can be chosen as a constant or a schedule, e.g., in DDPM paper, linear schedule from $\beta_1 = 0.0001$ to $\beta_T = 0.02$

Reverse diffusion

approximate $q(x_{t-1}|x_t)$ with a parameterized model p_θ (e.g. a neural network); since $q(x_{t-1}|x_t)$ will also be Gaussian for small enough β_t , we choose p_θ to be Gaussian and

parameterize the mean and variance:

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

from x_T to x_0 :

$$p_\theta(x_{0:T}) = p_\theta(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$$

Training a diffusion model

needs a reference image x_0 to draw an image, we can sample x_t at noise level t conditioned on x_0 :

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)$$

$$\text{where } \tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t$$

$$\text{and } \tilde{\mu}(x_t, x_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t} x_0 + \frac{\sqrt{\bar{\alpha}}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_t$$

we already saw in the reparameterization trick that

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1-\bar{\alpha}_t}\epsilon) \text{ where } \epsilon \sim \mathcal{N}(0, I)$$

combining last two equations:

$$\tilde{\mu}_t(x_t) = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon)$$

thus we can use a neural network $\epsilon_\theta(x_t, t)$ to approximate ϵ and consequently the mean:

$$\tilde{\mu}_\theta(x_t, t) = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(x_t, t))$$

thus the loss function (the denoising term in the ELBO) can be expressed as:

$$\begin{aligned} L_t &= \mathbb{E}_{x_0, t, \epsilon} \left[\frac{1}{2 \|\Sigma_\theta(x_t, t)\|_2^2} \|\tilde{\mu}_t - \mu_\theta(x_t, t)\|_2^2 \right] \\ &= \mathbb{E}_{x_0, t, \epsilon} \left[\frac{1}{2 \alpha_t (1 - \bar{\alpha}_t) \|\Sigma_\theta\|_2^2} \|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|_2^2 \right] \end{aligned}$$

in DDPM paper, simplified by ignoring the weighting term:

$$L_t = \mathbb{E}_{x_0, t, \epsilon} [\|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|_2^2]$$

training and sampling algorithms of DDPM paper:

Algorithm 1 Training

```
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
        $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$ 
6: until converged
```

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

Background Part of Our Course Reserach Project

Our researches belong to the large family of Diffusion models. The models are inspired by *diffusion process*, which is a kind of random process $(X_t)_{t \geq 0}$ that pictures the location of a particle moving at random but also governed by a drift and a random noise. Unlike other generative models which use discriminators (GAN) or encoders (VAE), etc., diffusion models use a Markov chain that gradually adds Gaussian noise to the data according to a variance schedule β_1, \dots, β_T as the approximate posterior $q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)$, called the *forward diffusion process*:

$$q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$$
$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

This way, the sample \mathbf{x}_0 gradually loses its features and as $T \rightarrow \infty$, \mathbf{x}_T becomes an isotropic Gaussian distribution. The variances β_t can be learned by reparameterization or held constant as hyperparameters, and expressiveness of the reverse process is ensured in part by the choice of Gaussian conditionals in $p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$, because both processes have the same functional form when β_t are small.

Now we want to reverse the above process to recreate the sample \mathbf{x}_0 from the heavily noised \mathbf{x}_T . The model then predicts the distriubtion $p_{\theta}(\mathbf{x}_0) := \int p_{\theta}(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T}$, where $\mathbf{x}_1, \dots, \mathbf{x}_T$ are latents of the same dimensionality as the data $\mathbf{x}_0 \sim q(\mathbf{x}_0)$. The joint distribution $p_{\theta}(\mathbf{x}_{0:T})$ is the *reverse process* as a Markov chain with learned Gaussian transitions starting at $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$:

$$p_{\theta}(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)$$
$$p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) := \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$$