

Диффузные модели, часть 1

дискретное время

Григорий Бартош

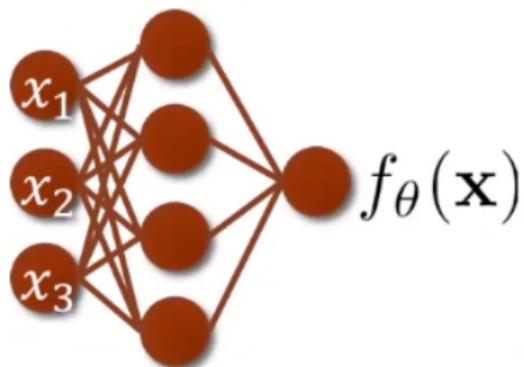
8 октября 2021

Energy-Based Models

- Energy-Based Models (EBMs)

$$p_{\theta}(x) = \frac{e^{-f_{\theta}(x)}}{Z_{\theta}}, \quad f_{\theta}(x) \in \mathbb{R}$$

$$Z_{\theta} = \int e^{-f_{\theta}(x)} dx$$



- Плюсы Простая модель
- Минусы Сложно обучать θ через максимизацию правдоподобия

$$\mathbb{E}_{p_{\text{data}}(x)}[-\log p_{\theta}(x)] = \mathbb{E}_{p_{\text{data}}(x)}[f_{\theta}(x) + \log Z_{\theta}]$$

Обучение EBMs через score matching

- Score function

$$s(x) = \nabla_x \log p(x)$$

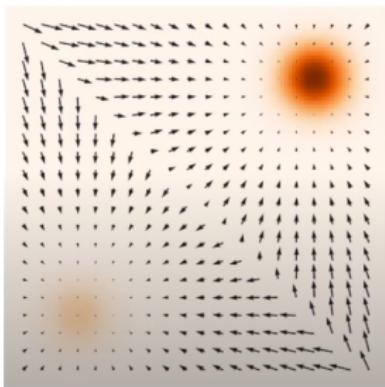
- Score function для EBM модели

$$s_\theta(x) = \nabla_x \log p_\theta(x) = -\nabla_x f_\theta(x) - \underbrace{\nabla_x \log Z_\theta}_{=0} = -\nabla_x f_\theta(x)$$

- Score matching (Fisher divergence)

$$\mathbb{E}_{p_{\text{data}}(x)} [\|s_\theta(x) - \nabla_x \log p_{\text{data}}(x)\|_2^2]$$

Score matching иллюстрация



$$\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x})$$

The equation illustrates the equivalence between the gradient of the log-data probability and the gradient of the log-model probability. On the left, a set of red arrows represents the gradient of the log-data probability. A blue double approximation symbol (\approx) connects it to a set of red arrows on the right, which represent the gradient of the log-model probability. This visualizes how the score function (gradient of the log-probability) is estimated by a neural network.

Denoising score matching

- Зашумим данные

$$p_\sigma(\tilde{x}) = \int_x q_\sigma(\tilde{x}|x)p_{\text{data}}(x)dx, \quad q_\sigma(\tilde{x}|x) = \mathcal{N}(\tilde{x}; x, \sigma^2 I)$$

- Score matching для зашумленных данных

$$\mathbb{E}_{p_\sigma(x)}[\|s_\theta(x) - \nabla_x \log p_\sigma(x)\|_2^2]$$

- Можно переписать как

$$\mathbb{E}_{q_\sigma(\tilde{x}|x)p_{\text{data}}(x)}[\|s_\theta(x) - \nabla_{\tilde{x}} \log q_\sigma(\tilde{x}|x)\|_2^2]$$

Sliced score matching – I

- Многомерный случай

$$\frac{1}{2} \mathbb{E}_{p_{\text{data}}(x)} [\|\nabla_x \log p_{\text{data}}(x) - s_\theta(x)\|_2^2]$$

- Одномерный случай

$$\begin{aligned}\mathcal{L}(\theta) &= \frac{1}{2} \mathbb{E}_{p_{\text{data}}(x)} [(\nabla_x \log p_{\text{data}}(x) - s_\theta(x))^2] \\ &= \frac{1}{2} \int p_{\text{data}}(x) [(\nabla_x \log p_{\text{data}}(x) - s_\theta(x))^2] dx \\ &= \frac{1}{2} \int p_{\text{data}}(x) (\nabla_x \log p_{\text{data}}(x))^2 dx - \\ &\quad \int p_{\text{data}}(x) \nabla_x \log p_{\text{data}}(x) s_\theta(x) dx + \\ &\quad \frac{1}{2} \int p_{\text{data}}(x) s_\theta(x)^2 dx\end{aligned}$$

Sliced score matching – II

$$\begin{aligned} \int p_{\text{data}}(x) \nabla_x \log p_{\text{data}}(x) s_\theta(x) dx &= \\ &= \int p_{\text{data}}(x) \frac{\nabla_x p_{\text{data}}(x)}{p_{\text{data}}(x)} s_\theta(x) dx \\ &= \int s_\theta(x) \nabla_x p_{\text{data}}(x) dx \\ &= s_\theta(x) p_{\text{data}}(x) \Big|_{-\infty}^{+\infty} - \int p_{\text{data}}(x) \nabla_x s_\theta(x) dx \\ &= - \int p_{\text{data}}(x) \nabla_x s_\theta(x) dx \end{aligned}$$

Sliced score matching – III

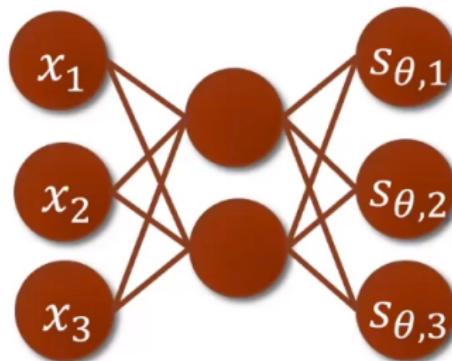
$$\begin{aligned}\mathcal{L}(\theta) &= C + \int p_{\text{data}}(x) \nabla_x s_\theta(x) dx + \frac{1}{2} \int p_{\text{data}}(x) s_\theta(x)^2 dx \\ &= \mathbb{E}_{p_{\text{data}}(x)} \left[\frac{1}{2} s_\theta(x)^2 + \nabla_x s_\theta(x) \right]\end{aligned}$$

- Обратно к многомерному случаю

$$\mathbb{E}_{p_{\text{data}}(x)} \left[\frac{1}{2} \|s_\theta(x)\|_2^2 + \text{tr}(\nabla_x s_\theta(x)) \right]$$

Sliced score matching – IV

- Score model



- Требует $\mathcal{O}(D)$ backpropagation

Sliced score matching – V

- Hutchinson trace estimator

$$\mathbb{E}_{p(v)} [v^T A v] = \text{tr}(A), \quad v \sim \{-1, 1\}^D$$

- Можно перейти к

$$\mathbb{E}_{p(v)} \mathbb{E}_{p_{\text{data}}(x)} \left[\frac{1}{2} \|s_\theta(x)\|_2^2 + v^T \nabla_x s_\theta(x) v \right]$$

- Vector jacobian product (vjp) можно посчитать за $\mathcal{O}(1)$ backpropagation

$$v^T \nabla_x s_\theta(x) = \nabla_x v^T s_\theta(x)$$

Langevin dynamics

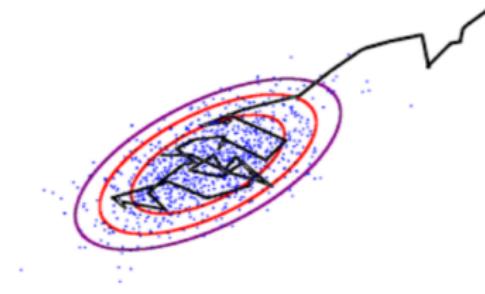
- Имеем доступ к

$$p(x), \quad \nabla_x \log p(x)$$

- Запустим динамику Ланжевена

$$x_{i+1} \leftarrow x_i + \epsilon \nabla_x \log p(x) + \sqrt{2\epsilon} z_i, \quad i = 0, 1, \dots, K,$$

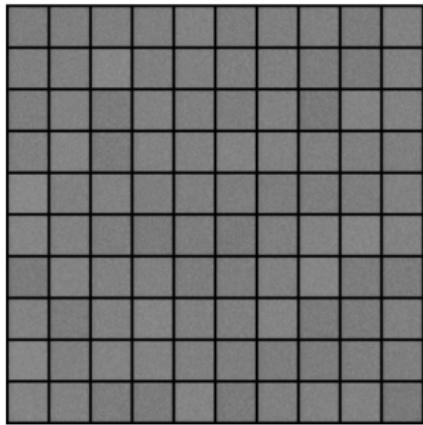
Где $z_i \sim \mathcal{N}(0, I)$ и $K \rightarrow \infty$.



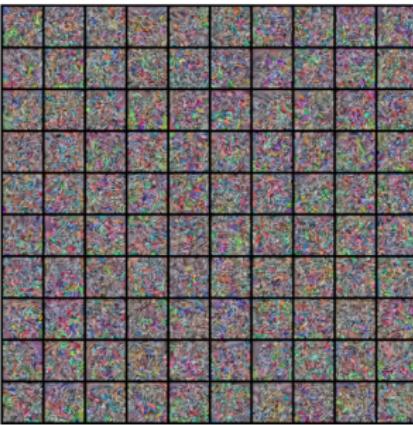
Score matching промежуточные итоги

- Задаем модель
 - Как EBM модель и переходим к score function
 - Задаем score function
- Обучаем score matching
 - Denoising score matching, зашумляя данные
 - Sliced score matching
- Сэмплируем
 - Динамикой Ланжевена

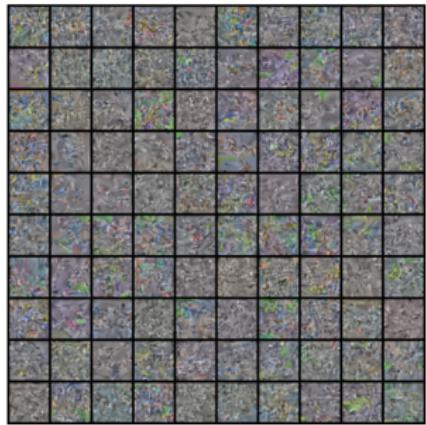
Найвный score mutching



(a) MNIST



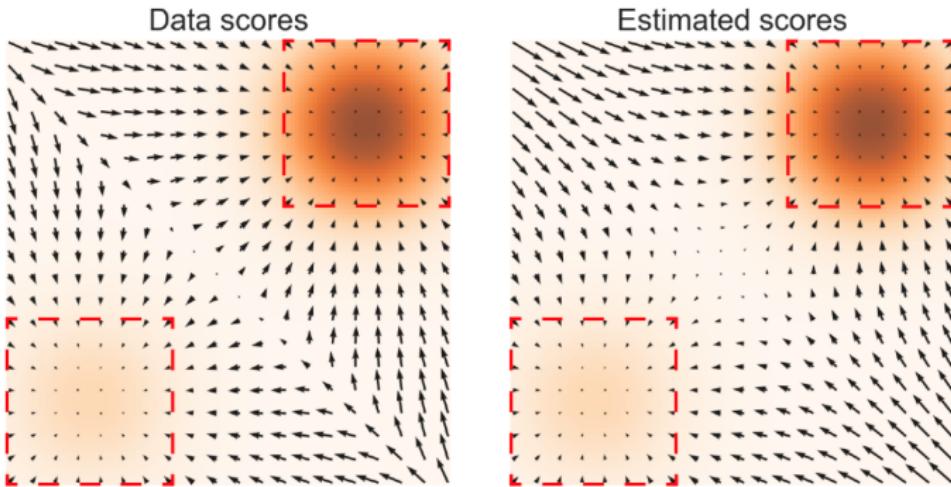
(b) CelebA



(c) CIFAR-10

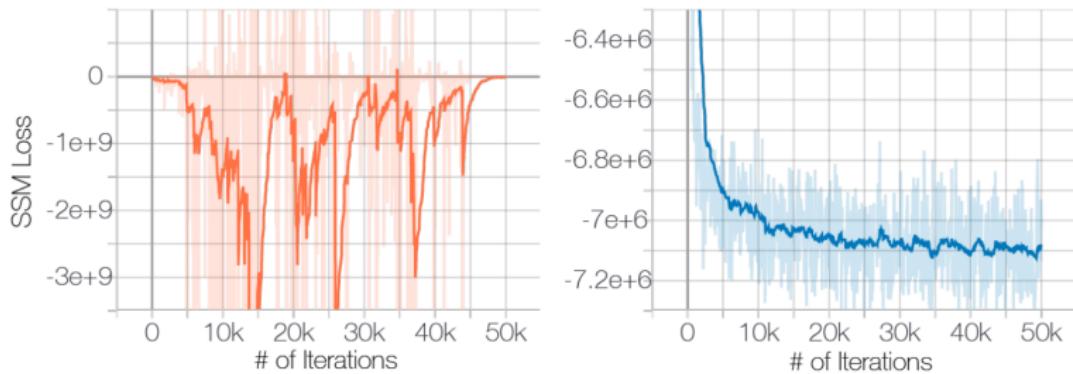
Score matching, регионы без данных

- Ничего не знаем про регионы, где нет данных



Score matching, нестабильность обучения

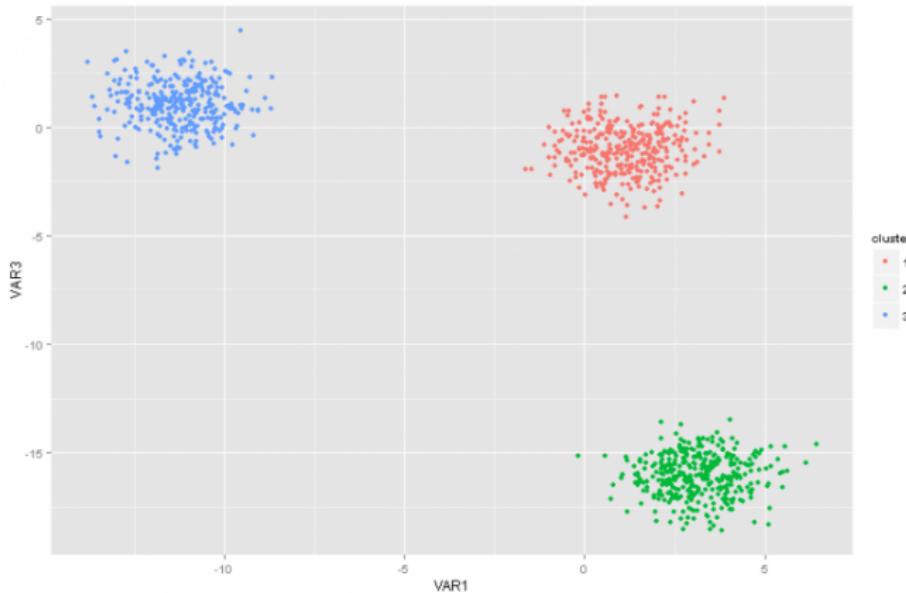
- Sliced score matching дает нестабильное обучение



- Добавление шума стабилизирует обучение

Score matching, ограничения динамики Ланжевена

- Динамика Ланжевена с трудом переходит между модами



Несколько уровней зашумления

- Введем несколько уровней зашумления

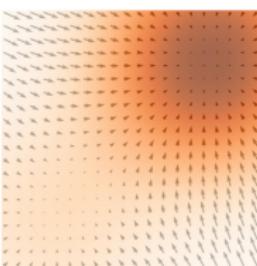
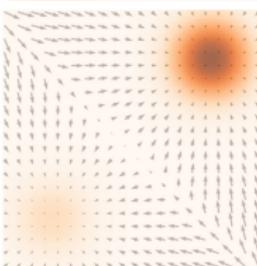
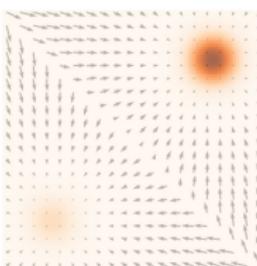
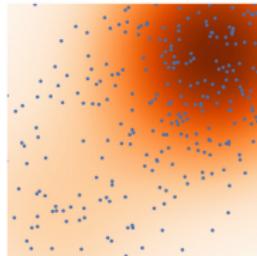
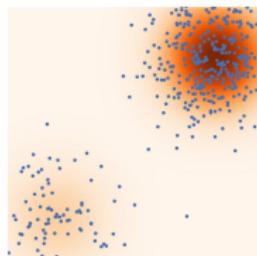
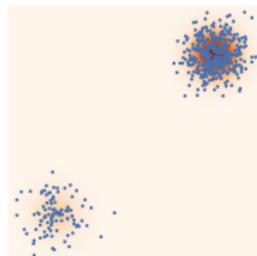
$$\sigma_1$$

<

$$\sigma_2$$

<

$$\sigma_3$$



Noise Conditiona Score Network NCSN

- Обусловим score function $s_\theta(x)$ на уровень i .
- На каждом уровне мы хотим

$$s_\theta(x, i) \approx \nabla_x \log p_{\sigma_i}(x)$$

- Будем обучать

$$\sum_{i=1}^L \lambda(i) \mathbb{E}_{p_{\sigma_i}(x)} [\|\nabla_x \log p_{\sigma_i}(x) - s_\theta(x, i)\|_2^2], \quad \lambda(i) = \sigma_i^2$$

- Авторы брали $L = 10$ и $T = 100$.

Algorithm 1 Annealed Langevin dynamics.

Require: $\{\sigma_i\}_{i=1}^L, \epsilon, T$.

```
1: Initialize  $\tilde{\mathbf{x}}_0$ 
2: for  $i \leftarrow 1$  to  $L$  do
3:    $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$        $\triangleright \alpha_i$  is the step size.
4:   for  $t \leftarrow 1$  to  $T$  do
5:     Draw  $\mathbf{z}_t \sim \mathcal{N}(0, I)$ 
6:      $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_{\theta}(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$ 
7:   end for
8:    $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$ 
9: end for
return  $\tilde{\mathbf{x}}_T$ 
```

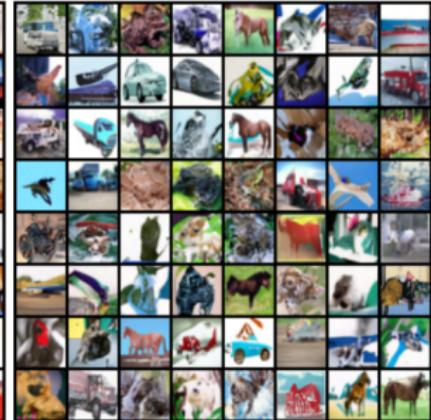
NCSN сэмплирование, примеры



(a) MNIST



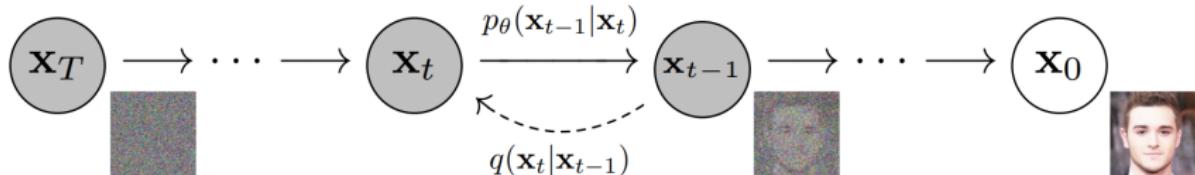
(b) CelebA



(c) CIFAR-10

DDPM, зашумление

- Общая схема



- Зашумление данных

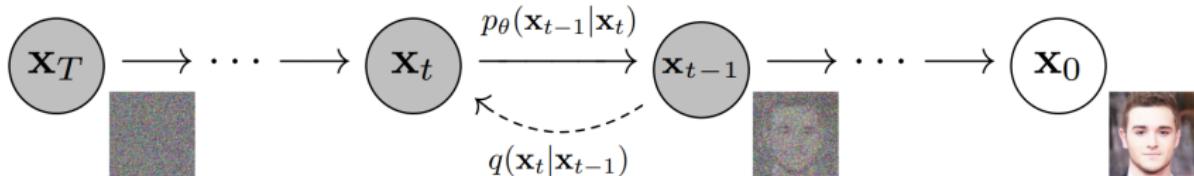
$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}), \quad q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$$

- Зашумление на несколько шагов

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

DDPM, удаление шума

- Общая схема



- Удаление шума

$$p_{\theta}(x_{0:T}|x_0) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t), \quad p_{\theta}(x_T) = \mathcal{N}(x_T; 0, I)$$

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

DDPM, вероятность в точки

$$\begin{aligned} p_\theta(x_0) &= \int p_\theta(x_{0:T}) dx_{1:T} \\ &= \int p_\theta(x_{0:T}) \frac{q(x_{1:T}|x_0)}{q(x_{1:T}|x_0)} dx_{1:T} \\ &= \int q(x_{1:T}|x_0) \frac{p_\theta(x_{0:T})}{q(x_{1:T}|x_0)} dx_{1:T} \\ &= \int q(x_{1:T}|x_0) p_\theta(x_T) \prod_{t=1}^T \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1}, x_0)} dx_{1:T} \end{aligned}$$

DDPM, вариационная нижняя оценка – I

$$\begin{aligned}\mathcal{L}_\theta &= \mathbb{E}_{q(x_0)} [-\log p_\theta(x_0)] \\&= \int q(x_0) \left[-\log \left[\int q(x_{1:T}|x_0) p_\theta(x_T) \prod_{t=1}^T \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1}, x_0)} dx_{1:T} \right] \right] dx_0 \\&\leq - \int q(x_0) \left[\int q(x_{1:T}|x_0) \log \left[p_\theta(x_T) \prod_{t=1}^T \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1}, x_0)} \right] dx_{1:T} \right] dx_0 \\&= - \int q(x_{0:T}) \log \left[p_\theta(x_T) \prod_{t=1}^T \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1}, x_0)} \right] dx_{0:T} \\&= \mathbb{E}_{q(x_{0:T})} \left[-\log p_\theta(x_T) - \sum_{t=1}^T \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1}, x_0)} \right]\end{aligned}$$

DDPM, вариационная нижняя оценка – II

$$\begin{aligned}&= \mathbb{E}_{q(x_0:T)} \left[-\log p_\theta(x_T) - \sum_{t=1}^T \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1}, x_0)} \right] \\&= \mathbb{E}_{q(x_0:T)} \left[-\log p_\theta(x_T) - \sum_{t>1} \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1}, x_0)} - \log \frac{p_\theta(x_0|x_1)}{q(x_t|x_1, x_0)} \right] \\&= \mathbb{E}_{q(x_0:T)} \left[-\log p_\theta(x_T) - \sum_{t>1} \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \frac{p_\theta(x_{t-1}|x_0)}{q(x_t|x_0)} - \log \frac{p_\theta(x_0|x_1)}{q(x_t|x_1, x_0)} \right] \\&= \mathbb{E}_{q(x_0:T)} \left[-\log \frac{p_\theta(x_T)}{q(x_T|x_0)} - \sum_{t>1} \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} - \log p_\theta(x_0|x_1) \right] \\&= \mathbb{E}_{q(x_0:T)} \left[D_{KL}(q(x_T|x_0) || p_\theta(x_T)) + \sum_{t>1} D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t)) - \log p_\theta(x_0|x_1) \right]\end{aligned}$$

DDPM, вариационная нижняя оценка – III

- ELBO

$$\mathbb{E}_{q(x_0:T)} \left[D_{KL}(q(x_T|x_0)||p_\theta(x_T)) + \sum_{t>1} D_{KL}(q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)) - \log p_\theta(x_0|x_1) \right]$$

- Распишем $q(x_{t-1}|x_t, x_0)$

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

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

$$\tilde{\beta}_t(x_t, x_0) = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t$$

DDPM, репараметризация – I

- Распишем \mathcal{L}_{t-1}

$$\mathcal{L}_{t-1} = \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t)\|^2 \right]$$

- Заметим

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

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t(x_0, \epsilon) - \sqrt{1 - \bar{\alpha}_t} \epsilon)$$

- Перепишем

$$\begin{aligned} \mathcal{L}_{t-1} &= \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \left\| \tilde{\mu}_t \left(x_t(x_0, \epsilon), \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t(x_0, \epsilon) - \sqrt{1 - \bar{\alpha}_t} \epsilon) \right) - \mu_\theta(x_t, t) \right\|^2 \right] \\ &= \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(x_t(x_0, \epsilon) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right) - \mu_\theta(x_t, t) \right\|^2 \right] \end{aligned}$$

DDPM, репараметризация – II

- Перепишем μ_θ

$$\mu_\theta = \frac{1}{\sqrt{\alpha_t}} \left(x_t(x_0, \epsilon) - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right)$$

- Снова распишем \mathcal{L}_{t-1}

$$\mathcal{L}_{t-1} = \mathbb{E}_{q, \epsilon} \left[\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} \|\epsilon - \epsilon_\theta(x_t, t)\|^2 \right]$$

DDPM сэмплирование, примеры

