

Генеративные модели

Лекция 8: *Diffusion Models*

Повторение

Energy – Based Models

Определим плотность $p_{\theta}(\mathbf{x})$ через энергию с помощью распределения Больцмана:

$$p_{\theta}(\mathbf{x}) = \frac{e^{-E_{\theta}(\mathbf{x})}}{Z_{\theta}}$$

- $E_{\theta}(\mathbf{x})$ – ненормализованная плотность
- Z_{θ} – нормировочная константа (статистическая сумма)

Z_{θ} – это интеграл по всему пространству данных:

$$Z_{\theta} = \int e^{-E_{\theta}(\mathbf{x})} d\mathbf{x}$$

- Для сложных высокоразмерных данных мы не можем вычислить этот интеграл
- Не зная Z_{θ} , мы не сможем найти $p_{\theta}(\mathbf{x})$ и его градиент для обучения

Energy – Based Models

$$\mathcal{L}_\theta = -\mathbb{E}_{p_{data}(\mathbf{x})} \nabla_\theta E_\theta(\mathbf{x}) + \mathbb{E}_{p_\theta(\mathbf{x})} \nabla_\theta E_\theta(\mathbf{x})$$

- Нам нужно уметь оценивать $\mathbb{E}_{p_\theta(\mathbf{x})}[\dots]$
- Будем оценивать этот член с помощью метода Монте-Карло (*Markov Chain Monte Carlo*, *MCMC*)

Введем понятие *функции оценки* (*score function*):

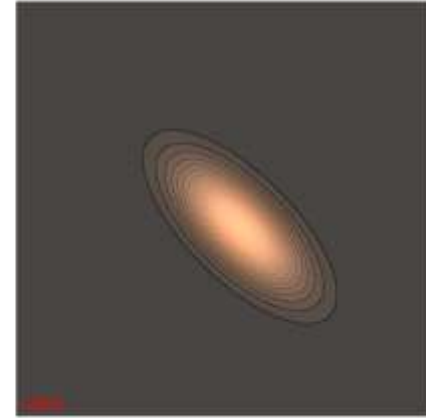
$$s_\theta(\mathbf{x}) = \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x})$$

Langevin Dynamics

Для сэмплирования из $p_\theta(\mathbf{x})$ будем использовать *динамику Ланжевена (Langevin Dynamics)*:

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\eta}{2} \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}_t) + \sqrt{\eta} \epsilon_t = \mathbf{x}_t + \frac{\eta}{2} s_\theta(\mathbf{x}_t) + \sqrt{\eta} \epsilon_t$$

- \mathbf{x}_t — текущая точка
- η — размер шага
- $\epsilon_t \sim \mathcal{N}(\mathbf{0}, I)$ — шум



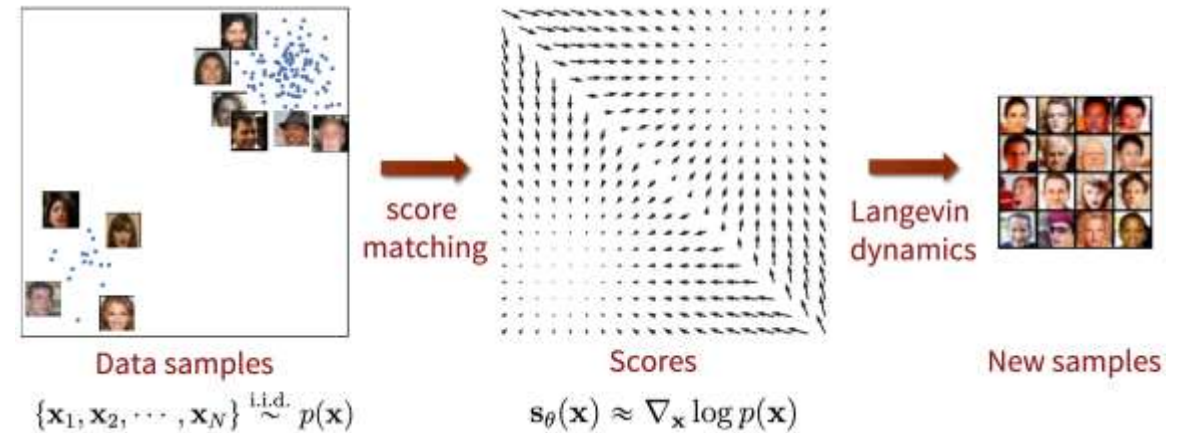
- Член $\frac{\eta}{2} \nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}_t)$ — это градиентный подъем, который толкает нашу точку \mathbf{x} в направлении наибольшего возрастания плотности
- Член $\sqrt{\eta} \epsilon_t$ добавляет случайности и не дает нашей точке \mathbf{x} застрять в локальном максимуме

Проблемы МСМС

- Обучение с использованием МСМС обычно очень медленное, не стабильное и даёт смещенные градиенты

Идея:

- Попробуем обучать *score function* $s_{\theta}(\mathbf{x})$ напрямую



Score Matching

Цель *score matching* (Hyvärinen, 2005) заключается в минимизации дивергенции Фишера:

$$D_F(p_{data}(\mathbf{x}) || p_{\theta}(\mathbf{x})) = \mathbb{E}_{p_{data}(\mathbf{x})} \left[\frac{1}{2} || \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_{data}(\mathbf{x}) ||^2 \right]$$

- Минимизируя эту функцию потерь, мы заставляем *score*-функцию модели $s_{\theta}(\mathbf{x})$ соответствовать *score*-функции реальных данных $s_{data}(\mathbf{x})$
- Если $s_{\theta}(\mathbf{x}) \approx s_{data}(\mathbf{x})$, то $p_{\theta}(\mathbf{x}) \approx p_{data}(\mathbf{x})$

Мы не знаем *score*-функцию реальных данных $s_{data}(\mathbf{x})$

Hyvärinen показал, что при определенных условиях регулярности функция потерь может быть переписана в виде:

$$D_F(p_{data}(\mathbf{x}) || p_{\theta}(\mathbf{x})) = \mathbb{E}_{p_{data}(\mathbf{x})} \left[\frac{1}{2} || s_{\theta}(\mathbf{x}) - s_{data}(\mathbf{x}) ||^2 \right] = \mathbb{E}_{p_{data}(\mathbf{x})} \left[\frac{1}{2} || s_{\theta}(\mathbf{x}) ||^2 + tr(\nabla_{\mathbf{x}} s_{\theta}(\mathbf{x})) \right]$$

Проблемы *Score Matching*

- Классический **score – matching** требует, чтобы $p_{data}(\mathbf{x})$ была гладкой и непрерывно дифференцируемой
- Реальные данные часто могут быть дискретны (изображения 0 – 255)
- Для таких данных $\log p_{data}(\mathbf{x})$ является разрывной функцией

Идея:

Добавим небольшой шум к нашим данным:

$$\mathbf{x}_\sigma = \mathbf{x} + \sigma \cdot \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, I)$$

Процесс добавления шума описывается условной вероятностью:

$$q(\mathbf{x}_\sigma | \mathbf{x}) = \mathcal{N}(\mathbf{x}_\sigma | \mathbf{x}, \sigma^2 \cdot I)$$

Denoising Score Matching

Мы получили новое гладкое распределение, к которому можно применять **score – matching**:

$$q(\mathbf{x}_\sigma) = \int q(\mathbf{x}_\sigma | \mathbf{x}) \cdot p_{data}(\mathbf{x}) d\mathbf{x}$$

Будем минимизировать дивергенцию Фишера между *score*-функцией модели $s_\theta(\mathbf{x})$ и *score*-функцией зашумленных данных $q(\mathbf{x}_\sigma)$:

$$\frac{1}{2} \mathbb{E}_{q(\mathbf{x}_\sigma)} \left[\left\| s_\theta(\mathbf{x}_\sigma) - \nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma) \right\|_2^2 \right] \rightarrow \min_\theta$$

Мы по-прежнему не знаем $\nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma)$, так как $q(\mathbf{x}_\sigma)$ зависит от $p_{data}(\mathbf{x})$:

Vincent (2010) показал, что минимизация дивергенции Фишера эквивалентна минимизации новой функции потерь:

$$\mathbb{E}_{q(\mathbf{x}_\sigma)} \left[\left\| s_\theta(\mathbf{x}_\sigma) - \nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma) \right\|_2^2 \right] = \mathbb{E}_{p_{data}(\mathbf{x})} \mathbb{E}_{q(\mathbf{x}_\sigma | \mathbf{x})} \left[\left\| s_\theta(\mathbf{x}_\sigma) - \nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma | \mathbf{x}) \right\|_2^2 \right] + const$$

Denoising Score Matching

Мы сами задаем $q(\mathbf{x}_\sigma|\mathbf{x})$:

$$q(\mathbf{x}_\sigma|\mathbf{x}) = \mathcal{N}(\mathbf{x}_\sigma|\mathbf{x}, \sigma^2 \mathbf{I})$$

В этом случае *score*-функция вычисляется аналитически:

$$\nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma|\mathbf{x}) = \nabla_{\mathbf{x}_\sigma} \left(-\frac{\|\mathbf{x}_\sigma - \mathbf{x}\|^2}{2\sigma^2} \right) = -\frac{\mathbf{x}_\sigma - \mathbf{x}}{\sigma^2} = -\frac{\boldsymbol{\epsilon}}{\sigma}$$

Подставим в нашу функцию потерь:

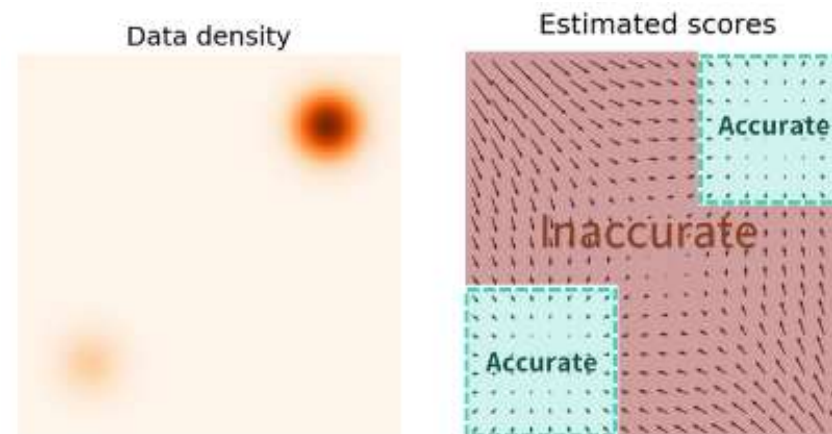
$$\frac{1}{2} \mathbb{E}_{p_{data}(\mathbf{x})} \mathbb{E}_{q(\mathbf{x}_\sigma|\mathbf{x})} \left[\left\| s_{\boldsymbol{\theta}}(\mathbf{x}_\sigma) - \nabla_{\mathbf{x}_\sigma} \log q(\mathbf{x}_\sigma|\mathbf{x}) \right\|^2 \right] = \frac{1}{2} \mathbb{E}_{p_{data}(\mathbf{x})} \mathbb{E}_{\mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\left\| s_{\boldsymbol{\theta}}(\mathbf{x} + \sigma \boldsymbol{\epsilon}) + \frac{\boldsymbol{\epsilon}}{\sigma} \right\|^2 \right] \rightarrow \min_{\boldsymbol{\theta}}$$

- Нейросеть $s_{\boldsymbol{\theta}}(\mathbf{x}_\sigma)$ учится напрямую предсказывать шум $\boldsymbol{\epsilon}$, который был добавлен к \mathbf{x}

Проблема *Denoising Score Matching*

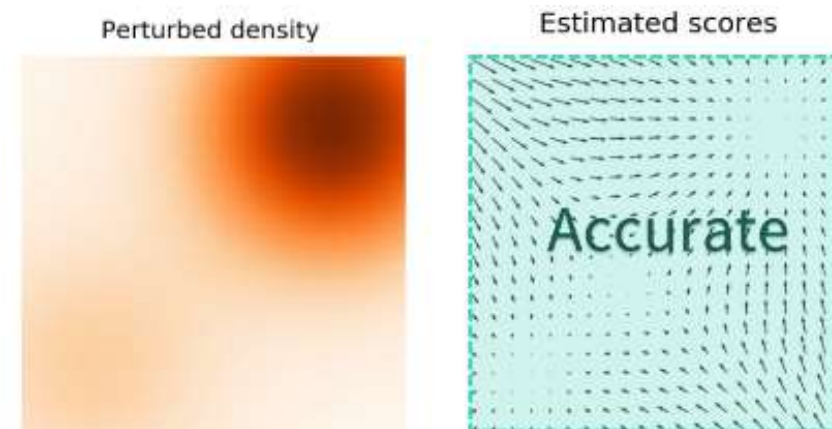
Если σ маленький:

- $q(\mathbf{x}_\sigma) \approx p_{data}(\mathbf{x})$
- Функция потерь *DSM* содержит члены $\frac{1}{\sigma}$ и $\frac{1}{\sigma^2}$, что приводит к высокой дисперсии градиентов
- Обучение нестабильно, динамика Ланжевена застревает в модах



Если σ большой:

- Функция потерь *DSM* с низкой дисперсией градиентов
- $q(\mathbf{x}_\sigma) \neq p_{data}(\mathbf{x})$ — сильно размытое распределение
- Модель хорошо работает в областях с низкой плотностью, но изучает искаженное распределение



Noise Conditioned Score Network

Идея:

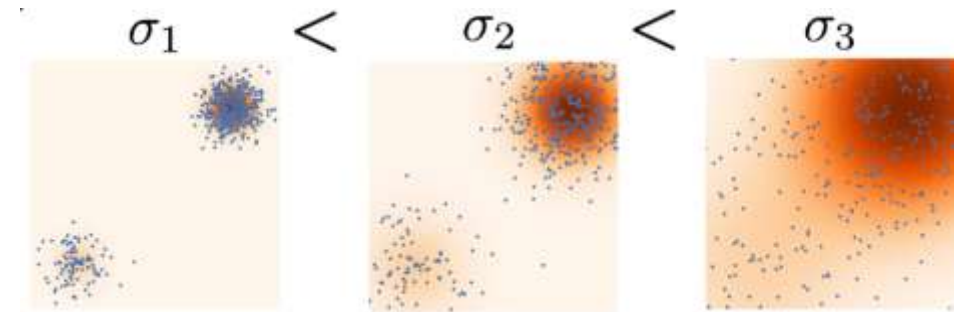
Будем одновременно использовать несколько уровней шума

- Определим последовательность уровней шума

$$\sigma_1 < \sigma_2 < \dots < \sigma_T:$$

- σ_1 — небольшой шум, $q(\mathbf{x}_1) \approx p_{data}(\mathbf{x})$
- σ_T — большой шум, $q(\mathbf{x}_T) \approx \mathcal{N}(0, \sigma_T^2 \mathbf{I})$

- Каждому σ_t будет соответствовать свое распределение $q(\mathbf{x}_t)$



Обучение *NCSN*

Будем обучать $s_{\theta, \sigma_t}(\mathbf{x}_t)$ на каждом уровне шума, используя взвешенную функцию потерь *DSM*:

$$\mathcal{L}_{\theta} = \sum_{t=1}^T \lambda(\sigma_t) \cdot \mathbb{E}_{p_{data}(\mathbf{x})} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x})} \left[\left\| s_{\theta, \sigma_t}(\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}) \right\|^2 \right],$$

где $\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}) = -\frac{\epsilon}{\sigma_t}$

Обычно веса $\lambda(\sigma_t)$ берут в виде σ_t^2 , чтобы сбалансировать члены с $\frac{1}{\sigma_t^2}$

$$\mathcal{L}_{\theta} = \sum_{t=1}^T \sigma_t^2 \cdot \mathbb{E}_{p_{data}(\mathbf{x})} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x})} \left[\left\| s_{\theta, \sigma_t}(\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}) \right\|^2 \right]$$

План

- *Diffusion Models*

- *Forward Diffusion Process*

- *Reverse Diffusion Process*

- *DDPM*

Diffusion Models

Идея диффузионных моделей

- *Sohl – Dickstein* (2015) задались вопросом как преобразовать одно распределение в другое

VAE:

- Кодировщик $q_\phi(\mathbf{z}|\mathbf{x})$ сжимает данные в латентный вектор \mathbf{z}
- Декодер $p_\theta(\mathbf{z}|\mathbf{x})$ восстанавливает \mathbf{x} из \mathbf{z}

Проблема:

Декодер $p_\theta(\mathbf{z}|\mathbf{x})$ должен за один шаг преобразовать вектор \mathbf{z} в сложный объект \mathbf{x}

Идея:

Попробуем вместо одного шага построить цепочку из T небольших шагов

Forward Diffusion Process

Forward Diffusion Process

Определим **прямой диффузионный процесс** как **марковскую цепь**, в которой мы добавляем шум к нашим данным:

$$q(\mathbf{x}_{0:T}) = q(\mathbf{x}_0) \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

- $q(\mathbf{x}_0)$ – распределение реальных данных
- $q(\mathbf{x}_t | \mathbf{x}_{t-1})$ – переходное распределение

Для описания добавления шума на каждом шаге t будем использовать **расписание шума** (**noise schedule**) β_t :

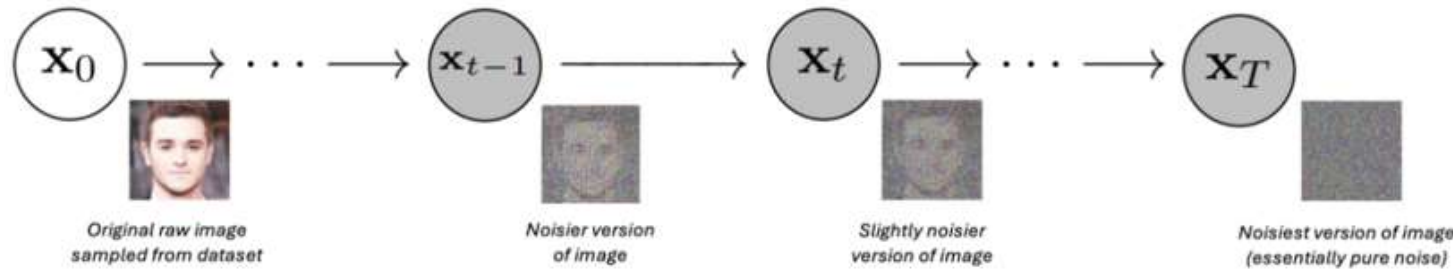
$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t | \sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

- $\sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1}$ – среднее значение
- $\beta_t \mathbf{I}$ – дисперсия

Forward Diffusion Process

Шаг получения \mathbf{x}_t из \mathbf{x}_{t-1} в явном виде:

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1} + \sqrt{\beta_t} \cdot \epsilon_t \quad \epsilon_t \sim \mathcal{N}(\mathbf{0}, I)$$



Иногда удобнее работать с параметром α_t , который имеет смысл **коэффициента сохранения сигнала**:

$$\alpha_t = 1 - \beta_t$$

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t | \sqrt{\alpha_t} \cdot \mathbf{x}_{t-1}, (1 - \alpha_t)I)$$

Forward Diffusion Process

- Можно ли сразу узнать зашумленное состояние \mathbf{x}_t , если нам известно начальное состояние \mathbf{x}_0 ?

Условное распределение $q(\mathbf{x}_t|\mathbf{x}_0)$, которое устанавливает связь между начальным состоянием \mathbf{x}_0 и любым последующим \mathbf{x}_t имеет вид:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}\left(\mathbf{x}_t \middle| \sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}\right)$$

- $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

Forward Diffusion Process

Подставим последовательно определение шагов диффузии:

$$\mathbf{x}_t = \sqrt{\alpha_t} \cdot \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \cdot \epsilon_t$$

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}} \cdot \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \cdot \epsilon_{t-1}$$

Получаем:

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \cdot (\sqrt{\alpha_{t-1}} \cdot \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \cdot \epsilon_{t-1}) + \sqrt{1 - \alpha_t} \cdot \epsilon_t = \\ &= \sqrt{\alpha_t \cdot \alpha_{t-1}} \cdot \mathbf{x}_{t-2} + \underbrace{\sqrt{\alpha_t(1 - \alpha_{t-1})} \cdot \epsilon_{t-1}}_{\mathbf{w}} + \sqrt{1 - \alpha_t} \cdot \epsilon_t \end{aligned}$$

$$\mathbb{E}[\mathbf{w}] = \mathbb{E}[\epsilon_{t-1}] + \mathbb{E}[\epsilon_t] = 0$$

$$\text{Var}[\mathbf{w}] = \alpha_t(1 - \alpha_{t-1}) \cdot \mathbf{I} + (1 - \alpha_t) \cdot \mathbf{I} = (1 - \alpha_t \cdot \alpha_{t-1}) \cdot \mathbf{I}$$

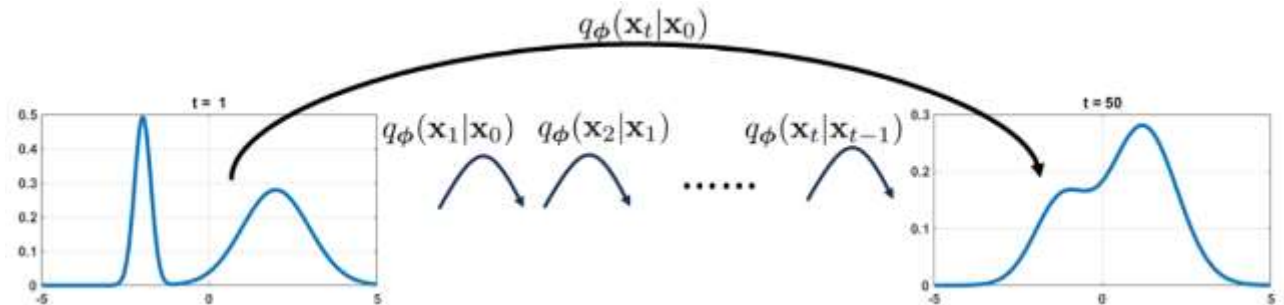
Forward Diffusion Process

Повторяя шаги по индукции, мы получим:

$$\mathbf{x}_t = \sqrt{\alpha_t \alpha_{t-1} \dots \alpha_1} \cdot \mathbf{x}_0 + \sqrt{1 - \alpha_t \alpha_{t-1} \dots \alpha_1} \cdot \epsilon_0 = \sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon_0$$

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$



Рассмотрим предел $t \rightarrow \infty$:

$$q(\mathbf{x}_t | \mathbf{x}_0) \rightarrow \mathcal{N}(\sqrt{0} \cdot \mathbf{x}_0, (1 - 0) \mathbf{I}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Forward Diffusion Process

Вспомним динамику Ланжевена:

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\eta}{2} \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}_t) + \sqrt{\eta} \epsilon_t \quad \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Прямой диффузионный процесс:

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1} + \sqrt{\beta_t} \cdot \epsilon_t$$

Воспользуемся приближением $\sqrt{1 - x} \approx 1 - \frac{x}{2}$ из ряда Тейлора для $x \ll 1$:

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1} + \sqrt{\beta_t} \epsilon_t \approx \left(1 - \frac{\beta_t}{2}\right) \cdot \mathbf{x}_{t-1} + \sqrt{\beta_t} \cdot \epsilon_t = \mathbf{x}_{t-1} + \frac{\beta_t}{2} (-\mathbf{x}_{t-1}) + \sqrt{\beta_t} \cdot \epsilon_t$$

Forward Diffusion Process

$$\begin{aligned}LD: \quad \mathbf{x}_{t+1} &= \mathbf{x}_t + \frac{\eta}{2} \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}_t) + \sqrt{\eta} \epsilon_t \\ FDP: \quad \mathbf{x}_t &\approx \mathbf{x}_{t-1} + \frac{\beta_t}{2} (-\mathbf{x}_{t-1}) + \sqrt{\beta_t} \cdot \epsilon_t\end{aligned}$$

Видим соответствие:

$$\beta_t = \eta$$

$$\nabla_{\mathbf{x}_{t-1}} \log p_{\theta}(\mathbf{x}_{t-1}) = -\mathbf{x}_{t-1}$$

Найдем *score* – функцию стандартного нормального распределения:

$$\nabla_{\mathbf{x}} \log \mathcal{N}(\mathbf{0}, I) = \nabla_{\mathbf{x}} \log C e^{-\frac{\mathbf{x}^2}{2}} = \nabla_{\mathbf{x}} \log C + \nabla_{\mathbf{x}} \log e^{-\frac{\mathbf{x}^2}{2}} = -\mathbf{x}$$

Forward Diffusion Process

$$LD: \mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\eta}{2} \nabla_{\mathbf{x}} \log p_{\theta}(\mathbf{x}_t) + \sqrt{\eta} \epsilon_t$$

$$FDP: \mathbf{x}_t \approx \mathbf{x}_{t-1} + \frac{\beta_t}{2} (-\mathbf{x}_{t-1}) + \sqrt{\beta_t} \cdot \epsilon_t$$

Вывод:

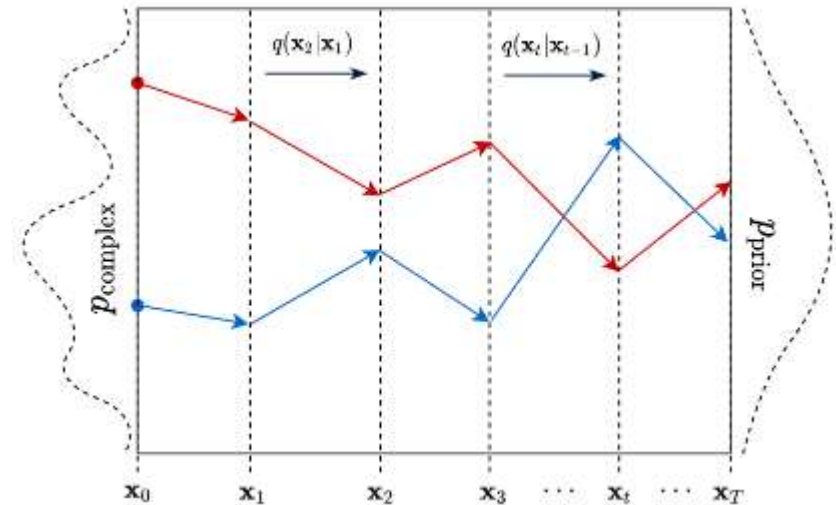
- Прямой диффузионный процесс – это динамика Ланжевена, которая тянет наши данные \mathbf{x} к распределению $\mathcal{N}(\mathbf{0}, I)$, используя его *score* – функцию
- Это даёт теоретическое обоснование, что \mathbf{x}_T сойдется к $\mathcal{N}(\mathbf{0}, I)$

Математика прямого процесса

Диффузия – это физический процесс перехода частиц из области высокой плотности в область низкой плотности

Прямой процесс:

- Берём $\mathbf{x}_0 \sim p_{data}(\mathbf{x})$
- Добавляем к ней шум $\mathbf{x}_t = \sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1} + \sqrt{\beta_t} \cdot \epsilon_t$
- После T шагов получим $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, I)$



- Прямой процесс фиксирован и не обучается
- Наша цель – попытаться обратить этот процесс

Reverse Diffusion Process

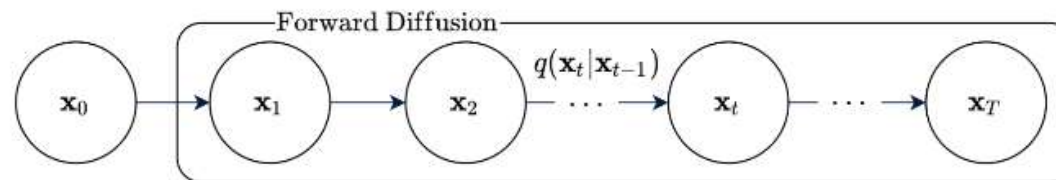
Diffusion Model as VAE

- Построим декодер p_θ , который будет обращать прямой диффузионный процесс q

Идея:

Рассмотрим весь диффузионный процесс как иерархический VAE:

- $\mathbf{x} = \mathbf{x}_0$ — наблюдаемые данные
- $\mathbf{z} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$ — скрытые переменные



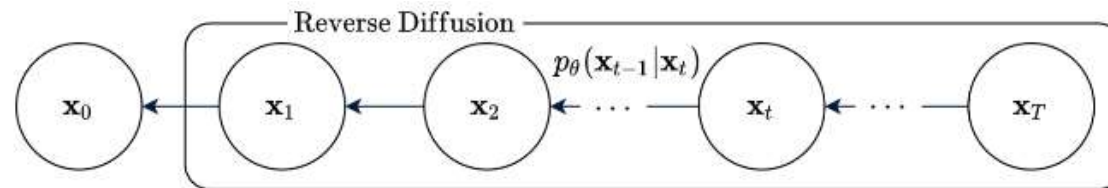
- Кодировщик q_ϕ — это наш прямой процесс q , который не является обучаемым
- В этом случае вариационное распределение $q_\phi(\mathbf{z}|\mathbf{x})$ имеет вид:

$$q(\mathbf{z}|\mathbf{x}) = q(\mathbf{x}_1, \dots, \mathbf{x}_T|\mathbf{x}_0) = q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

Diffusion Model as VAE

В VAE декодер $p_{\theta}(\mathbf{x}|\mathbf{z})$ – это шаг $\mathbf{z} \rightarrow \mathbf{x}$

В нашем случае $\mathbf{z} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$, поэтому $p_{\theta}(\mathbf{x}|\mathbf{z})$ – это последний шаг:



$$p_{\theta}(\mathbf{x}|\mathbf{z}) = p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)$$

Априорное распределение $p_{\theta}(\mathbf{z})$ – это вся остальная обратная цепь:

$$p_{\theta}(\mathbf{z}) = p_{\theta}(\mathbf{x}_{1:T}) = p_{\theta}(\mathbf{x}_T) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

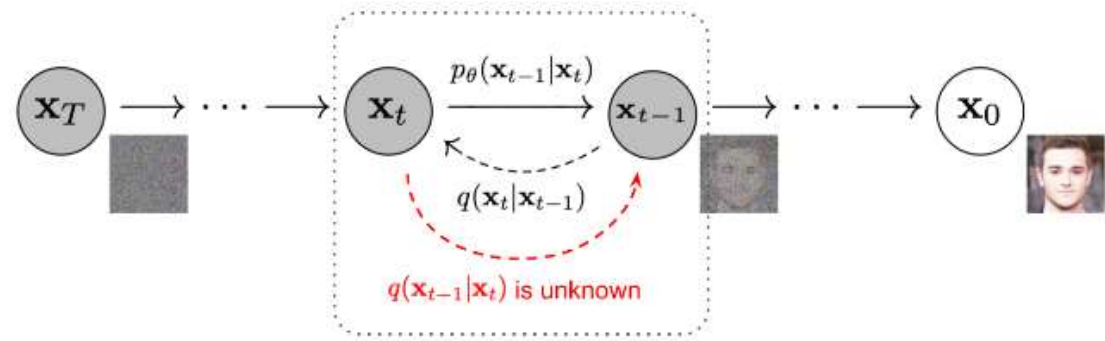
- **VAE**: простой *prior* $p(\mathbf{z})$ + сложный декодер $p_{\theta}(\mathbf{x}|\mathbf{z})$
- **Diffusion**: сложный *prior* $p_{\theta}(\mathbf{z})$ + простой декодер $p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)$

Diffusion Model as VAE

Мы хотим обучить $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$, чтобы он был похож на истинный обратный шаг $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$

Можем найти $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ по теореме Байеса:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}) \cdot q(\mathbf{x}_{t-1})}{q(\mathbf{x}_t)}$$



$q(\mathbf{x}_t)$ – маргинальное распределение всех возможных зашумленных изображений на шаге t :

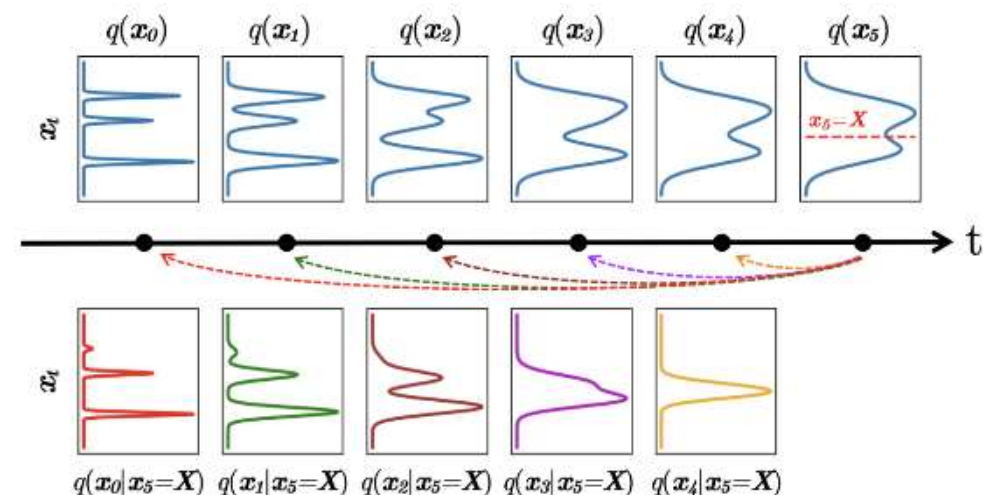
$$q(\mathbf{x}_t) = \int q(\mathbf{x}_t|\mathbf{x}_0) \cdot p_{data}(\mathbf{x}_0) d\mathbf{x}_0$$

Мы не можем вычислить этот интеграл из-за $p_{data}(\mathbf{x}_0)$

Теорема Феллера

Теорема (Feller, 1949):

Если β_t на каждом шаге достаточно мал, то обратный процесс $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ тоже будет гауссовским



Это объясняет почему мы должны брать много шагов диффузии ($T \approx 1000$)

Будем учить декодер p_θ предсказывать параметры \mathbf{x}_{t-1} , зная \mathbf{x}_t и t :

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\sigma}_\theta^2(\mathbf{x}_t, t))$$

ELBO for Diffusion Model

Как и в VAE, будем обучать модель, максимизируя *ELBO*:

$$\log p_{\theta}(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})} = \mathcal{L}_{\phi, \theta}(\mathbf{x}) \rightarrow \max_{\phi, \theta}$$

Подставим наши распределения:

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \log \frac{p_{\theta}(\mathbf{x}_0, \mathbf{x}_{1:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})}$$

- В числителе $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$, в знаменателе $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ — напоминает KL — дивергенцию
- Попробуем перевернуть $q(\mathbf{x}_t|\mathbf{x}_{t-1}) \rightarrow q(\mathbf{x}_{t-1}|\mathbf{x}_t)$

Условное ядро

Рассмотрим условное ядро $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) \cdot q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)}$$

- $q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = q(\mathbf{x}_t|\mathbf{x}_{t-1})$ – ядро перехода из \mathbf{x}_{t-1} в \mathbf{x}_t не зависит от \mathbf{x}_0
- $q(\mathbf{x}_{t-1}|\mathbf{x}_0)$ – ядро перехода из \mathbf{x}_0 в \mathbf{x}_{t-1}
- $q(\mathbf{x}_t|\mathbf{x}_0)$ – ядро перехода из \mathbf{x}_0 в \mathbf{x}_t

Подставим всё известное формулу:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{\mathcal{N}(\sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1}, \beta_t I) \cdot \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}} \cdot \mathbf{x}_0, (1 - \bar{\alpha}_{t-1})I)}{\mathcal{N}(\sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0, (1 - \bar{\alpha}_t)I)} = \mathcal{N}(\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \cdot I)$$

Условное ядро

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{\mathcal{N}(\sqrt{1 - \beta_t} \cdot \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \cdot \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}} \cdot \mathbf{x}_0, (1 - \bar{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(\sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})} = \mathcal{N}(\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \cdot \mathbf{I})$$

Параметры нового распределения:

$$\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \cdot \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)}{1 - \bar{\alpha}_t} \cdot \mathbf{x}_0$$

$$\tilde{\beta}_t = \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}$$

- Теперь у нас есть параметры распределения, которые должен предсказывать декодер $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$

ELBO for Diffusion Model

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0) = \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \cdot q(\mathbf{x}_t | \mathbf{x}_0)}{q(\mathbf{x}_{t-1} | \mathbf{x}_0)}$$

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})} = \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \log \frac{p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{\prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0)} =$$

$$\mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \log \frac{p(\mathbf{x}_T) \cdot p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1) \cdot \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_1 | \mathbf{x}_0) \cdot \prod_{t=2}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0)} = \mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \log \frac{p(\mathbf{x}_T) \cdot p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1) \cdot \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_1 | \mathbf{x}_0) \cdot \prod_{t=2}^T \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \cdot q(\mathbf{x}_t | \mathbf{x}_0)}{q(\mathbf{x}_{t-1} | \mathbf{x}_0)}} =$$

$$\mathbb{E}_{q(\mathbf{x}_{1:T} | \mathbf{x}_0)} \log \frac{p(\mathbf{x}_T) \cdot p_{\theta}(\mathbf{x}_0 | \mathbf{x}_1) \cdot \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)}{q(\mathbf{x}_T | \mathbf{x}_0) \cdot \prod_{t=2}^T q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)}$$

$$\prod_{t=2}^T \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \cdot q(\mathbf{x}_t | \mathbf{x}_0)}{q(\mathbf{x}_{t-1} | \mathbf{x}_0)} = \frac{q(\mathbf{x}_1 | \mathbf{x}_2, \mathbf{x}_0) \cdot q(\mathbf{x}_2 | \mathbf{x}_0)}{q(\mathbf{x}_1 | \mathbf{x}_0)} \cdot \frac{q(\mathbf{x}_2 | \mathbf{x}_3, \mathbf{x}_0) \cdot q(\mathbf{x}_3 | \mathbf{x}_0)}{q(\mathbf{x}_2 | \mathbf{x}_0)} \cdot \frac{q(\mathbf{x}_3 | \mathbf{x}_4, \mathbf{x}_0) \cdot q(\mathbf{x}_4 | \mathbf{x}_0)}{q(\mathbf{x}_3 | \mathbf{x}_0)}$$

ELBO for Diffusion Model

$$\mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \log \frac{p(\mathbf{x}_T) \cdot p_\theta(\mathbf{x}_0|\mathbf{x}_1) \cdot \prod_{t=2}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_T|\mathbf{x}_0) \cdot \prod_{t=2}^T q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} =$$

$$\mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[\log p_\theta(\mathbf{x}_0|\mathbf{x}_1) + \log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \sum_{t=2}^T \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right] =$$

$$\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) + \mathbb{E}_{q(\mathbf{x}_T|\mathbf{x}_0)} \log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_0)} + \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_t|\mathbf{x}_0)} \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} =$$

$$\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) - KL(q(\mathbf{x}_T|\mathbf{x}_0) || p(\mathbf{x}_T)) - \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} KL(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))$$

$$\mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_t|\mathbf{x}_0)} \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} = \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\mathbb{E}_{q(\mathbf{x}_{t-1}|\mathbf{x}_0)} \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)} \right]$$

ELBO for Diffusion Model

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} \log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)}_{\mathcal{L}_0} - \underbrace{KL(q(\mathbf{x}_T|\mathbf{x}_0) || p(\mathbf{x}_T))}_{\mathcal{L}_T} - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} KL(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{\mathcal{L}_t}$$

\mathcal{L}_0 – единственный шаг декодера $p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)$, отвечает за восстановление изображения \mathbf{x}_0 из \mathbf{x}_1 :

$$\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) = \log \mathcal{N}(\mathbf{x}_0 | \boldsymbol{\mu}_{\theta,t}(\mathbf{x}_1), \boldsymbol{\sigma}_{\theta,t}^2(\mathbf{x}_1)), \quad \mathbf{x}_1 \sim q(\mathbf{x}_1|\mathbf{x}_0)$$

\mathcal{L}_T – KL -дивергенция между двумя необучаемыми распределениями:

$$KL(q(\mathbf{x}_T|\mathbf{x}_0) || p(\mathbf{x}_T)) = KL\left(\mathcal{N}\left(\sqrt{\bar{\alpha}_T} \cdot \mathbf{x}_0, (1 - \bar{\alpha}_T)\right) || \mathcal{N}(\mathbf{0}, \mathbf{I})\right) = const$$

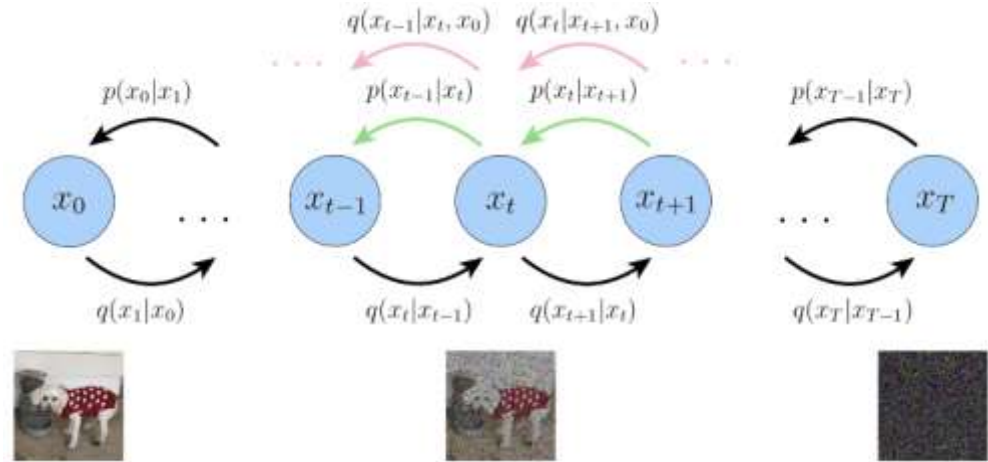
- \mathcal{L}_t – потеря на каждом шаге t
- Каждый член заставляет нейросеть p_{θ} быть похожим на истинное распределение q

ELBO for Diffusion Model

$$\mathcal{L}_t = \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} KL(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))$$

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}|\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\boldsymbol{\beta}}_t \cdot \mathbf{I})$$

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\mu}_{\theta,t}(\mathbf{x}_t), \boldsymbol{\sigma}_{\theta,t}^2(\mathbf{x}_t))$$



ELBO for Diffusion Model

- Не будем обучать дисперсию $\sigma_{\theta,t}^2$
- Зафиксируем $\sigma_{\theta,t}^2$, приравняв к известной дисперсии $\tilde{\beta}$:

$$\sigma_{\theta,t}^2 = \tilde{\beta}_t \cdot I$$

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\mu}_{\theta,t}(\mathbf{x}_t), \sigma_{\theta,t}^2(\mathbf{x}_t)) = \mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\mu}_{\theta,t}(\mathbf{x}_t), \tilde{\beta}_t \cdot I)$$

KL — дивергенция между двумя гауссианами с одинаковой дисперсией – это *MSE* между их средними:

$$\begin{aligned}\mathcal{L}_t &= \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} KL(\mathcal{N}(\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t I) || \mathcal{N}(\boldsymbol{\mu}_{\theta,t}(\mathbf{x}_t), \tilde{\beta}_t \cdot I)) = \\ &= \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{2\tilde{\beta}_t} \left\| \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_{\theta,t}(\mathbf{x}_t) \right\|^2 \right]\end{aligned}$$

Training

Обучение:

- Берем объект $\mathbf{x}_0 \sim p_{data}(\mathbf{x})$
- Получаем зашумленный объект $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon$, $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$
- Вычисляем лосс:

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} \log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) - KL(q(\mathbf{x}_T|\mathbf{x}_0)||p(\mathbf{x}_T)) - \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{2\tilde{\beta}_t} \left\| \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_{\theta, t}(\mathbf{x}_t) \right\|^2 \right]$$

Inference

Сэмплинг:

- Сэмплируем $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, I)$
- Делаем обратный проход $range(T, 1)$:

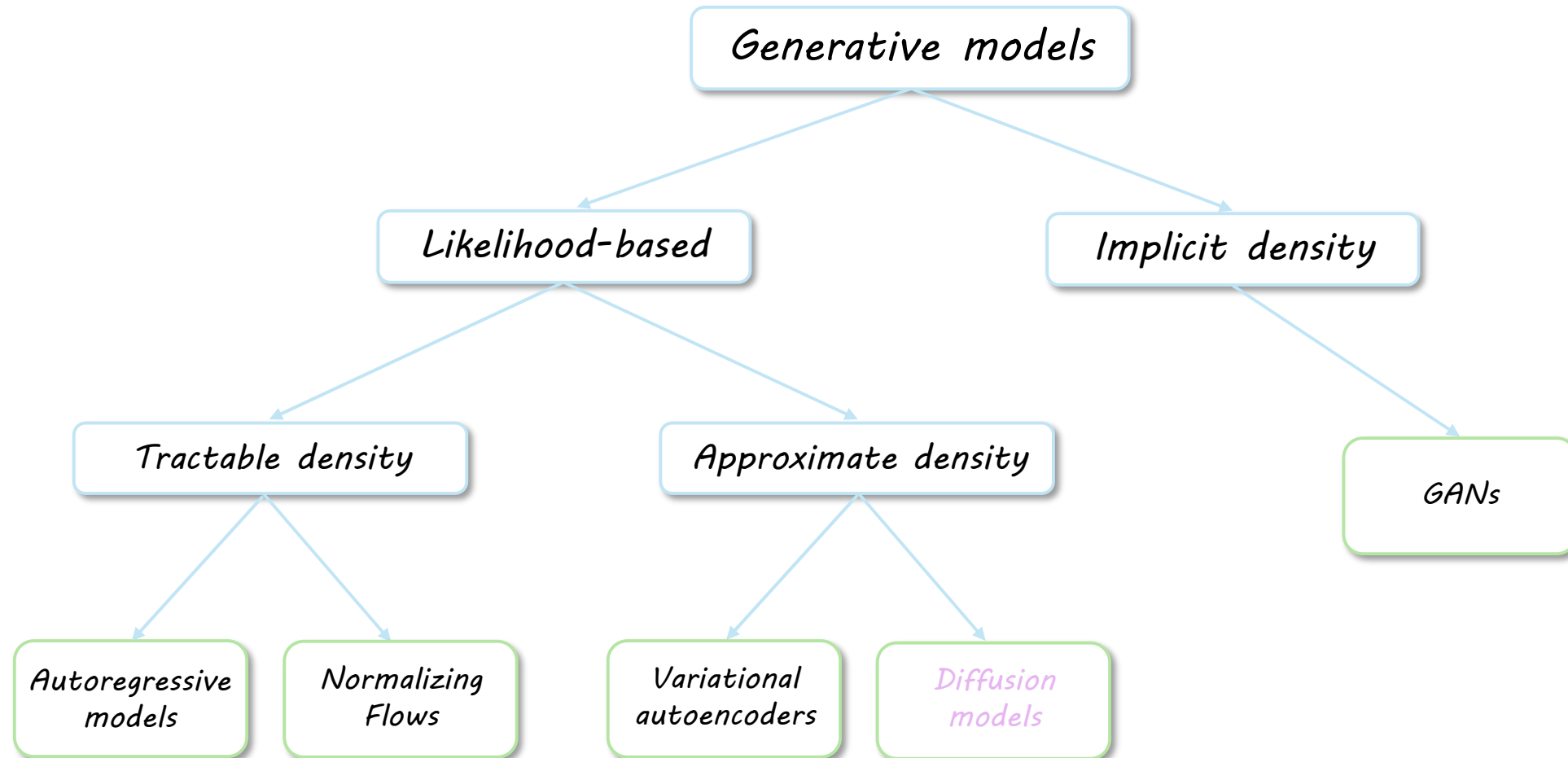
Подаем \mathbf{x}_t и t в обученную нейросеть, получаем $\boldsymbol{\mu}_{\theta,t}(\mathbf{x}_t)$

Сэмплируем \mathbf{x}_{t-1} из обратного процесса $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$:

$$\mathbf{x}_{t-1} = \boldsymbol{\mu}_{\theta,t}(\mathbf{x}_t) + \sqrt{\tilde{\beta}_t} \cdot \boldsymbol{\epsilon}$$

Denoising Diffusion Probabilistic Model

Зоопарк генеративных моделей



Reparameterization

Проблема:

Обучать нейросеть предсказывать μ_θ — не всегда стабильно

Существует линейная зависимость между \mathbf{x}_0 , \mathbf{x}_t и ϵ :

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon \qquad \mathbf{x}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon)$$

Подставим в формулу для $\tilde{\mu}_t$:

$$\begin{aligned} \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \cdot \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)}{1 - \bar{\alpha}_t} \cdot \mathbf{x}_0 = \\ &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \cdot \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)}{1 - \bar{\alpha}_t} \cdot \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon) = \frac{1}{\sqrt{\alpha_t}} \cdot \mathbf{x}_t - \frac{(1 - \alpha_t)}{\sqrt{\alpha_t}(1 - \bar{\alpha}_t)} \cdot \epsilon \end{aligned}$$

Reparameterization

$$\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}} \cdot \mathbf{x}_t - \frac{(1 - \alpha_t)}{\sqrt{\alpha_t(1 - \bar{\alpha}_t)}} \cdot \boldsymbol{\epsilon}$$

Сделаем репараметризацию и заставим нейросеть иметь такую же форму:

$$\boldsymbol{\mu}_{\theta,t}(\mathbf{x}_t) = \frac{1}{\sqrt{\alpha_t}} \cdot \mathbf{x}_t - \frac{(1 - \alpha_t)}{\sqrt{\alpha_t(1 - \bar{\alpha}_t)}} \cdot \boldsymbol{\epsilon}_{\theta,t}(\mathbf{x}_t)$$

Теперь наша модель будет предсказывать шум $\boldsymbol{\epsilon}_{\theta,t}(\mathbf{x}_t)$, а не среднее $\boldsymbol{\mu}_{\theta,t}(\mathbf{x}_t)$:

$$\mathcal{L}_t = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, I)} \left[\frac{(1 - \alpha_t)^2}{2\tilde{\beta}_t\alpha_t(1 - \bar{\alpha}_t)} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta,t}(\mathbf{x}_t) \right\|^2 \right] =$$
$$\mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, I)} \left[\frac{(1 - \alpha_t)^2}{2\tilde{\beta}_t\alpha_t(1 - \bar{\alpha}_t)} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta,t}(\sqrt{\alpha_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \boldsymbol{\epsilon}) \right\|^2 \right]$$

Функция потерь *DDPM*

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} \log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1) - KL(q(\mathbf{x}_T|\mathbf{x}_0)||p(\mathbf{x}_T)) - \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} KL(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{\mathcal{L}_t}$$

$$\mathcal{L}_t = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, I)} \left[\frac{(1 - \alpha_t)^2}{2\tilde{\beta}_t \alpha_t (1 - \bar{\alpha}_t)} \left\| \epsilon - \epsilon_{\theta, t}(\sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon) \right\|^2 \right]$$

Отбросим весовой коэффициент:

$$\mathcal{L}_t = \mathbb{E}_{t \sim U(2, T)} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, I)} \left[\left\| \epsilon - \epsilon_{\theta, t}(\sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon) \right\|^2 \right]$$

Training

Обучение:

- Берем объект $\mathbf{x}_0 \sim p_{data}(\mathbf{x})$
- Выбираем номер шага $t \sim U(2, T)$ и шум $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$
- Получаем зашумленный объект $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \cdot \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \cdot \epsilon$
- Подаем \mathbf{x}_t и t в нейросеть
- Вычисляем MSE между реальным ϵ и предсказанным шумом $\epsilon_{\theta, t}$

Inference

Сэмплинг:

- Сэмплируем $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, I)$
- Делаем обратный проход $range(T, 1)$:

Подаем \mathbf{x}_t и t в обученную нейросеть, получаем $\epsilon_{\theta,t}(\mathbf{x}_t)$

Вычисляем среднее:

$$\mu_{\theta,t}(\mathbf{x}_t) = \frac{1}{\sqrt{\alpha_t}} \cdot \mathbf{x}_t - \frac{(1 - \alpha_t)}{\sqrt{\alpha_t(1 - \bar{\alpha}_t)}} \cdot \epsilon_{\theta,t}(\mathbf{x}_t)$$

Сэмплируем \mathbf{x}_{t-1} из обратного процесса $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$:

$$\mathbf{x}_{t-1} = \mu_{\theta,t}(\mathbf{x}_t) + \sqrt{\tilde{\beta}_t} \cdot \epsilon$$

DDPM Samples



Спасибо за внимание!