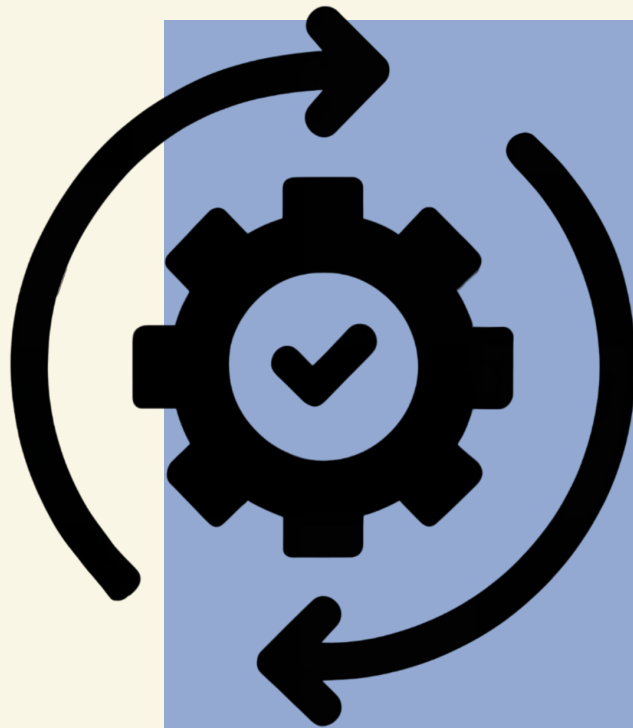


# Эффективные системы машинного обучения



Лекция 10

Преподаватель

29 ноября 2024

Представление диффузионных  
моделей. Решения ОДУ

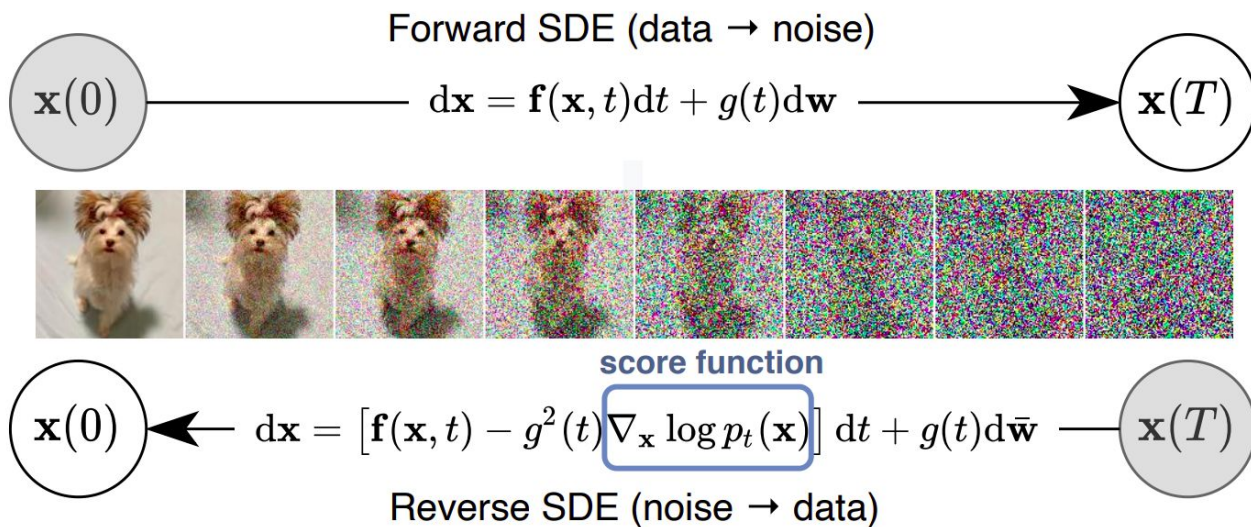
Оганов Александр

ВМК МГУ

# Напоминание

Из чего стоят диффузионные модели?

- Процесс разрушения (прямой)
- Процесс генерации (обратный)



# Непрерывные модели

Далее мы будем рассматривать только модели с непрерывным временем. В качестве процесса зашумления мы можем взять VP-SDE или VE-SDE

$$\theta^* = \arg \min_{\theta} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{\mathbf{x}(0)} \mathbb{E}_{\mathbf{x}(t) | \mathbf{x}(0)} \left[ \left\| \mathbf{s}_{\theta}(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_{0t}(\mathbf{x}(t) | \mathbf{x}(0)) \right\|_2^2 \right] \right\}$$

$$p_{0t}(\mathbf{x}(t) | \mathbf{x}(0)) = \begin{cases} \mathcal{N}(\mathbf{x}(t); \mathbf{x}(0), [\sigma^2(t) - \sigma^2(0)]\mathbf{I}), & \text{(VE SDE)} \\ \mathcal{N}(\mathbf{x}(t); \mathbf{x}(0)e^{-\frac{1}{2} \int_0^t \beta(s) ds}, \mathbf{I} - \mathbf{I}e^{-\int_0^t \beta(s) ds}) & \text{(VP SDE)} \\ \mathcal{N}(\mathbf{x}(t); \mathbf{x}(0)e^{-\frac{1}{2} \int_0^t \beta(s) ds}, [1 - e^{-\int_0^t \beta(s) ds}]^2 \mathbf{I}) & \text{(sub-VP SDE)} \end{cases}$$

# Параметризация

Можно учить:

- предсказание шума
- score функцию
- предсказание чистого объекта  $x_0$

$$p_{0t}(x(t) \mid x(0)) := \mathcal{N}(x(t) \mid x(0)\sqrt{\bar{\alpha}_t}, (1 - \bar{\alpha}_t)I)$$
$$\log p_{0t}(x(t) \mid x(0)) = \text{const} - \frac{1}{2} \frac{(x(t) - \sqrt{\bar{\alpha}_t}x(0))^2}{1 - \bar{\alpha}_t}$$
$$\nabla_x \log p_{0t}(x(t) \mid x(0)) = - \frac{x(t) - \sqrt{\bar{\alpha}_t}x(0)}{1 - \bar{\alpha}_t}$$
$$\nabla_x \log p_{0t}(x(t) \mid x(0)) = - \frac{\epsilon(x(t), x(0))}{\sqrt{1 - \bar{\alpha}_t}}$$

# Как генерировать?

Для генерации необходимо решить либо обратный СДУ, либо соответствующее ОДУ

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t) d\bar{\mathbf{w}}$$

$$d\mathbf{x} = \left[ \mathbf{f}(\mathbf{x}, t) - \frac{1}{2} g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt$$

# Открытые вопросы

- 1) Как решать ОДУ/СДУ учитывая условия задачи?
- 2) Как же правильно задавать лосс?
- 3) Как выбрать  $f(x,t)$  и  $g(t)$ ?
- 4) Как не сломать себе голову, когда читаешь статьи и увидеть единый формализм

# Как повысить качество?

Обычно, новые модели учат достаточно плохо и неправильно: выбирают не ту архитектуру, оптимизируют не тот лосс и тд.

Какие гипер-параметры у нас есть?

- Зашумление
- Вес в лоссе
- Архитектура

Как их выбирать?

$$\theta^* = \arg \min_{\theta} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{\mathbf{x}(0)} \mathbb{E}_{\mathbf{x}(t)|\mathbf{x}(0)} \left[ \left\| \mathbf{s}_{\theta}(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_{0t}(\mathbf{x}(t) | \mathbf{x}(0)) \right\|_2^2 \right] \right\}$$

# Ответ на все вопросы

Для начала, стоит подумать, как записать все в едином виде, поэтому обратимся к статье Elucidating the Design Space of Diffusion-Based Generative Models (EDM) (<https://arxiv.org/abs/2206.00364>).

1) Как влияет СДУ на переходные вероятности?

$$d\mathbf{x} = f(t) \mathbf{x} dt + g(t) d\omega_t$$

$$s(t) = \exp\left(\int_0^t f(\xi) d\xi\right), \quad \text{and} \quad \sigma(t) = \sqrt{\int_0^t \frac{g(\xi)^2}{s(\xi)^2} d\xi}$$

$$p_{0t}(\mathbf{x}(t) \mid \mathbf{x}(0)) = \mathcal{N}(\mathbf{x}(t); s(t) \mathbf{x}(0), s(t)^2 \sigma(t)^2 \mathbf{I})$$



# Обратный процесс

Зная прямой процесс, мы можем найти и обратный. Попробуем  
выписать обратный процесс в общем виде используя определения  $f(t)$  и  
 $g(t)$

$$f(t) = \dot{s}(t)/s(t) \quad g(t) = s(t) \sqrt{2 \dot{\sigma}(t) \sigma(t)}$$

$$d\mathbf{x} = \left[ [f(t)] \mathbf{x} - \frac{1}{2} [g(t)]^2 \nabla_{\mathbf{x}} \log p(\mathbf{x}/s(t); \sigma(t)) \right] dt \quad (35)$$

$$= \left[ [\dot{s}(t)/s(t)] \mathbf{x} - \frac{1}{2} \left[ s(t) \sqrt{2 \dot{\sigma}(t) \sigma(t)} \right]^2 \nabla_{\mathbf{x}} \log p(\mathbf{x}/s(t); \sigma(t)) \right] dt \quad (36)$$

$$= \left[ [\dot{s}(t)/s(t)] \mathbf{x} - \frac{1}{2} \left[ 2 s(t)^2 \dot{\sigma}(t) \sigma(t) \right] \nabla_{\mathbf{x}} \log p(\mathbf{x}/s(t); \sigma(t)) \right] dt \quad (37)$$

$$= \left[ \frac{\dot{s}(t)}{s(t)} \mathbf{x} - s(t)^2 \dot{\sigma}(t) \sigma(t) \nabla_{\mathbf{x}} \log p\left(\frac{\mathbf{x}}{s(t)}; \sigma(t)\right) \right] dt. \quad (38)$$

# Итог

Нам достаточно задать  $s(t)$  и  $\sigma(t)$  и мы сможем получить рассчитать переходные вероятности и обратный процесс!

Если мы рассматриваем VE-SDE, то  $s(t) = 1$  и получаем:

$$dx = -\dot{\sigma}(t) \sigma(t) \nabla_x \log p(x; \sigma(t)) dt$$

# Как задать все процессы?

	VP [49]	VE [49]	iDDPM [37] + DDIM [47]	Ours (“EDM”)
<b>Sampling (Section 3)</b>				
ODE solver	Euler	Euler	Euler	2 <sup>nd</sup> order Heun
Time steps $t_{i < N}$	$1 + \frac{i}{N-1}(\epsilon_s - 1)$	$\sigma_{\max}^2 (\sigma_{\min}^2 / \sigma_{\max}^2)^{\frac{i}{N-1}}$	$u_{\lfloor j_0 + \frac{M-1-j_0}{N-1}i + \frac{1}{2} \rfloor}$ , where $u_M = 0$ $u_{j-1} = \sqrt{\frac{u_j^2 + 1}{\max(\bar{\alpha}_{j-1}/\bar{\alpha}_j, C_1)} - 1}$	$(\sigma_{\max}^{\frac{1}{\rho}} + \frac{i}{N-1}(\sigma_{\min}^{\frac{1}{\rho}} - \sigma_{\max}^{\frac{1}{\rho}}))^{\rho}$
Schedule $\sigma(t)$	$\sqrt{e^{\frac{1}{2}\beta_d t^2 + \beta_{\min} t} - 1}$	$\sqrt{t}$	$t$	$t$
Scaling $s(t)$	$1/\sqrt{e^{\frac{1}{2}\beta_d t^2 + \beta_{\min} t}}$	1	1	1
<b>Network and preconditioning (Section 5)</b>				
Architecture of $F_{\theta}$	DDPM++	NCSN++	DDPM	(any)
Skip scaling $c_{\text{skip}}(\sigma)$	1	1	1	$\sigma_{\text{data}}^2 / (\sigma^2 + \sigma_{\text{data}}^2)$
Output scaling $c_{\text{out}}(\sigma)$	$-\sigma$	$\sigma$	$-\sigma$	$\sigma \cdot \sigma_{\text{data}} / \sqrt{\sigma_{\text{data}}^2 + \sigma^2}$
Input scaling $c_{\text{in}}(\sigma)$	$1/\sqrt{\sigma^2 + 1}$	1	$1/\sqrt{\sigma^2 + 1}$	$1/\sqrt{\sigma^2 + \sigma_{\text{data}}^2}$
Noise cond. $c_{\text{noise}}(\sigma)$	$(M-1)\sigma^{-1}(\sigma)$	$\ln(\frac{1}{2}\sigma)$	$M-1 - \arg \min_j  u_j - \sigma $	$\frac{1}{4} \ln(\sigma)$
<b>Training (Section 5)</b>				
Noise distribution	$\sigma^{-1}(\sigma) \sim \mathcal{U}(\epsilon_t, 1)$	$\ln(\sigma) \sim \mathcal{U}(\ln(\sigma_{\min}), \ln(\sigma_{\max}))$	$\sigma = u_j, j \sim \mathcal{U}\{0, M-1\}$	$\ln(\sigma) \sim \mathcal{N}(P_{\text{mean}}, P_{\text{std}}^2)$
Loss weighting $\lambda(\sigma)$	$1/\sigma^2$	$1/\sigma^2$	$1/\sigma^2$ (note: *)	$(\sigma^2 + \sigma_{\text{data}}^2) / (\sigma \cdot \sigma_{\text{data}})^2$
<b>Parameters</b>				
	$\beta_d = 19.9, \beta_{\min} = 0.1$	$\sigma_{\min} = 0.02$	$\bar{\alpha}_j = \sin^2(\frac{\pi}{2} \frac{j}{M(C_2+1)})$	$\sigma_{\min} = 0.002, \sigma_{\max} = 80$
	$\epsilon_s = 10^{-3}, \epsilon_t = 10^{-5}$	$\sigma_{\max} = 100$	$C_1 = 0.001, C_2 = 0.008$	$\sigma_{\text{data}} = 0.5, \rho = 7$
	$M = 1000$		$M = 1000, j_0 = 8^{\dagger}$	$P_{\text{mean}} = -1.2, P_{\text{std}} = 1.2$

\* iDDPM also employs a second loss term  $L_{\text{vib}}$

<sup>†</sup> In our tests,  $j_0 = 8$  yielded better FID than  $j_0 = 0$  used by iDDPM

# Параметризация

Мы знаем, что нейронные сети очень важно хорошо параметризовать. Например, если данные и таргет нормализованные, то модель будет проще учить.

Далее мы считаем, что учим диффузионную модель на предсказание  $x_0$  в следующем виде:

$$D_{\theta}(\mathbf{x}; \sigma) = c_{\text{skip}}(\sigma) \mathbf{x} + c_{\text{out}}(\sigma) F_{\theta}(c_{\text{in}}(\sigma) \mathbf{x}; c_{\text{noise}}(\sigma))$$

Обучаем минимизируя:

$$\mathbb{E}_{\sigma, \mathbf{y}, \mathbf{n}} [\lambda(\sigma) \|D(\mathbf{y} + \mathbf{n}; \sigma) - \mathbf{y}\|_2^2]$$

$$\sigma \sim p_{\text{train}}, \mathbf{y} \sim p_{\text{data}}, \text{ and } \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

# Как учить?

Все гипер параметры авторы либо находили из перебора (подробнее на семинаре) или из теории.

$$\mathbb{E}_{\sigma, \mathbf{y}, \mathbf{n}} \left[ \lambda(\sigma) \|D(\mathbf{y} + \mathbf{n}; \sigma) - \mathbf{y}\|_2^2 \right]$$

$$\sigma \sim p_{\text{train}}, \mathbf{y} \sim p_{\text{data}}, \text{ and } \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$\mathbb{E}_{\sigma, \mathbf{y}, \mathbf{n}} \left[ \underbrace{\lambda(\sigma) c_{\text{out}}(\sigma)^2}_{\text{effective weight}} \left\| \underbrace{F_{\theta}(c_{\text{in}}(\sigma) \cdot (\mathbf{y} + \mathbf{n}); c_{\text{noise}}(\sigma))}_{\text{network output}} - \underbrace{\frac{1}{c_{\text{out}}(\sigma)} (\mathbf{y} - c_{\text{skip}}(\sigma) \cdot (\mathbf{y} + \mathbf{n}))}_{\text{effective training target}} \right\|_2^2 \right]$$

$$\lambda(\sigma) = 1/c_{\text{out}}(\sigma)^2$$

# Итоги

Мы научились:

- задавать зашумления в общем виде
- выбирать правильную параметризацию
- эффективно обучать

Мы не умеем:

- Правильно и эффективно решать ОДУ/СДУ

# Как инферить?

Мы умеем круто обучать диффузионные модели, но все еще не умеем быстро генерировать качественные изображения.

Есть следующие подходы:

- Дистилляции (<https://arxiv.org/abs/2303.01469>)
- Кэширования
- Квантизации
- Солверы ОДУ/СДУ
- есть и другие, но все основные направления современной науке описаны выше

# Солверы

Одним из современных направлений в диффузионных моделях это изучение, оптимизация и разработка солверов.

Что такое солвер?

Солвер – метод/алгоритм, который решает конкретную численную задачу. Например, метод Эйлера для решения ОДУ

Зачем нужны солверы?

Чем более точное решение ОДУ/СДУ мы получим, тем лучше будут изображения, при этом часто методы 2 порядка при меньшем числе шагов получают решения сильно лучше, чем методы 1 порядка.



# Солверы

Уже существует много работ о солверах, но мы поговорим про основные:

- DDIM (метод придуманный для дискретных моделей)
- DPM-Solver (солвер для решения ОДУ, которые возникают в диффузионках) (<https://arxiv.org/abs/2206.00927>)

# Как решать ОДУ?

В статье EDM также предложен метод решения ОДУ 2 порядка Heun, но мы посмотрим чуть с другой стороны.

---

**Algorithm 1** Deterministic sampling using Heun's 2<sup>nd</sup> order method with arbitrary  $\sigma(t)$  and  $s(t)$ .

---

```

1: procedure HEUNAMPLER( $D_\theta(\mathbf{x}; \sigma)$ ,  $\sigma(t)$ ,  $s(t)$ ,  $t_{i \in \{0, \dots, N\}}$ )
2:   sample  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \sigma^2(t_0) s^2(t_0) \mathbf{I})$                                 ▷ Generate initial sample at  $t_0$ 
3:   for  $i \in \{0, \dots, N - 1\}$  do                                              ▷ Solve Eq. 4 over  $N$  time steps
4:      $\mathbf{d}_i \leftarrow \left( \frac{\dot{\sigma}(t_i)}{\sigma(t_i)} + \frac{\dot{s}(t_i)}{s(t_i)} \right) \mathbf{x}_i - \frac{\dot{\sigma}(t_i) s(t_i)}{\sigma(t_i)} D_\theta \left( \frac{\mathbf{x}_i}{s(t_i)}; \sigma(t_i) \right)$   ▷ Evaluate  $d\mathbf{x}/dt$  at  $t_i$ 
5:      $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + (t_{i+1} - t_i) \mathbf{d}_i$                                 ▷ Take Euler step from  $t_i$  to  $t_{i+1}$ 
6:     if  $\sigma(t_{i+1}) \neq 0$  then                                              ▷ Apply 2nd order correction unless  $\sigma$  goes to zero
7:        $\mathbf{d}'_i \leftarrow \left( \frac{\dot{\sigma}(t_{i+1})}{\sigma(t_{i+1})} + \frac{\dot{s}(t_{i+1})}{s(t_{i+1})} \right) \mathbf{x}_{i+1} - \frac{\dot{\sigma}(t_{i+1}) s(t_{i+1})}{\sigma(t_{i+1})} D_\theta \left( \frac{\mathbf{x}_{i+1}}{s(t_{i+1})}; \sigma(t_{i+1}) \right)$   ▷ Eval.  $d\mathbf{x}/dt$  at  $t_{i+1}$ 
8:        $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + (t_{i+1} - t_i) \left( \frac{1}{2} \mathbf{d}_i + \frac{1}{2} \mathbf{d}'_i \right)$   ▷ Explicit trapezoidal rule at  $t_{i+1}$ 
9:   return  $\mathbf{x}_N$                                                                 ▷ Return noise-free sample at  $t_N$ 

```

---

# Солверы или как решать ОДУ

*Forward process*

$$q(x_t | x_0) = \mathcal{N}(x_t | \alpha_t x_0, \sigma_t^2 I)$$

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

$$f(t) = \frac{\dot{\alpha}_t}{\alpha_t}, \quad g^2(t) = 2\sigma_t \dot{\sigma}_t - \frac{2\sigma_t^2 \dot{\alpha}_t}{\alpha_t}$$

*Backward process*

$$dx_t = (f(t)x_t - g^2(t)\nabla_x \log p(x_t, t))dt + g(t)dW_t$$

$$\frac{dx_t}{dt} = f(t)x_t - \frac{1}{2}g^2(t)\nabla_x \log p(x_t, t)$$

# Качество генерации

Мы знаем: чем точнее решаем ОДУ/СДУ, тем выше будет качество изображения (следует из теории). Вопрос ли в том, что решать?

# Качество генерации

Мы знаем: чем точнее решаем ОДУ/СДУ, тем выше будет качество изображения (следует из теории). Вопрос ли в том, что решать?

Ответ: Будем решать вариант с ОДУ, а не СДУ, так как при больших шагах дисперсия при  $dW_t$  слишком высокая, что будет мешать сходимости. Запустим на ОДУ численный метод и все готово, например Рунге–Кутта. Правда же???

# Напоминание с курса ЧМов

Пусть задан ОДУ

$$\mathbf{x}'(t) = \mathbf{f}(t, \mathbf{x}(t), \theta)$$

Обычно такие дифференциальные уравнения называют линейными с переменными коэффициентами. Для них существует аналитическое решение, например, с помощью метода вариации постоянной

$$\mathbf{x}_t = e^{\int_s^t \mathbf{f}(\tau) d\tau} \mathbf{x}_s + \int_s^t \left( e^{\int_\tau^t \mathbf{f}(r) dr} \frac{g^2(\tau)}{2\sigma_\tau} \epsilon_\theta(\mathbf{x}_\tau, \tau) \right) d\tau$$

# Что мы получили?

Если посмотрим внимательней, то некоторые коэффициенты мы можем рассчитать аналитически. То есть использовать знания о задаче, а не кидать все в готовый солвер! (или же...)

Попробуем упростить задачу и уменьшить число сложных букв.

$$\mathbf{x}_t = e^{\int_s^t f(\tau) d\tau} \mathbf{x}_s + \int_s^t \left( e^{\int_\tau^t f(r) dr} \frac{g^2(\tau)}{2\sigma_\tau} \boldsymbol{\epsilon}_\theta(\mathbf{x}_\tau, \tau) \right) d\tau$$

# Полезные формулы

$$\log SNR = \log \frac{\alpha_t^2}{\sigma_t^2}$$

$$\lambda_t := \log \frac{\alpha_t}{\sigma_t}, \quad \frac{d\lambda_t}{dt} = \frac{\dot{\alpha}_t}{\alpha_t} - \frac{\dot{\sigma}_t}{\sigma_t}$$

$$\log p_\theta(x_t, t) = -\frac{\epsilon_\theta(x_t, t)}{\sigma_t}$$



# Чуть упростим

$$\mathbf{x}_t = e^{\int_s^t f(\tau) d\tau} \mathbf{x}_s + \int_s^t \left( e^{\int_\tau^t f(r) dr} \frac{g^2(\tau)}{2\sigma_\tau} \boldsymbol{\epsilon}_\theta(\mathbf{x}_\tau, \tau) \right) d\tau$$

$$\int f(t) dt = \int \frac{d \log \alpha_t}{dt} dt = \log \alpha_t + C$$

$$g^2(t) = 2\sigma_t \dot{\sigma}_t - \frac{2\sigma_t^2 \dot{\alpha}_t}{\alpha_t} = 2\sigma_t^2 \left( \frac{\dot{\sigma}_t}{\sigma_t} - \frac{\dot{\alpha}_t}{\alpha_t} \right) = -2\sigma_t^2 \frac{d\lambda_t}{dt}$$

$$\mathbf{x}_t = \frac{\alpha_t}{\alpha_s} \mathbf{x}_s - \alpha_t \int_s^t \left( \frac{d\lambda_\tau}{d\tau} \right) \frac{\sigma_\tau}{\alpha_\tau} \boldsymbol{\epsilon}_\theta(\mathbf{x}_\tau, \tau) d\tau$$

# Еще проще

Сделаем замену переменных

$$\mathbf{x}_t = \frac{\alpha_t}{\alpha_s} \mathbf{x}_s - \alpha_t \int_s^t \left( \frac{d\lambda_\tau}{d\tau} \right) \frac{\sigma_\tau}{\alpha_\tau} \boldsymbol{\epsilon}_\theta(\mathbf{x}_\tau, \tau) d\tau$$

$$\mathbf{x}_t = \frac{\alpha_t}{\alpha_s} \mathbf{x}_s - \alpha_t \int_{\lambda_s}^{\lambda_t} e^{-\lambda} \hat{\boldsymbol{\epsilon}}_\theta(\hat{\mathbf{x}}_\lambda, \lambda) d\lambda$$

# Что получили?

Мы аналитически можем подсчитать коэффициент перед линейной частью, что раньше мы кидали в солвер не думая. Причем перед этим коэффициентом у нас  $e^{\int \dots}$ , так что ошибка накапливалась очень весомая.

Теперь мы свели задачу к оценке интеграла, который известен в литературе как exponentially weighted integral и сводиться к exponential integrators

$$\int_{\lambda_s}^{\lambda_t} e^{-\lambda} \hat{\epsilon}_{\theta}(\hat{x}_{\lambda}, \lambda) d\lambda$$

# Что дает Тейлор

$$\mathbf{x}_{t_{i-1} \rightarrow t_i} = \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \alpha_{t_i} \int_{\lambda_{t_{i-1}}}^{\lambda_{t_i}} e^{-\lambda} \hat{\mathbf{e}}_{\theta}(\hat{\mathbf{x}}_{\lambda}, \lambda) d\lambda$$

$$\hat{\mathbf{e}}_{\theta}(\hat{\mathbf{x}}_{\lambda}, \lambda) = \sum_{n=0}^{k-1} \frac{(\lambda - \lambda_{t_{i-1}})^n}{n!} \hat{\mathbf{e}}_{\theta}^{(n)}(\hat{\mathbf{x}}_{\lambda_{t_{i-1}}}, \lambda_{t_{i-1}}) + \mathcal{O}((\lambda - \lambda_{t_{i-1}})^k)$$

$$\mathbf{x}_{t_{i-1} \rightarrow t_i} = \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \alpha_{t_i} \sum_{n=0}^{k-1} \hat{\mathbf{e}}_{\theta}^{(n)}(\hat{\mathbf{x}}_{\lambda_{t_{i-1}}}, \lambda_{t_{i-1}}) \int_{\lambda_{t_{i-1}}}^{\lambda_{t_i}} e^{-\lambda} \frac{(\lambda - \lambda_{t_{i-1}})^n}{n!} d\lambda + \mathcal{O}(h_i^{k+1})$$

# Производные...

Пока забудем о производных и возьмем  $k = 1$  тогда, получим:

$$\tilde{\mathbf{x}}_{t_i} = \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \sigma_{t_i} (e^{h_i} - 1) \epsilon_{\theta}(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1}), \quad \text{where } h_i = \lambda_{t_i} - \lambda_{t_{i-1}}$$

Назовем метод DPM-Solver-1, порядок сходимости будет  $O(\max_i h_i)$ , то есть чем меньше  $h_i$ , тем лучше сходимость.

Идея: тогда разумнее взять  $h_i = \text{const}$ , получим расписание  $\log\text{SNR}$

# Где-то это уже было... DDIM и DPM-Solver-1

Denoising Diffusion Implicit Models (DDIM) [19] design a deterministic method for fast sampling from DPMs. For two adjacent time steps  $t_{i-1}$  and  $t_i$ , assume that we have a solution  $\tilde{\mathbf{x}}_{t_{i-1}}$  at time  $t_{i-1}$ , then a single step of DDIM from time  $t_{i-1}$  to time  $t_i$  is

$$\tilde{\mathbf{x}}_{t_i} = \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \alpha_{t_i} \left( \frac{\sigma_{t_{i-1}}}{\alpha_{t_{i-1}}} - \frac{\sigma_{t_i}}{\alpha_{t_i}} \right) \epsilon_{\theta}(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1}). \quad (4.1)$$

Although motivated by entirely different perspectives, we show that the updates of DPM-Solver-1 and Denoising Diffusion Implicit Models (DDIM) [19] are identical. By the definition of  $\lambda$ , we have  $\frac{\sigma_{t_{i-1}}}{\alpha_{t_{i-1}}} = e^{-\lambda_{t_{i-1}}}$  and  $\frac{\sigma_{t_i}}{\alpha_{t_i}} = e^{-\lambda_{t_i}}$ . Plugging these and  $h_i = \lambda_{t_i} - \lambda_{t_{i-1}}$  to Eq. (4.1) results in exactly a step of DPM-Solver-1 in Eq. (3.7). However, the semi-linear ODE formulation of DPM-Solver allows for principled generalization to higher-order solvers and convergence order analysis.

# Мы не умеем считать производные...

Когда речь идет о численных методах и возникают производные  
высокого порядка, то значит пора переходить к аппроксимациям)))

Причем чем больше производных мы приблизим, тем выше порядок  
метода

# Зоопарк

DPM-Solver-1

DPM-Solver-2

DPM-Solver-2 (general version)

DPM-Solver-3

DPM-Solver-12

DPM-Solver-23

DPM-Solver-fast

## Algorithm 1 DPM-Solver-2.

**Require:** initial value  $\mathbf{x}_T$ , time steps  $\{t_i\}_{i=0}^M$ , model  $\epsilon_\theta$

- 1:  $\tilde{\mathbf{x}}_{t_0} \leftarrow \mathbf{x}_T$
- 2: **for**  $i \leftarrow 1$  to  $M$  **do**
- 3:  $s_i \leftarrow t_\lambda \left( \frac{\lambda_{t_{i-1}} + \lambda_{t_i}}{2} \right)$
- 4:  $\mathbf{u}_i \leftarrow \frac{\alpha_{s_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \sigma_{s_i} \left( e^{\frac{h_i}{2}} - 1 \right) \epsilon_\theta(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1})$
- 5:  $\tilde{\mathbf{x}}_{t_i} \leftarrow \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \sigma_{t_i} (e^{h_i} - 1) \epsilon_\theta(\mathbf{u}_i, s_i)$
- 6: **end for**
- 7: **return**  $\tilde{\mathbf{x}}_{t_M}$

## Algorithm 2 DPM-Solver-3.

**Require:** initial value  $\mathbf{x}_T$ , time steps  $\{t_i\}_{i=0}^M$ , model  $\epsilon_\theta$

- 1:  $\tilde{\mathbf{x}}_{t_0} \leftarrow \mathbf{x}_T, r_1 \leftarrow \frac{1}{3}, r_2 \leftarrow \frac{2}{3}$
- 2: **for**  $i \leftarrow 1$  to  $M$  **do**
- 3:  $s_{2i-1} \leftarrow t_\lambda (\lambda_{t_{i-1}} + r_1 h_i), \quad s_{2i} \leftarrow t_\lambda (\lambda_{t_{i-1}} + r_2 h_i)$
- 4:  $\mathbf{u}_{2i-1} \leftarrow \frac{\alpha_{s_{2i-1}}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \sigma_{s_{2i-1}} (e^{r_1 h_i} - 1) \epsilon_\theta(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1})$
- 5:  $D_{2i-1} \leftarrow \epsilon_\theta(\mathbf{u}_{2i-1}, s_{2i-1}) - \epsilon_\theta(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1})$
- 6:  $\mathbf{u}_{2i} \leftarrow \frac{\alpha_{s_{2i}}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \sigma_{s_{2i}} (e^{r_2 h_i} - 1) \epsilon_\theta(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1}) - \frac{\sigma_{s_{2i}} r_2}{r_1} \left( \frac{e^{r_2 h_i} - 1}{r_2 h_i} - 1 \right) D_{2i-1}$
- 7:  $D_{2i} \leftarrow \epsilon_\theta(\mathbf{u}_{2i}, s_{2i}) - \epsilon_\theta(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1})$
- 8:  $\tilde{\mathbf{x}}_{t_i} \leftarrow \frac{\alpha_{t_i}}{\alpha_{t_{i-1}}} \tilde{\mathbf{x}}_{t_{i-1}} - \sigma_{t_i} (e^{h_i} - 1) \epsilon_\theta(\tilde{\mathbf{x}}_{t_{i-1}}, t_{i-1}) - \frac{\sigma_{t_i}}{r_2} \left( \frac{e^{h_i} - 1}{h_i} - 1 \right) D_{2i}$
- 9: **end for**
- 10: **return**  $\tilde{\mathbf{x}}_{t_M}$



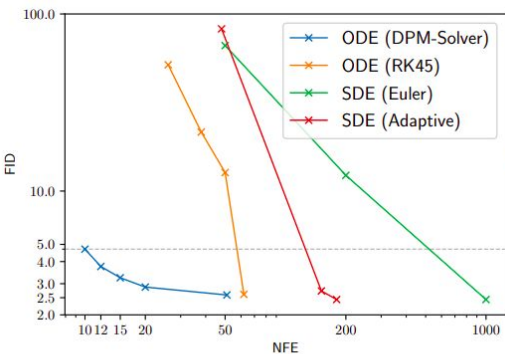
# Explicit Exponential Runge Kutta

$$\frac{d\mathbf{x}_t}{dt} = \alpha \mathbf{x}_t + \mathbf{N}(\mathbf{x}_t, t)$$

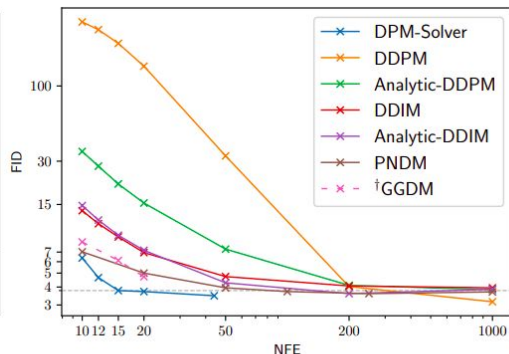
$$\mathbf{x}_{t+h} = e^{\alpha h} \mathbf{x}_t + e^{\alpha h} \int_0^h e^{-\alpha \tau} \mathbf{N}(\mathbf{x}_{t+\tau}, t + \tau) d\tau$$

for approximating the same integral with  $\alpha = 1$  and  $\mathbf{N} = \tilde{\epsilon}_\theta$ . However, DPM-Solver is different from the expRK methods, because their linear term  $e^{\alpha h} \mathbf{x}_t$  is different from our linear term  $\frac{\alpha_{t+h}}{\alpha_t} \mathbf{x}_t$ . In summary, DPM-Solver is inspired by the same technique of expRK for deriving high-order approximations of the exponentially weighted integral, but the formulation of DPM-Solver is different from expRK, and DPM-Solver is customized for the specific formulation of diffusion ODEs.

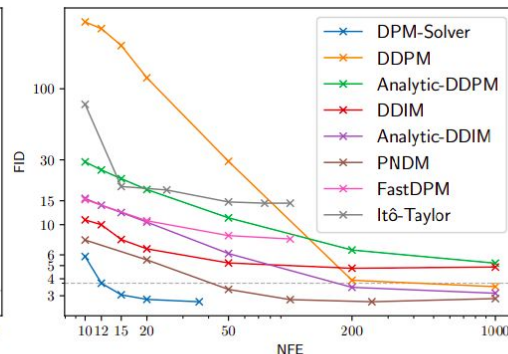
# Результаты



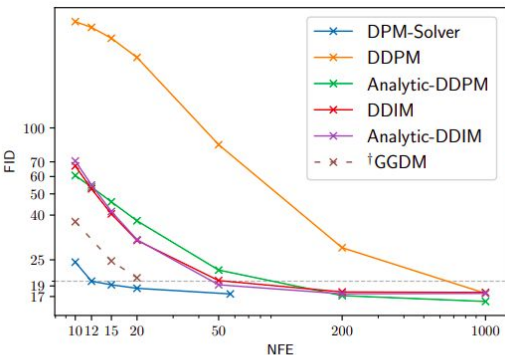
(a) CIFAR-10 (continuous)



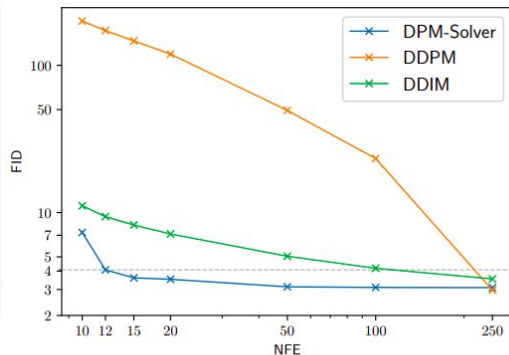
(b) CIFAR-10 (discrete)



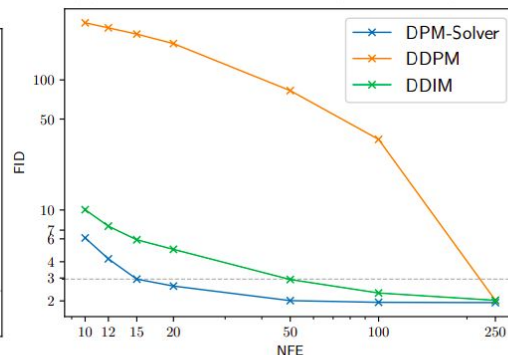
(c) CelebA 64x64 (discrete)



(d) ImageNet 64x64 (discrete)



(e) ImageNet 128x128 (discrete)



(f) LSUN bedroom 256x256 (discrete)

# Итоги

Мы умеем:

- Эффективно учить и параметризовать диффузионные модели
- Использовать более умные солверы, которые учитывают условия задачи
- Получили качественную генерацию в 10–20 шагов
- Умеем использовать одну диффузионную модель для разных параметризаций

Мы не умеем:

...

# Заключения теории

На самом деле, мы многое еще не умеем и не знаем.

Есть подходы основанные на оптимальном транспорте и flow matching лоссами.

Есть подходы со спрямлением траектории, но это отдельные.

Много разной математики и эвристик, много с генерации дискретных объектов.

Есть еще много чего! Но нам не хватит и 10 лекций, поэтому пора чуть остановиться.

Все желающие могут пройти отдельный курс про диффузионные модели "[Генеративные модели на основе ODE и SDE](#)" от Ракитина Дениса Романовича