

Порождающие модели белковых структур

Роман Сергеевич Клыпа

Научный руководитель: к.ф.-м.н. С. В. Грудинин

Кафедра интеллектуальных систем ФПМИ МФТИ

Специализация: Интеллектуальный анализ данных

Направление: 03.04.01 Прикладные математика и физика

2024

Постановка задачи

Задача генерации трехмерных объектов

Необходимо построить модель G , порождающую объекты согласно $p(\mathcal{M})$, где $\mathcal{M} - \mathbb{R}^{3M}$.

Уменьшение размерности

Наличие связей в \mathbb{R}^{3M} позволяет уменьшить пространство до $SE(3)^N$ ($N < M$), которое можно отождествить с $SO(3)^N \times \mathbb{R}^{3N}$.

Score Matching

Вместо $p(\mathbf{x})$ моделируется $\nabla_{\mathbf{x}} \log p(\mathbf{x})$:

$$\mathcal{L}(\theta) = \mathbb{E}_{p(\mathbf{x})} \left[\|s_{\theta}(\mathbf{x}) - \nabla_{\mathbf{x}} \log p(\mathbf{x})\|^2 \right]. \quad (1)$$

Генерация возможна с помощью алгоритма Ланжевена:

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \epsilon \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \sqrt{2\epsilon} \mathbf{z}_i, \quad i = 0, 1, \dots, K, \quad \mathbf{z}_i \sim \mathcal{N}(0, I). \quad (2)$$

Цель исследования

Поставленные цели

- ▶ Предложить прямой и обратный диффузионные процессы для $SO(3)$.
- ▶ Предложить способ обеспечить эквивариантность процессов ($g_1 f(g_2) = f(g_1 g_2)$).
- ▶ Показать преимущества предложенных методов на реальной задаче.

Существующие подходы и их недостатки

- ▶ Предложенные ранее методы не гарантируют сходимость прямого процесса к шуму (Yim et al. 2023).
- ▶ Предложенные ранее методы являются эмпирическими.

Предлагаемый подход: прямой процесс

Здесь и далее $\mathbf{R} \in \text{SO}(3)$, $\log \mathbf{R} = \mathbf{r} \in \mathfrak{so}(3)$, $*$ - композиция элементов $\mathfrak{so}(3)$.

Предложение 1 (Клыпа, 2024)

Прямой процесс на $\text{SO}(3)$:

$$\mathbf{R}_t = \exp[d(t)\mathbf{r}_0 * \tilde{\mathbf{r}}(t)], \quad \tilde{\mathbf{r}}(t) \sim \text{IGSO}_3(\mathbf{Id}, \sigma^2(t)) \quad (3)$$

где $d(t)$ – коэффициент дрифта, $d(0) = 1$, $d(1) = 0$.

Теорема 1 (Клыпа, 2024)

Для прямого процесса 3:

$$p_t(\mathbf{R}_t | \mathbf{R}_0) = \text{IGSO}_3(\mathbf{R}_t | \exp[d(t)\mathbf{r}_0], \sigma^2(t)), \quad (4)$$

в частности $p_1(\mathbf{R}_1 | \mathbf{R}_0) = \text{IGSO}_3(\mathbf{R}_1 | \mathbf{Id}, \sigma^2(1))$.

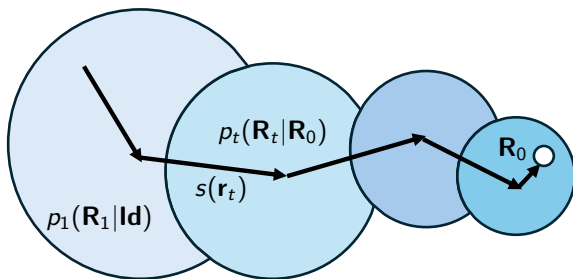
Предлагаемый подход: обратный процесс

Предложение 2 (Клыпа, 2024)

Обратный процесс на $SO(3)$, соответствующий 3:

$$\mathbf{R}_{t-dt} = \exp[\mathbf{r}_t * g^2(t)s(\mathbf{r}_t)dt * g(t)\sqrt{dt}\tilde{\mathbf{r}}(t)], \quad (5)$$

где $s(\mathbf{r}_t) = \nabla_{\mathbf{r}} \log p_t(\mathbf{r}_t | \mathbf{r}_0)$, $g(t) = \sqrt{\frac{d}{dt}\sigma^2(t)}$.



Динамика Ланжевена с отжигом.

Эквивариантность и сходимость процессов

Лемма 1 (Клыпа, 2024)

При $p(\mathbf{r}|\mathbf{r}_0) = \text{IGSO}_3(\mathbf{r}|\mathbf{r}_0, \sigma^2)$, процесс

$$\mathbf{r}_{i+1} = \mathbf{r}_i * \epsilon \nabla_{\mathbf{r}} \log p(\mathbf{r}_i|\mathbf{r}_0) * \sqrt{2\epsilon} \mathbf{z}_i, \quad \mathbf{z}_i \sim \text{IGSO}_3(\mathbf{Id}, \sigma^2) \quad (6)$$

сходится к $p(\mathbf{r}|\mathbf{r}_0)$ при $i \rightarrow \infty$.

Теорема 2 (Клыпа, 2024)

Прямой процесс 3 эквивариантен относительно $\text{SO}(3)$:

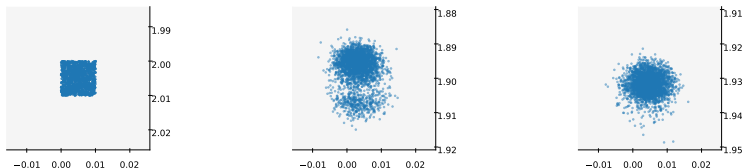
$$\mathbf{R}\mathbf{R}_t(\mathbf{R}_0) = \mathbf{R}_t(\mathbf{R}\mathbf{R}_0) \quad \forall \mathbf{R} \in \text{SO}(3). \quad (7)$$

Если $s_\theta(\mathbf{r}_t)$ инвариантен относительно $\text{SO}(3)$, то обратный процесс 5 также эквивариантен:

$$\mathbf{R}\mathbf{R}_t(\mathbf{R}_1) = \mathbf{R}_t(\mathbf{R}\mathbf{R}_1) \quad \forall \mathbf{R} \in \text{SO}(3). \quad (8)$$

Вычислительный эксперимент

Элементы группы $SO(3)$ изображены в трехмерном пространстве при использовании представления векторов Эйлера. Для наглядности полученные распределения были спроецированы на плоскость.

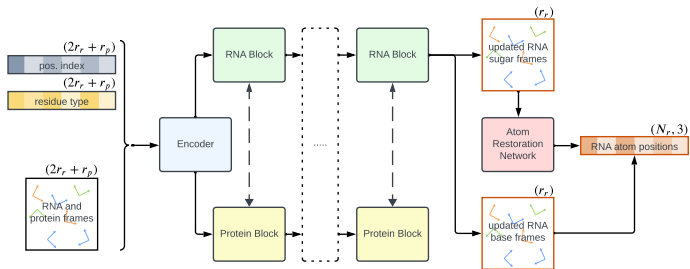


Слева изображено целевое распределение, по центру - результат генерации метода Yim et al. 2023., справа - результат предлагаемого нами метода.

Генерация трехмерных структур РНК

Постановка задачи

Генерация трехмерных структур молекул РНК, при взаимодействии с протеином: $P(\text{geom}_{RNA} | \text{geom}_{pr}, \text{seq}_{pr}, \text{seq}_{RNA})$. Предполагается, что $\text{geom}_{RNA}, \text{geom}_{pr} \in \text{SE}(3)^N$.



Архитектура модели MolBindDif. Здесь r_r - количество остатков РНК, r_p - количество остатков белка, а N_r - общее количество атомов в РНК.

Результаты MolBindDif

Теорема 3 (Клыпа, 2024)

Архитектура модели MolBindDif эквивариантна относительно $SO(3)$.

Результаты эксперимента

Процесс	$\downarrow \text{rRMSD}_{rr}, \text{\AA}$	$\downarrow \text{rRMSD}_{rp}, \text{\AA}$	$\uparrow \text{IDDT}_{rr}$	$\uparrow \text{IDDT}_{rp}$
Yim et al. 2023	10.7 ± 5.5	14.1 ± 7.8	0.17 ± 0.06	0.10 ± 0.05
Klypa, 2024	11.1 ± 5.1	13.4 ± 8.2	0.15 ± 0.06	0.11 ± 0.05
RoseTTaFoldNA	13.3 ± 6.3	17.6 ± 7.8	0.19 ± 0.10	0.08 ± 0.04

Результаты генерации трехмерных структур РНК. RoseTTaFoldNA не является генеративным процессом.

Выносятся на защиту

1. Предложен новый порождающий процесс на $SO(3)$.
2. Доказана эквивариантность процесса.
3. Продемонстрированы его преимущества относительно прошлых работ на синтетических данных.
4. Продемонстрирована его сходимость на реальных данных.

Публикации

1. Roman Klypa, Kliment Olechnovič, Ben Shor, Dina Schneidman-Duhovny, Sergei Grudinin. *MolBindDif: Protein-conditioned RNA structure diffusion*
OpenReview preprint, ICML 2024 submission