

СОДЕРЖАНИЕ

1	Практическая часть	3
1.1	Цель	3
1.2	Используемые данные	3
1.2.1	Informer	3
1.2.1.1	ProbSparse Self-Attention	4
1.2.1.2	Self-Attention Distilling & Кодировщик	5
1.2.1.3	Генеративный декодер	5
1.2.2	Performer	7
1.2.3	Autoformer	7
1.3	Методология	7
1.3.1	Повышение эффективности извлечения локальных паттернов	7
1.3.2	Заимствование механизма внимания из Performer	7
1.3.3	Внедрение модуля декомпозиции ряда из Autoformer	7
1.4	Эксперимент	7
1.4.1	Датасет	7
1.4.2	Детали	7
1.4.2.1	Training	7
1.4.2.2	Baselines	7
1.4.2.3	Hyper-parameter tuning	7
1.4.2.4	Setup	7
1.4.2.5	Metrics	7
1.4.2.6	Platform	7
1.5	Результаты	7
1.6	Заключение	7

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	8
---	----------

1 Практическая часть

Abstract

1.1 Цель

С момента своей публикации, модель Трансформера [4] завоевала огромное признание. Однако у нее есть несколько серьезных проблем, которые усложняют работу с длинными временными последовательностями (LSTF). Последующие исследования предложили различные методы решения данных и связующих проблем (см. Informer [6], Performer [1], Autoformer [5], PatchTST [3], TFT [2] и др.). В данной работе мы сфокусируемся на первых трех.

В данной работе, мы предлагаем заменить слой эмбединга в модели Informer компактным двухслойным сверточным блоком с целью повышения эффективности извлечения локальных паттернов, внедрить модуль декомпозиции ряда из Autoformer для явного разделения трендовых и сезонных компонент и заменить ProbSparse-внимание на линейное FAVOR+ из Performer для учёта глобальных зависимостей при низких вычислительных затратах.

1.2 Используемые данные

Прежде чем перейти к методологии нашей работы, рассмотрим модели и понятия, которыми далее будем пользоваться.

1.2.1 Informer

В 2021 году, Zhou et al., опубликовали свою статью «**Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting**», в которой представили новую, основанную на Трансформере [4] модель под названием **Informer** [6].

Informer был создан для решения задачи **Long-sequence time-series forecasting (LSTF)**. Zhou et al. поставили перед собой следующий вопрос: Можем ли мы построить модель, основанную на трансформере, которая

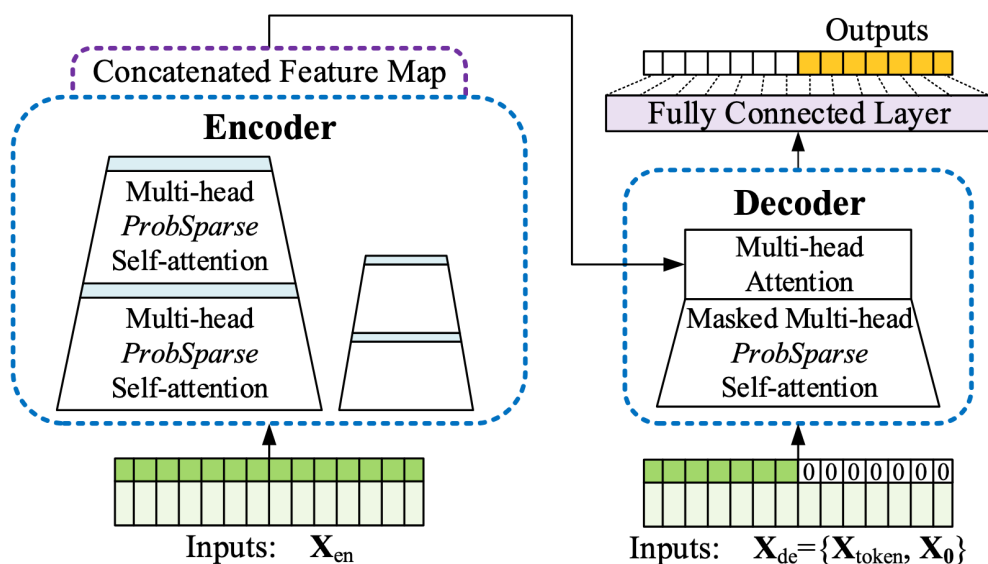


Рис. 1.1: Обзор модели Informer [6].

- (a) захватывает очень длинные зависимости
- (b) эффективно работает для тысячи временных шагов

Ключевые рассматриваемые слабости трансформера:

- Квадратичная стоимость механизма само-внимания (self-attention) для последовательностей длины L .
- Накладывание слоев умножает эту стоимость, достигая ограничений по памяти.
- Пошаговое (динамичное) декодирование работает медленно и накапливает ошибки.

Информер отвечает на каждый из этих вопросов, реконструируя механизм внимания, кодировщик и декодер.

1.2.1.1 ProbSparse Self-Attention

Зачастую, в длинных временных рядах, большинство скалярных произведений запросов с ключами пренебрежимо малы и лишь некоторые из них достаточно больши. Вместо того, чтобы считать все возможные пары запрос-ключ, informer предлагает следующее:

1. Измерить разброс каждого запроса q_i :

$$M(q_i, K) = \max_j \frac{q_i k_j^T}{\sqrt{d}} - \frac{1}{L_K} \sum_{j=1}^{L_K} \frac{q_i k_j^T}{\sqrt{d}}$$

2. Выбрать u лучших запросов, исходя из $M(q_i, K)$, где $u \propto \ln L$.
3. Вычислить полное внимание только для этих u рядов и аппроксимировать остальные через среднее значение.

Что сокращает как время, так и память с $O(L^2)$ до $O(L \log L)$, при этом сохраняя всю важную информацию.

1.2.1.2 Self-Attention Distilling & Кодировщик

Даже после предыдущей операции, каждый слой все равно производит карты признаков длины L , многие из которых повторяют похожие паттерны. Мы можем «дистиллировать» сильнейшие сигналы и сократить последовательность по мере продвижения.

- После каждого блока с механизмом внимания, применяем:
 1. 1-D свертку + ELU активацию,
 2. max-pool со страйдом 2

Что уменьшает размерность в два раза на каждом слое, результируя в пирамиде стэков, чьи выходы в конечном итоге конкатенируют.

Self-attention distilling фокусируется на доминирующих паттернах, при этом сокращая память до $O((2 - \epsilon)L \log L)$.

1.2.1.3 Генеративный декодер

Вместо того, чтобы генерировать токены по очереди один за другим, заимствуем трюк с начальными токеном из NLP:

- Взять срез известной истории (например 5 дней перед целевыми 7ю днями) в качестве начального токена

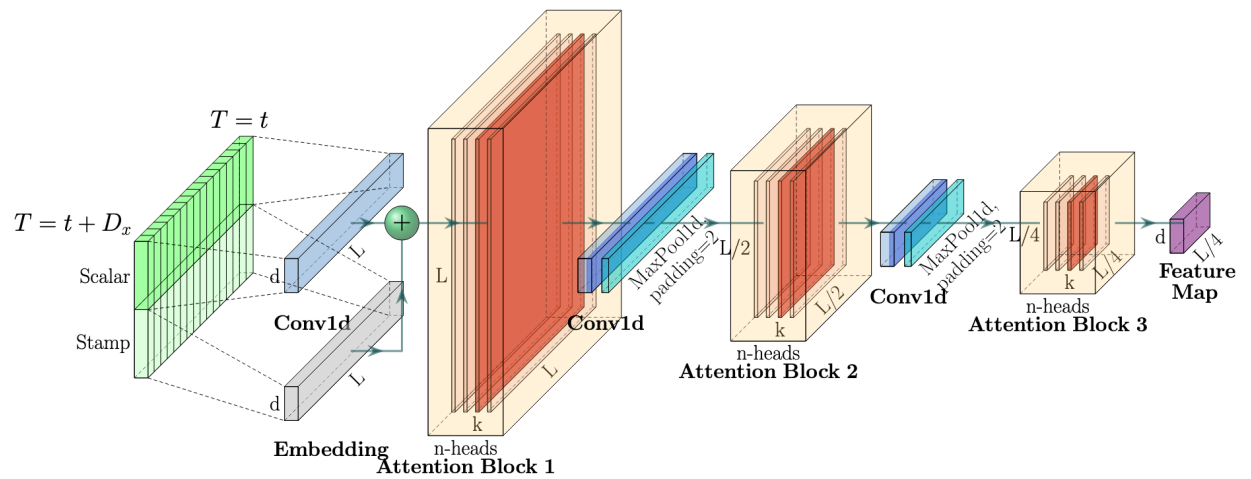


Рис. 1.2: Обзор одного стэка в кодировщике Informer-a [6].

- Разместить плейсхолдеры для всех L_y будущих значение (используя только их временные метки для позиционного контекста)
- За один прямой проход одновременно заполнить все L_y выходов через маскированное ProbSparse внимание.

Что избегает накапливания ошибки и работает гораздо быстрее.

1.2.2 Performer

1.2.3 Autoformer

1.3 Методология

1.3.1 Повышение эффективности извлечения локальных паттернов

1.3.2 Заимствование механизма внимания из Performer

1.3.3 Внедрение модуля декомпозиции ряда из Autoformer

1.4 Эксперимент

1.4.1 Датасет

1.4.2 Детали

1.4.2.1 Training

1.4.2.2 Baselines

1.4.2.3 Hyper-parameter tuning

1.4.2.4 Setup

1.4.2.5 Metrics

1.4.2.6 Platform

1.5 Результаты

1.6 Заключение

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- [1] Krzysztof Choromanski, Valentin Likhoshesterov, David Dohan, Xingyou Song, Alex Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Łukasz Kaiser, David Belanger, Lucy Colwell, and Adrian Weller. Rethinking attention with performers. In Proceedings of the International Conference on Learning Representations (ICLR), 2021. arXiv:2009.14794.
- [2] Bryan Lim, Sercan Ö. Arik, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. In Proceedings of the International Conference on Learning Representations (ICLR), 2021. arXiv:1912.09363.
- [3] Wenjie Nie, Di He, Tao Qin, Ming Zhou, and Tie-Yan Liu. A time series is worth 64 words: Long-term forecasting with transformers. In Proceedings of the International Conference on Learning Representations (ICLR), 2023. arXiv:2211.14730.
- [4] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems (NeurIPS), volume 30, 2017.
- [5] Haixu Wu, Yao Xu, Jindong Wang, Guodong Long, Xingquan Jiang, Chengqi Zhang, and Lina Yao. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. In Advances in Neural Information Processing Systems (NeurIPS), volume 34, pages 22419–22430, 2021.
- [6] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 11106–11115, 2021.