

# Эффективное агрегирование по меткам для задачи последовательностей событий

Галина Леонидовна Боева  
Научный руководитель: к.ф.-м.н. А. А. Зайцев

Кафедра интеллектуальных систем ФПМИ МФТИ  
Специализация: Интеллектуальный анализ данных  
Направление: 09.04.01 Информатика и вычислительная техника

2025

# Агрегирование по меткам для задачи последовательностей событий

## Проблема

Современные подходы фокусируются на архитектуре преобразования последовательных данных, агрегируя данные по временным меткам, но теряя информацию о взаимозависимостях меток.

## Цель работы

Создание подхода, основанного на механизме собственного внимания над метками, предшествующими прогнозируемому шагу.

## Задачи

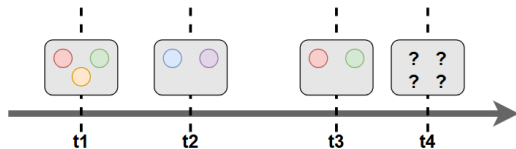
1. разработка метода на основе внимания для предсказания множества меток
2. формирование теоретического обоснования эффективности применения агрегирования по меткам
3. обоснование причинно-следственных связей с помощью построения графа на основе внимания

# Постановка задачи предсказания временных наборов меток

1.  $U = \{u_1, u_2, \dots, u_N\}$  — множество из  $N$  элементов
2.  $Y = \{y_1, y_2, \dots, y_M\}$  — словарь меток размера  $M$
3.  $s \subset Y$  — произвольное подмножество меток
4.  $S_i = \{s_i^1, s_i^2, \dots, s_i^T\}$  — временная последовательность из  $T$  множеств для элемента  $u_i$

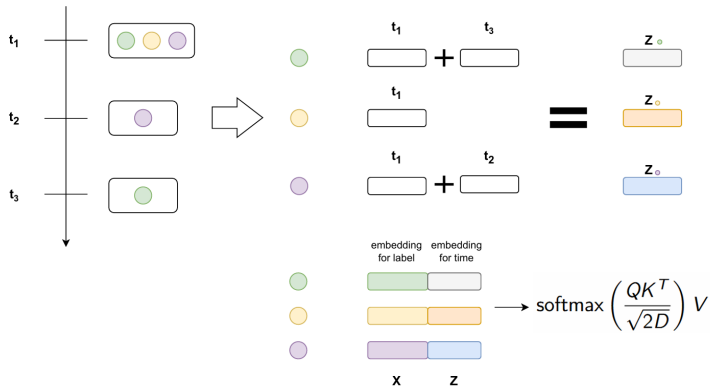
## Задача:

Для заданной последовательности исторических множеств  $S_i = \{s_i^1, s_i^2, \dots, s_i^T\}$  элемента  $u_i \in U$  предсказать следующее множество меток  $\hat{s}_i^{T+1}$ .



# Предложенный метод на основе внимания на метках

1.  $X \in \mathbb{R}^{M \times D}$  - матрица эмбеддингов меток
2.  $t_j \in \mathbb{R}^D$  - временные эмбеддинги (позиционное кодирование)
3.  $Z \in \mathbb{R}^{M \times D}$  - временные признаки:  $Z^{(m,:)} = \sum_{j|y_m \in s_i^j} t_j$
4.  $G = X \oplus Z \in \mathbb{R}^{M \times 2D}$  - конкатенированные признаки



## Постановка задачи агрегирования по меткам

Агрегация по меткам:  $N_m = \sum_{t=1}^T \mathbf{1}(y_m \in s_t)$ ,  $\mathbf{Z}_{\text{label}}(m, :) = N_m \cdot \mathbf{x}_m$

Агрегация по времени:  $\mathbf{Z}_{\text{time}}(m, :) = \sum_{t: y_m \in s_t} \mathbf{t}_t$

Теорема 1 (Боева, 2025, об эффективности агрегирования по меткам)

Если  $M \ll T$ , то:

$$T_{\text{label}} = O(M^2 D) \ll T_{\text{time}} = O(T^2 D)$$

## Эффективное агрегирование по меткам

**Лемма:** Пусть  $G \in \mathbb{R}^{L \times D}$  — матрица входных представлений, где  $L$  — количество объектов (меток или временных событий), а  $D$  — размерность эмбединга. Тогда время выполнения одного трансформерного слоя внимания:

$$T_{\text{attn}} = O(L^2 D).$$

- для агрегирования по меткам:  $T_{\text{label}} = O(M^2 D)$

- для агрегирования по времени:  $T_{\text{time}} = O(T^2 D)$

**Сравнение вычислительной сложности:**

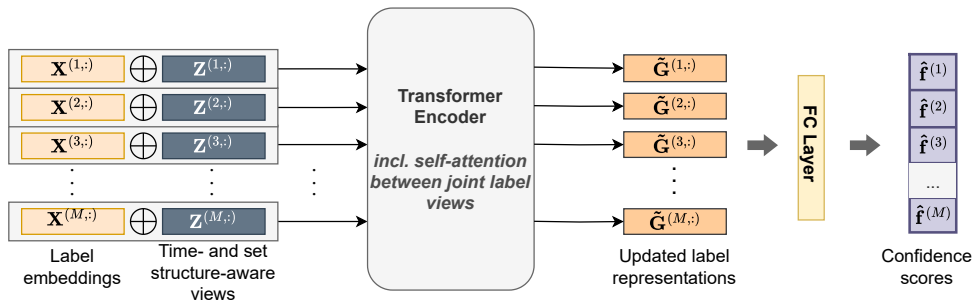
Рассмотрим отношение времени работы моделей:

$$\frac{T_{\text{label}}}{T_{\text{time}}} = \frac{O(M^2 D)}{O(T^2 D)} = O\left(\frac{M^2}{T^2}\right)$$

Если выполняется условие  $M \ll T$ , то есть  $\lim_{T \rightarrow \infty} \frac{M}{T} = 0$ , то:

$$\frac{T_{\text{label}}}{T_{\text{time}}} \rightarrow 0 \quad \Rightarrow \quad T_{\text{label}} \ll T_{\text{time}}$$

## Предложенный метод на основе внимания на метках



Общий пайплайн получения глобальных представлений

## Вычислительный эксперимент: Данные

Статистика наборов данных для прогнозирования временных наборов.

Dataset	#Sets	MdnSS	MaxSS	Vocab	MnLen	#Seqs
Mimic III	17 849	5	23	169	2.7	6636
Instacart	115 604	6	43	134	16.5	7000

- ▶ **Mimic III** — датасет, состоящий из медицинских карт пациентов из отделения интенсивной терапии. Событие, связанное с пациентом, включает в себя время поступления в больницу и набор классификационных кодов заболеваний.
- ▶ **Instacart** — набор данных содержит записи о заказах товаров пользователями. Товары из маркетплейсов и магазинов.

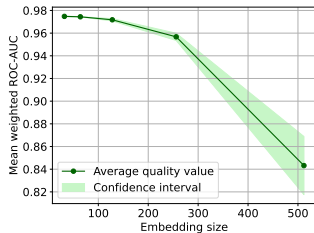


## Вычислительный эксперимент: Основные результаты

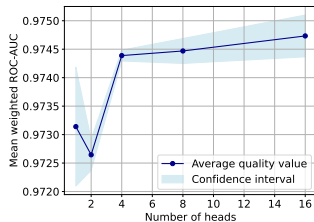
Сравнение подхода our LANET с существующими моделями для прогнозирования временных наборов на основе четырех наборов данных. Выделены наилучшие значения, а вторые по значению подчеркнуты.

Data	Model	Weighted F1 $\uparrow$	Weighted ROC-AUC $\uparrow$	Hamming Loss $\downarrow$
Mim	SFCNTSP	$0.3791 \pm 0.0081$	$0.7034 \pm 0.0024$	$0.0377 \pm 0.0004$
	DNNTSP	$0.3928 \pm 0.0030$	$0.6926 \pm 0.0003$	$0.0365 \pm 0.0003$
	GPTopFreq	$0.4291 \pm 0.0073$	$0.6912 \pm 0.0028$	$0.0398 \pm 0.0005$
	TCMBN	<u><math>0.4979 \pm 0.0180</math></u>	<u><math>0.8670 \pm 0.0095</math></u>	<u><math>0.0305 \pm 0.0008</math></u>
	LANET(ours)	<b><math>0.8214 \pm 0.0224</math></b>	<b><math>0.9852 \pm 0.0023</math></b>	<b><math>0.0220 \pm 0.0001</math></b>
Ins	SFCNTSP	$0.1672 \pm 0.0112$	$0.6852 \pm 0.0448$	$0.0581 \pm 0.0004$
	DNNTSP	<u><math>0.4160 \pm 0.0009</math></u>	$0.7913 \pm 0.0004$	$0.0541 \pm 0.0002$
	GPTopFreq	$0.4087 \pm 0.0079$	$0.7736 \pm 0.0039$	<u><math>0.0529 \pm 0.0008</math></u>
	TCMBN	$0.3687 \pm 0.0065$	<u><math>0.8187 \pm 0.0030</math></u>	$0.0530 \pm 0.0005$
	LANET(ours)	<b><math>0.6159 \pm 0.0029</math></b>	<b><math>0.9445 \pm 0.0008</math></b>	<b><math>0.0474 \pm 0.0003</math></b>

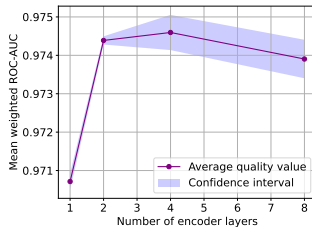
# Вычислительный эксперимент: Дополнительные исследования



Зависимость качества LANET от размера векторных представлений.



Зависимость качества LANET от количества голов во внимании.



Зависимость качества LANET от количества слоев энкодера.

## Графовая интерпретация внимания на метках

Пусть  $G = (V, E)$  — исходный граф, где:  $V$  — множество вершин (меток товаров),  $E$  — множество рёбер (зависимостей между метками).

После удаления вершины  $v_{\max} \in V$ , соответствующей метке с наибольшим суммарным весом внимания, получаем новый граф:

$G' = (V', E')$ , где:  $V' = V \setminus \{v_{\max}\}$ ,  $E' = E \cap (V' \times V')$ .

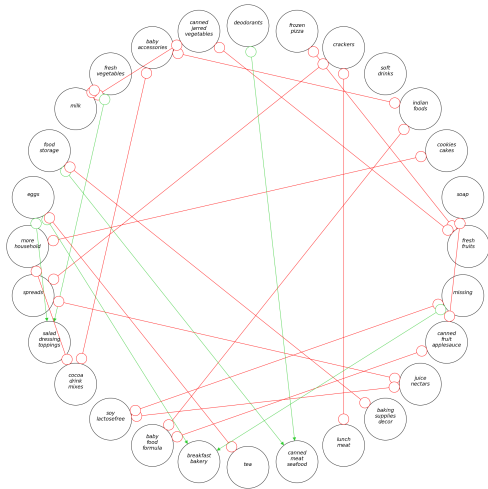
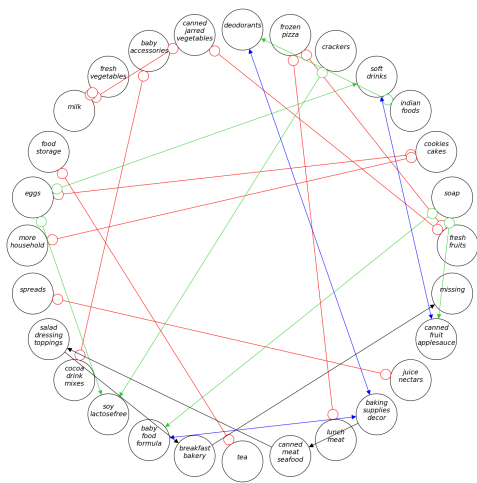
**Красные линии** — обозначают близость или частое совместное появление меток в корзинах покупок. Такие связи могут быть не причинными, а статистическими.

**Синие двунаправленные стрелки** — указывают на сложные взаимодействия между метками, часто через скрытые общие причины (например, сезонность, вкусовые предпочтения).

**Черные односторонние стрелки** — представляют собой направленные причинные связи: если  $A \rightarrow B$ , это означает, что метка  $A$  влияет на выбор  $B$ .

**Зеленые стрелки** — обратные связи, где метка является дочерней по отношению к другой.

# Графовая интерпретация внимания на метках



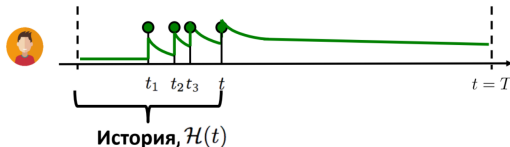
# Временные точечные процессы

Рассмотрим **точечный процесс**  $\{N_t\}_{t \geq 0}$ , где  $N_t$  — количество событий на интервале  $[0, t]$ . Точечный процесс — это стохастический процесс, описывающий случайные моменты времени  $t_1 < t_2 < \dots < t_n$ , в которые происходят события.

$$\lambda_t = \lim_{\Delta t \rightarrow 0^+} \frac{\mathbb{E}[N_{t+\Delta t} - N_t \mid \mathcal{H}_{t-}]}{\Delta t}, \quad (1)$$

где  $\mathcal{H}_{t-} = \sigma(N_s : 0 \leq s < t)$  — предистория до момента  $t$

**Процесс Хокса** — это разновидность точечного процесса, в котором каждое событие увеличивает вероятность возникновения будущих событий:  $\lambda_t \mu + \int_0^t g(t-s) dN_s$ , где  $\mu > 0$  — базовая интенсивность,  $g : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  — ядро возбуждения.



## Агрегирование по меткам в процессах Хокса

Рассмотрим многомерный процесс Хокса  $\{N_t^m\}_{t \geq 0}$ , где  $m = 1, \dots, M$  обозначает тип события (метку), и интенсивность определяется как:

$$\lambda_t^m = \mu_m + \sum_{n=1}^M \int_0^t g_{mn}(t-s) dN_s^n.$$

Рассмотрим агрегированный сигнал:  $S_t^m = \frac{1}{t} \int_0^t \phi(m, N_s) ds$ .

### Теорема 2 (Боева, 2025, применение к процессам Хокса)

Если  $\phi(m, N_s)$  является достаточной статистикой для  $\lambda_t^m$ , то существует последовательность оценок  $\{\hat{\lambda}_t^m\}_{t \geq 0}$ , основанных на  $\{S_t^m\}_{t \geq 0}$ , такая что:

$$\lim_{t \rightarrow \infty} \mathbb{E} \left[ \ell(\hat{\lambda}_t^m, \lambda_t^m) \right] = 0,$$

и эта сходимость равномерна по  $m = 1, \dots, M$ .

## Выносятся на защиту

1. Доказана теорема об эффективности агрегирования по меткам над агрегированием по времени.
2. Разработана модель LANET, которая основана на агрегировании данных по меткам.
3. Выполнены ряд экспериментов для задачи классификации с несколькими метками на двух различных выборках и сравнение с базовыми подходами в данной области.
4. Проанализированы причинно-следственные связи в self-attention, где используется графовый подход на основе построения PAG для взаимосвязи меток.
5. Сформированы выводы о зависимости метрик от гиперпараметра, отвечающего за размер входных представлений, количество голов во внимании и также количества слоев энкодера.

## Список работ автора по теме диплома

Статья опубликована в октябре 2024 года на конференцию ранга A ECAI.

1. Elizaveta Kovtun, Galina Boeva, Andrey Shulga, and Alexey Zaytsev. Label Attention Network for Temporal Sets Prediction: You Were Looking at a Wrong Self-Attention, IOS Press, October 2024.
2. Vladislav Zhuzhel, Galina Boeva, Vsevolod Grabar, Artem Zabolotnyi, Alexander Stepikin, Vladimir Zholobov, Maria Ivanova, Mikhail Orlov, Ivan Kireev, Evgeny Burnaev, Rodrigo Rivera-Castro, Alexey Zaytsev. Continuous-time convolutions model of event sequences (2023). Статья подана в журнал Experts Systems With Applications.
3. Ilya Kuleshov, Galina Boeva, Vladislav Zhuzhel, Evgeni Vorsin, Evgenia Romanenkova, Alexey Zaytsev. DeNOTS: Stable Deep Neural ODEs for Time Series (2024). Статья подана на NeurIPS 2025.

**Вклад:** создание теоретического обоснования подхода в виде двух теорем, разработка идеи статьи, базовые подходы, исследование устойчивости модели и графовая интерпретация внимания.