

Addition of external information for enhancement of local embeddings for event sequences data models

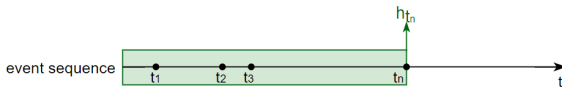
Maria Kovaleva
Research advisor: Ph.D. Alexey Zaytsev

Department of Intelligent Systems
Moscow Institute of Physics and Technology
&
Data Science
Skolkovo Institute of Science and Technology

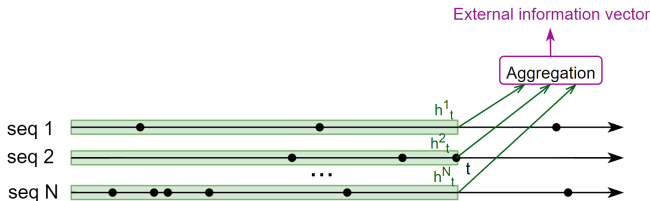
2024

Introduction. External information

General problem: building of embeddings \mathbf{h}_t for event sequences



Gaps: self-supervised models for representation learning ignore external information



Idea: the external information is contained in sequences themselves and can be represented as their aggregation

Aim and Objectives

The aim of the work:

To enhance the embeddings for event sequences data models using aggregation of external information.

Objectives:

1. Development of aggregation methods for external aggregation accounting
2. Validation of developed methods on bank transactions data

Problem statement and baseline approaches

$D = \{S^i\}_{i=1}^n$ — set of n sequences

$S^i = \left\{ \left(t_j^i, \mathbf{z}_j^i \right) \right\}_{j=0}^{T^i}$ — event sequence

$t_j^i \in [0, T^i]$ — time of the event; $\mathbf{z}_j^i \in \mathbb{R}^d$ — description of the event

Embeddings construction: $e(S^i) = H^i$

- ▶ e is encoder: usually dense NN for event description encoding + recurrent NN
- ▶ $H^i = \left\{ \left(t_j^i, \mathbf{h}_j^i \right) \right\}$ — embeddings

Contrastive learning for encoder:

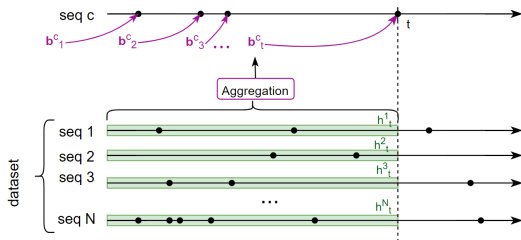
$L_{km} = I_{k=m} d(\mathbf{h}^k, \mathbf{h}^l)^2 + \frac{1}{2} (1 - I_{k=m}) \max \{0, \rho - d(\mathbf{h}^k, \mathbf{h}^l)\}^2$,
where d is a distance between embeddings and ρ is a hyperparameter

Autoregressive learning for encoder:

A loss function consists of the cross-entropy for the categorical features and MSE for the continuous features

Receiving representations of external information and taxonomy of aggregation methods

Construction of a vector of external information



Taxonomy of aggregation methods

When constructing aggregation methods, one can use the similarity of the current sequence and sequences from the training set

Similarity of sequences:

1. by embeddings
2. by time

General formula for aggregation: $\mathbf{b}_\tau^c = f(H, f_e(H, \mathbf{h}_{t<\tau}^c), f_t(\tau, T))$

Where $H = [\mathbf{h}_{t<\tau}^1, \dots, \mathbf{h}_{t<\tau}^n]$ is a matrix with embeddings of all sequences at current time,

$T = [t_{t<\tau}^1, \dots, t_{t<\tau}^n]$ is a vector of the last event times for all sequences;
 f_e and f_t are functions to measure similarity by embeddings and time,
 f is aggregation function, usually weighed sum of vectors from H .

Proposed aggregation methods

Classic:

1 **Mean:** $\mathbf{b}_\tau^c = \frac{1}{n} \sum_{i=1}^n \mathbf{h}_{t < \tau}^i$

2 **Max:** $\mathbf{b}_\tau^c = \max(H)$

Inspired by the Hawkes process*:

3 **Exp Hawkes:**

$$\mathbf{b}_\tau^c = H \exp(-(\tau \mathbf{1} - T))$$

4 **Exp learnable Hawkes:** $\mathbf{b}_\tau^c =$

$$\phi_{NN}(\text{concat}(H, \mathbf{h}_\tau^c)) \exp(-(\tau \mathbf{1} - T))$$

5 **Attention Hawkes:** $\mathbf{b}_\tau^c =$

$$H(\text{softmax}(H^T \mathbf{h}_\tau^c) \odot \exp(-(\tau \mathbf{1} - T)))$$

Attention based:

6 **Attention:**

$$\mathbf{b}_\tau^c = H \text{softmax}(H^T \mathbf{h}_\tau^c)$$

7 **Learnable attention:**

$$\mathbf{b}_\tau^c = H \text{softmax}(H^T A \mathbf{h}_\tau^c)$$

8 **Symmetrical attention:**

$$\mathbf{b}_\tau^c = H \text{softmax}(H^T S^T S \mathbf{h}_\tau^c)$$

9 **Kernel attention:**

$$\mathbf{b}_\tau^c = H \text{softmax}(\phi(H^T) \phi(\mathbf{h}_\tau^c))$$

Where A and S are matrices with learnable parameters,
 ϕ and ϕ_{NN} are learnable transformation (two-layer NN).

* Laub P. J., Taimre T., Pollett P. K. Hawkes processes. arXiv preprint arXiv:1507.02822. – 2015.

Validation of proposed methods

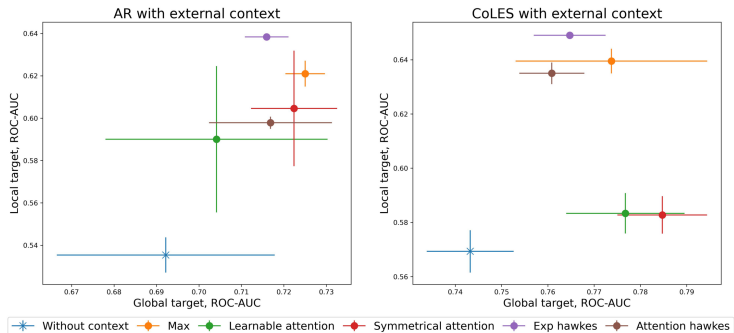
Dataset of the bank transactions:

1. Each sequence S^i : one bank client transactions
2. Each event \mathbf{Z}_j^i : transaction (merchant category code + amount)
3. Target: whether the client left the bank

Validation:

1. **Global** : the model inference on the whole sequence to check its global patterns
 - a binary downstream task
 - b one output vector for a boosting model
2. **Local**: the model inference on sliding windows to check local properties.
 - a next event type prediction
 - b local target prediction
 - c MLP head for prediction

Results: external information addition enhances metrics



The best models are to the right and higher.

1. Accounting of the **external information** improve metrics
2. **Exp Hawkes** method is the best for the local task
3. **Classic** and **attention based** methods are the best for the global task

Defence statements

1. It was suggested to use embeddings aggregations to account for external information in event sequences.
2. A taxonomy of aggregation methods is proposed and specific methods are implemented.
3. It was shown that addition of external information improves the quality of embeddings when used in various applied problems with real data.

Publications

1. Bazarova, A.*, Kovaleva, M.[†]*, Kuleshov, I.*, Romanenkova, E.*, Stepikin, A.*, Yugay, A.*, Mollaev, D., Kireev, I., Savchenko, A., and Zaytsev, A. Universal representations for financial transactional data: embracing local, global, and external contexts. arXiv preprint arXiv:2404.02047 (2024)

*Equal contribution

[†]Contribution: analysis of the external information addition

Additional slides. Results.

	Global target			
	Contrastive learning		Autoregressive learning	
	ROC-AUC	PR-AUC	ROC-AUC	PR-AUC
Without contex	0.743 \pm 0.009	0.792 \pm 0.014	0.692 \pm 0.025	0.734 \pm 0.032
Mean	0.773 \pm 0.004	<u>0.828 \pm 0.003</u>	<u>0.722 \pm 0.007</u>	<u>0.776 \pm 0.005</u>
Max	0.774 \pm 0.021	<u>0.818 \pm 0.032</u>	0.725 \pm 0.005	0.777 \pm 0.002
Attention	0.760 \pm 0.014	0.808 \pm 0.017	0.696 \pm 0.014	0.744 \pm 0.017
Learn. attention	0.777 \pm 0.013	0.830 \pm 0.013	0.704 \pm 0.026	0.751 \pm 0.020
Sym. attention	0.785 \pm 0.010	0.835 \pm 0.005	<u>0.722 \pm 0.010</u>	0.769 \pm 0.004
Kernel attention	<u>0.775 \pm 0.003</u>	0.824 \pm 0.002	<u>0.709 \pm 0.019</u>	<u>0.760 \pm 0.003</u>
Exp Hawkes	0.765 \pm 0.008	0.814 \pm 0.009	0.716 \pm 0.005	<u>0.767 \pm 0.013</u>
Exp learn. Hawkes	0.764 \pm 0.008	0.812 \pm 0.008	0.714 \pm 0.025	0.758 \pm 0.020
Attention Hawkes	0.761 \pm 0.007	0.796 \pm 0.009	<u>0.717 \pm 0.014</u>	0.751 \pm 0.023
	Local target			
Without contex	0.569 \pm 0.008	0.321 \pm 0.003	0.535 \pm 0.008	0.299 \pm 0.011
Mean	0.592 \pm 0.005	0.342 \pm 0.005	0.543 \pm 0.006	0.312 \pm 0.006
Max	<u>0.640 \pm 0.005</u>	0.400 \pm 0.006	<u>0.621 \pm 0.006</u>	0.256 \pm 0.008
Attention	0.600 \pm 0.009	0.348 \pm 0.010	<u>0.534 \pm 0.016</u>	0.301 \pm 0.007
Learn. attention	0.583 \pm 0.007	0.330 \pm 0.008	0.590 \pm 0.035	<u>0.338 \pm 0.025</u>
Sym. attention	0.583 \pm 0.007	0.329 \pm 0.007	<u>0.605 \pm 0.027</u>	<u>0.350 \pm 0.020</u>
Kernel attention	0.582 \pm 0.007	0.329 \pm 0.007	<u>0.572 \pm 0.021</u>	0.330 \pm 0.023
Exp Hawkes	0.649 \pm 0.000	<u>0.366 \pm 0.003</u>	0.638 \pm 0.001	0.351 \pm 0.001
Exp learn. Hawkes	0.581 \pm 0.012	0.322 \pm 0.013	0.539 \pm 0.034	0.293 \pm 0.025
Attention Hawkes	<u>0.635 \pm 0.004</u>	<u>0.359 \pm 0.005</u>	0.598 \pm 0.003	0.331 \pm 0.001

Results of validation for different methods. The best values are **bolded**, the second values are underlined, the third values are double underlined.