# Cohortney: the framework for clustering of event sequences

Taras Khakhulin, Mikhail Pautov

# Problem statement

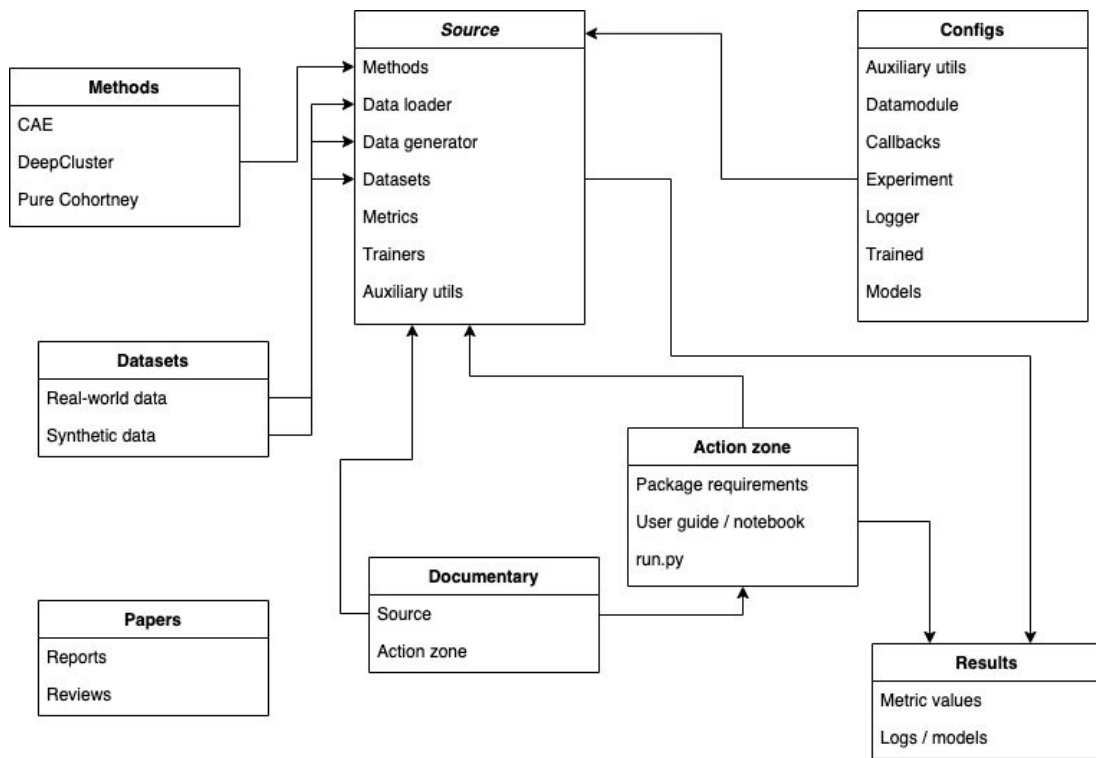The problem of clustering of event sequences is formulated as follows:
- Set of event sequences $S$ contains of objects $s = \{(t, c)\}$ as sequences of occurrences of $c$-type events (out of C possible types, $\{1,2,...,C\}$) at time points $t$.

- The goal is to cluster set $S$ by splitting it onto an optimal (in some sense) number of non-intersecting subsets in which there are similar (in some sense) event sequences.

- It is supposed that the characterization of point processes is done via intensity function – expected grow of $c$-type event at time $t$ given number of its occurrences before $t$.

Our goal is to build a **framework** for methods solving such a problem.

Skoltech

# Potential application / use-cases

- The framework as a tool: testing and evaluation of different methods of event clustering

- The methods' application
  - Banking sphere – bank transactions, user actions monitoring
  - Business processes – prediction / preempt of failure regarding certain actions
  - Healthcare system – monitoring of medical manifestations / features of patients

- Clustering in a nutshell
  - Websites' usage monitoring
  - Clustering visitors of different institutions

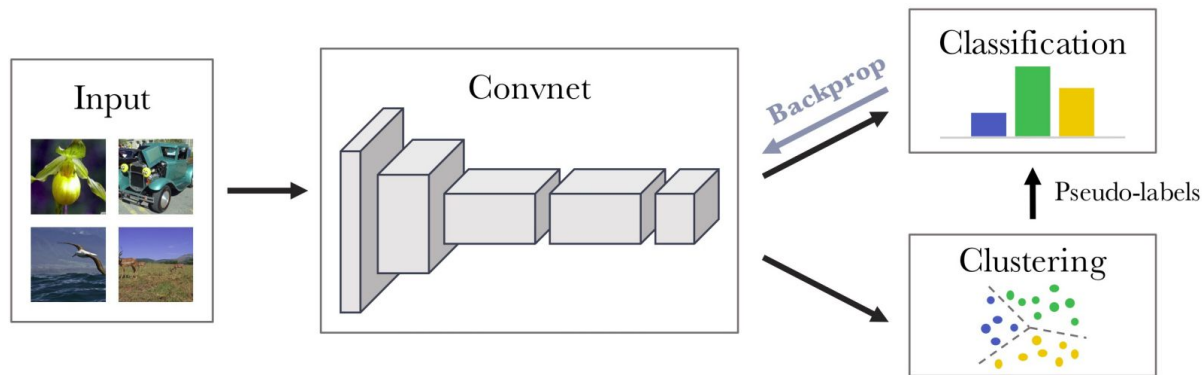Skoltech

# Structure of framework



- modular structure

- all modules are configurable

- different methods are implemented

- single API structure

- built-in data loader, downloader and preprocessor

- properly documented

- transparent pipeline

- user guide inside

Skoltech

# Methods implemented

- Pure Cohortney

  - The main idea of this method is the following: clustering based on the LSTM Model
    We have **K** clusters and **N** realizations of event sequences.
    Given these training sequences, we want to find the clusters among them and the corresponding intensity functions.

- Deep Clustering over Cohortney

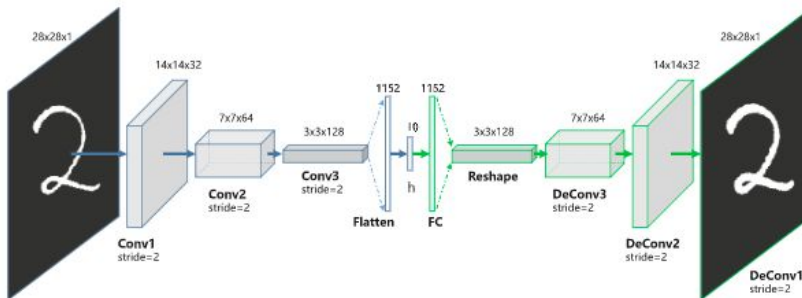- Convolutional Autoencoders over Cohortney

# Methods implemented

- Pure Cohortney

- Deep Clustering over Cohortney
  - Given partitions obtained from Cohortney for a particular **N** we may proceed with clustering algorithms suitable for sequences of features of the same length.
    Therefore we use combination Cohortney with a deep learning approaches for object clustering, such method is called Deep Clustering.

- Convolutional Autoencoders over Cohortney

# Methods implemented

- Pure Cohortney


- Deep Clustering over Cohortney

- Convolutional Autoencoders over Cohortney
  - The clustering of representation with convolutional autoencoders is well-known and widely used method. we apply 1D convolutions over cohortney estimated over event sequences.

# Datasets description

- Synthetic
  - sine-like impact function, which piecewise constant
  - exponential kernel function

- Real Datasets
  - AgeDataset
    - Contains sequences of transaction records stemming from clients at financial institutions. For each client, we have a sequence of transactions. Each transaction is described with a discrete code, identifying the type of transaction, the Merchant Category Code, such as a bookstore, ATM, drugstore, and more. According to transactions' history, dataset is categorized clients into age categories by the authors. The data provides each client's age bin with four different bins in total.
  - LinkedIn
    - Contain users' history of employers. There are 82 different companies, comprising different event types, and with 2439 users in total. Based on this data, we can categorize users into different professional areas. There is information about industry sectors for each user to compare those sectors with real ones. The event sequences in these three datasets have strong but structural triggering patterns, which we can model via different Hawkes processes.

**Skoltech**

[1] Zhuzhel V. et al. COHORTNEY: Deep Clustering for Heterogeneous Event Sequences //arXiv preprint arXiv:2104.01440. – 2021.

# Validation procedure

We follow the evaluation procedure described in the original paper [1]: when the labels are obtained by a certain method, we use one out of two metrics to evaluate the results:

- *Consistency,* or minimum proportion of preserved pairs over trials (for real-world data):

$$\text{Consistency} = \min_{j \in \{1,\ldots,J\}} \sum_{j' \neq j} \sum_{(n,n') \in \mathcal{M},} \frac{1\left\{k_n^{j'} = k_{n'}^{j'}\right\}}{(J-1)\,|\mathcal{M}_j|}.$$

- *Purity,* (for synthetic data) defined as $\quad \text{Purity} = \frac{1}{N} \sum_{k=1}^{K} \max_{j \in \{1,\ldots,K'\}} |\mathcal{W}_k \cap \mathcal{C}_j|$

Here $\quad \mathcal{M}_j = \left\{(n, n') \mid k_n^j = k_{n'}^j\right\}\quad$ is the set of sequence pairs within the same cluster in j-th trial, $k_j^n$ is index of cluster of n-th sequence in the j-th trial, $W_k$ is the learned index set of sequences belonging to k-th cluster, $C_j$ is the real index set of sequences belonging to j-th cluster and N is the total number of sequences

Skoltech

[1] Zhuzhel V. et al. COHORTNEY: Deep Clustering for Heterogeneous Event Sequences //arXiv preprint arXiv:2104.01440. – 2021.

# Results: excerpt

Table 1. Quantitative results of considered methods

| Dataset / Method | CAE | DeepCluster |
|---|---|---|
| K2_C5 | $0.83 \pm 0.09$ | $0.71 \pm 0.06$ |
| K5_C5 | $0.37 \pm 0.04$ | $0.32 \pm 0.11$ |
| LinkedIn | $0.20 \pm 0.01$ | $0.19 \pm 0.02$ |
| sin_K2_C5 | $0.95 \pm 0.02$ | $0.91 \pm 0.06$ |

- Clustering with use of autoencoders over Cohortney is better than with DeepCluster over Cohortney
- The synthetic data is easier to be clusterized
- Increasing the  number of GT clusters affects performance of a model
- Runs of the methods are hyperparameters-dependent

- The results are mainly illustrative
- The columns correspond to methods, the rows correspond to datasets
- The value of Purity metric is on intersection

Skoltech

# Documentation



- The whole project is properly documented

- Documentation provides information about functionality and description of all the modules in a project

- Documentation also consists of guides and pieces of advices for framework's users

Skoltech

# Conclusion

- proposed a framework for event sequences clustering methods which have a lot of industrial applications.
- easy reproducibility for such unique task.
- extends for new dataset without any problems
- adding new methods is now possible without difficulties, as the functionality has become more unified

# Conclusion

- proposed a framework for event sequences clustering methods which have a lot of industrial applications.
- easy reproducibility for such unique task.
- extends for new dataset without any problems
- adding new methods is now possible without difficulties, as the functionality has become more unified

Possible Improvements:
Add More:
- metrics
- methods
- datasets

# Conclusion

- proposed a framework for event sequences clustering methods which have a lot of industrial applications.
- easy reproducibility for such unique task.
- extends for new dataset without any problems
- adding new methods is now possible without difficulties, as the functionality has become more unified

Possible Improvements:
*Add More:*
- metrics
- methods
- datasets

*Takeaways*:
- ML is not about fit and predict
- some tasks has unique specification and difficult to understanding

Skoltech