

Attributed Sequence Embedding

Zhongfang Zhuang, Xiangnan Kong, Elke Rundensteiner
Worcester Polytechnic Institute
{zzhuang, xkong, rundenst}@wpi.edu

Jihane Zouaoui, Aditya Arora
Amadeus IT Group
{jihane.zouaoui, aditya.arora}@amadeus.com

Abstract—Mining tasks over sequential data, such as clickstreams and gene sequences, require a careful design of embeddings usable by learning algorithms. Recent research in feature learning has been extended to sequential data, where each instance consists of a sequence of heterogeneous items with a variable length. However, many real-world applications often involve *attributed sequences*, where each instance is composed of both a sequence of categorical items and a set of attributes. In this paper, we study this new problem of *attributed sequence embedding*, where the goal is to learn the representations of attributed sequences in an unsupervised fashion. This problem is core to many important data mining tasks ranging from user behavior analysis to the clustering of gene sequences. This problem is challenging due to the dependencies between sequences and their associated attributes. We propose a deep multimodal learning framework, called **NAS**, to produce embeddings of attributed sequences. The embeddings are task independent and can be used on various mining tasks of attributed sequences. We demonstrate the effectiveness of our embeddings of attributed sequences in various unsupervised learning tasks on real-world datasets.

Index Terms—Sequence, Embedding, Attributed sequence.

I. INTRODUCTION

Sequential data arise naturally in a wide range of applications [1], [2], [3], [4]. Examples of sequential data include click streams of web users, purchase histories of online customers, and DNA sequences of genes. Different from conventional multidimensional data [5], the sequential data [6] are not represented as feature vectors of continuous values, but as sequences of categorical items with variable-lengths.

Many real-world applications involve mining tasks over sequential data [4], [7], [3]. For example, in online ticketing systems, administrators are interested in finding fraudulent sequences from the clickstreams of users. In user profiling systems, researchers are interested in grouping purchase histories of customers into clusters. Motivated by these real-world applications, sequential data mining has received considerable attention in recent years [2], [1].

Sequential data usually requires a careful design of its embedding before being fed to data mining algorithms. One of the feature learning problems on sequential data is called sequence embedding [8], [9], where the goal is to transform a sequence into a fixed-length embedding.

Conventional methods on sequence embedding focus on learning from sequential data alone [10], [8], [9], [11]. However, in many real-world applications, sequences are often associated with a set of attributes. We define such data as *attributed sequences*, where each instance is represented by a

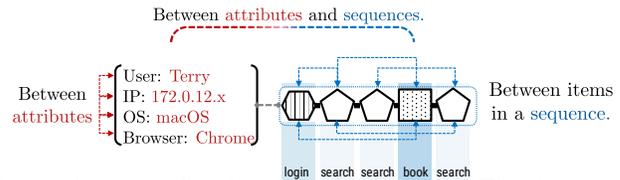


Fig. 1: An example of attributed sequences. The three types of dependencies in an attributed sequence: **item** dependencies, **attribute** dependencies and **attribute-sequence** dependencies.

set of *attributes* associated with a *sequence*. For example, in online ticketing systems as shown in Fig. 1, each user transaction includes both a sequence of user actions (e.g., “login”, “search” and “pick seats”) and a set of attributes (e.g., “user name”, “browser” and “IP address”) indicating the context of the transaction. In gene function analysis, each gene can be represented by both a DNA sequence and a set of attributes indicating the expression levels of the gene in different types of cells. Motivated by the recent success in attributed graph embedding [12], [13], in this paper, we study the problem of attributed sequence embedding.

Building embedding for attributed sequences (as shown in Fig. 2d corresponds to transforming an attributed sequence into a fixed-length embedding with continuous values. Different from the work in [14], [15], we do not have labels for any attributed sequence instances in the embedding task. Sequence embedding problems are particularly challenging with additional attributes. In sequence embedding problems (as shown in Fig. 2a, conventional methods focus on modeling the *item dependencies*, i.e., the dependencies between different items within a sequence. However, in attributed sequences, the dependencies between items can be different if the sequence is observed under different contexts (attributes). Even the same ordering of the items can have different meanings if associated with different attribute values. In this paper, instead of building embeddings to model only the dependencies between items in each single sequence, we aim to model three types of dependencies in an attributed sequence jointly: (1) *item dependencies*, (2) *attribute dependencies* (i.e., the dependencies between different attributes) and (3) *attribute-sequence dependencies* (i.e., the dependencies between attributes and items in a sequence).

Despite its relevance, the problem of producing attributed sequence embeddings in an unsupervised setting remains open. We summarize the major research challenges as follows:

- 1) **Heterogeneous Dependencies.** The bipartite structure of attributed sequences poses unique challenges in feature

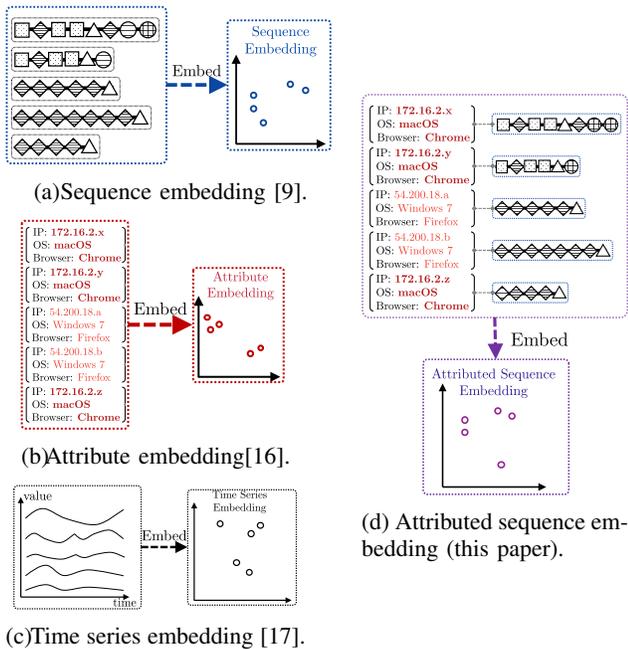


Fig. 2: Comparison of different embedding problems.

learning. As shown in Fig. 1, there exist three types of possible dependencies in an attributed sequence: item dependencies, attribute dependencies and attribute-sequence dependencies.

Motivating Example 1. In Fig. 3, we present an example of fraud detection from a user privilege management system in Amadeus [18]. This system logs each user session as an attributed sequence (denoted as $J_1 \sim J_5$). Each attributed sequence consists of a sequence of user’s activities and a set of attributes derived from metadata values. The attributes (e.g., “IP”, “OS” and “Browser”) are recorded when a user logs into the system and remain unchanged during each user session. We use different shapes and colors to denote different user activities, e.g., “Reset password”, “Delete a user”. In real-world applications like this, the attributes and the associated sequences are already saved within one integrated record. An important step in this fraud detection system is to “red flag” suspicious user sessions for potential security breaches. In Fig. 3, we observe three groups of embeddings learned from the Amadeus application logs. For each group, we use a dendrogram to demonstrate the similarities between embeddings within that group. Neither of the embeddings using only sequences or only attributes detects any outliers due to the lacking of considerations of attribute-sequence dependencies. However, user session J_5 is discovered to be fraudulent using a learning algorithm that incorporates all three types of dependencies.

2) **Lack of Labeled Data.** With the continuously incoming volume of data and the high labor cost of manually labeling data, it is rare to find attributed sequences from many real-world applications with labels (e.g., *fraud*,

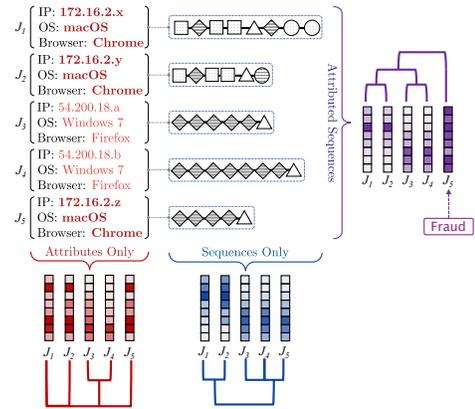


Fig. 3: Dendrograms of embeddings learned from attributed sequences for fraud detection tasks. J_5 is a user committing fraud. However, it is considered a normal user session by the embedding generated using either only attributes or only sequences. J_5 can only be caught as a fraud instance using the embedding learned using both attributes and sequences.

normal) attached. Without proper labels, it is challenging to learn an embedding function that is capable of transforming attributed sequences into compact embeddings concerning the three types of dependencies.

Motivating Example 2. Continuing with our Motivating Example 1, the Amadeus records user activities and their session metadata in the log files. Due to the large volume of entries and complex user sessions, the log files do not have labels depicting whether one user session is fraudulent or not. Only when an embedding function that is capable of transforming unlabeled user sessions $J_1 \sim J_5$ respecting the differences between them, an anomaly detection algorithm can identify J_5 as a fraudulent session.

In this paper, we focus on the *generic* problem of embedding attributed sequences in an unsupervised fashion. We propose a novel framework (called NAS) using deep learning models to address the above challenges. This paper offers the following contributions:

- We study the problem of attributed sequence embedding without any labels available.
- We propose a framework and a training strategy to exploit the dependencies among the attributed sequences.
- We evaluate the embeddings generated by NAS framework on real-world datasets using outlier detection tasks. We also conduct case studies of user behavior analysis and demonstrate the usefulness of NAS in real-world applications.

II. PROBLEM FORMULATION

A. Preliminaries

Definition 1 (Sequence): Given a set of r categorical items $\mathcal{I} = \{e_1, \dots, e_r\}$, the k -th sequence in the dataset $S_k = (\alpha_k^{(1)}, \dots, \alpha_k^{(l_k)})$ is an ordered list of items, where $\alpha_k^{(t)} \in \mathcal{I}, \forall t = 1, \dots, l_k$.

Different sequences can have a varying number of items. For example, the number of user click activities varies between different user sessions. The meanings of items are different in different datasets. For example, in user behavior analysis from clickstreams, each item represents one action in user's click stream (e.g., $\mathcal{I} = \{\text{search}, \text{select}\}$, where $r = 2$). Similarly in DNA sequencing, each item represents one canonical base (e.g., $\mathcal{I} = \{\text{A}, \text{T}, \text{G}, \text{C}\}$, where $r = 4$).

There are dependencies between items in a sequence. Without loss of generality, we use the one-hot encoding of S_k , denoted as $\mathbf{S}_k = (\alpha_k^{(1)}, \dots, \alpha_k^{(l_k)}) \in \mathbb{R}^{l_k \times r}$ where each item $\alpha_k^{(t)} \in \mathbb{R}^r$ in \mathbf{S}_k is a one-hot vector corresponding to the original item $\alpha_k^{(t)}$ in the sequence S_k .

Additionally, each sequence is associated with a set of *attributes*. Each attribute value can be either categorical or numerical. The attribute values are denoted using a vector $\mathbf{x}_k \in \mathbb{R}^u$, where u is the number of attributes in \mathbf{x}_k . For example, in a dataset where each instance has two attributes "IP" and "OS", $u = 2$.

With the attributes and sequences, we now formally define the *attributed sequences* (Def. 2) and the *attribute-sequence dependencies* (Def. 3).

Definition 2 (Attributed Sequence): Given a vector of attribute values \mathbf{x}_k and a sequence \mathbf{S}_k , an attributed sequence $J_k = (\mathbf{x}_k, \mathbf{S}_k)$ is an ordered pair of the attribute value vector \mathbf{x}_k and the sequence \mathbf{S}_k .

Definition 3 (Attribute-Sequence Dependencies): Given an attributed sequence $J_k = (\mathbf{x}_k, \mathbf{S}_k)$, the log likelihood of J_k is $\log \Pr(\mathbf{x}_k, \mathbf{S}_k)$.

B. Problem Definition

The goal of attributed sequence embedding is to learn an embedding function that transforms each attributed sequence with a variable-length sequence of categorical items and a set of attributes into a compact representation in the form of a vector. However, these representations are only valuable if an embedding function is capable of learning all three types of dependencies. Hence, given a set of attributed sequences, we define the learning objective of the embedding function as a minimization of the aggregated negative log likelihood of all three types of dependencies.

Definition 4 (Attributed Sequence Embedding): Given a dataset of attributed sequences $\mathcal{J} = \{J_1, \dots, J_n\}$, the problem of attributed sequence embedding is to find an embedding function Θ with a set of parameters (denoted as θ) that produces embeddings for J_k in the form of vectors. The problem is formulated as:

$$\underset{\theta}{\text{minimize}} \sum_{k=1}^n \sum_{t=1}^{l_k} -\log \Pr(\alpha_k^{(t)} | \beta_k^{(t)}, \mathbf{x}_k) \quad (1)$$

where $\beta_k^{(t)} = (\alpha_k^{(t-1)}, \dots, \alpha_k^{(1)})$, $\forall t = 2, \dots, l_k$ represents the items prior to $\alpha_k^{(t)}$ in the sequence.

Our problem can be interpreted as: we want to minimize the prediction error of the $\alpha_k^{(t)}$ in each attributed sequence given the attribute values \mathbf{x}_k and all the items prior to $\alpha_k^{(t)}$.

III. THE NAS FRAMEWORK

A. Attribute Network

Fully connected neural network [19] is capable of modeling the dependencies of the inputs and at the same time reduce the dimensionality. Fully connected neural network has been widely used [20], [19], [21] for unsupervised data representations learning, including tasks such as dimensionality reduction and generative data modeling. With the high-dimensional sparse input attribute values $\mathbf{x}_k \in \mathbb{R}^u$, it is ideal to use such a network to learn the attribute dependencies. We design our attribute network as:

$$\begin{aligned} \mathbf{V}_k^{(1)} &= \rho(\mathbf{W}_A^{(1)} \mathbf{x}_k + \mathbf{b}_A^{(1)}) \\ &\vdots \\ \mathbf{V}_k^{(M+1)} &= \sigma(\mathbf{W}_A^{(M+1)} \mathbf{V}_k^{(M)} + \mathbf{b}_A^{(M+1)}) \\ &\vdots \\ \widehat{\mathbf{x}}_k &= \sigma(\mathbf{W}_A^{(2M)} \mathbf{V}_k^{(2M-1)} + \mathbf{b}_A^{(2M)}) \end{aligned} \quad (2)$$

where ρ and σ are two activation functions. In this attribute network, we use the ReLU function proposed in [22] (defined as $\rho(z) = \max(0, z)$) and sigmoid function (defined as $\sigma(z) = \frac{1}{1+e^{-z}}$).

The attribute network is an encoder-decoder stack with $2M$ layers, where the first M layers composed of the *encoder* while the next M layers work as the *decoder*. With d_M hidden units in the M -th layer, the input attribute vector $\mathbf{x}_k \in \mathbb{R}^u$ is first transformed to $\mathbf{V}_k^{(M)} \in \mathbb{R}^{d_M}$, $d_M \ll u$ by the encoder. Then the decoder attempts to reconstruct the input and produce the reconstruction result $\widehat{\mathbf{x}}_k \in \mathbb{R}^u$. An ideal attribute network should be able to reconstruct the input from the $\mathbf{V}_k^{(M)}$. The smallest attribute network is built with $M = 1$, where there are one layer of encoder and one layer of decoder.

B. Sequence Network

The proposed sequence network is a variation of the long short-term memory model (LSTM) [23]. The sequence network takes advantage of the conventional LSTM to learn the dependencies between items in sequences. [23] defines the conventional LSTM model is defined as

$$\begin{aligned} \mathbf{i}_k^{(t)} &= \sigma(\mathbf{W}_i \alpha_k^{(t)} + \mathbf{U}_i \mathbf{h}_k^{(t-1)} + \mathbf{b}_i) \\ \mathbf{f}_k^{(t)} &= \sigma(\mathbf{W}_f \alpha_k^{(t)} + \mathbf{U}_f \mathbf{h}_k^{(t-1)} + \mathbf{b}_f) \\ \mathbf{o}_k^{(t)} &= \sigma(\mathbf{W}_o \alpha_k^{(t)} + \mathbf{U}_o \mathbf{h}_k^{(t-1)} + \mathbf{b}_o) \\ \mathbf{g}_k^{(t)} &= \sigma(\mathbf{W}_g \alpha_k^{(t)} + \mathbf{U}_g \mathbf{h}_k^{(t-1)} + \mathbf{b}_g) \\ \mathbf{c}_k^{(t)} &= \mathbf{f}_k^{(t)} \odot \mathbf{c}_k^{(t-1)} + \mathbf{i}_k^{(t)} \odot \mathbf{g}_k^{(t)} \\ \mathbf{h}_k^{(t)} &= \mathbf{o}_k^{(t)} \odot \tanh(\mathbf{c}_k^{(t)}) \end{aligned} \quad (3)$$

where \odot denotes element-wise product, σ is a sigmoid activation function, $\mathbf{i}_k^{(t)}$, $\mathbf{f}_k^{(t)}$, $\mathbf{o}_k^{(t)}$ and $\mathbf{g}_k^{(t)}$ are the internal gates of an LSTM. The cell states (denoted as $\mathbf{c}_k^{(t)}$) and hidden states (denoted as $\mathbf{h}_k^{(t)}$), which store the information of the sequential data, are two important components in the LSTM model. Without loss of generality, we denote the dimension of the cell states and the hidden states as d .

Integration of Attribute Network and Sequence Network. Different from the conventional LSTM, our proposed sequence network also accepts the output from the attribute network to condition the sequence network. In particular, we have redesigned the function of the *hidden states* to integrate the information from the attribute network by conditioning the sequence network at the first time step as

$$\mathbf{h}_k^{(t)} = \mathbf{o}_k^{(t)} \odot \tanh(\mathbf{c}_k^{(t)}) + \mathbb{1}(t=1) \odot \mathbf{V}_k^{(M)} \quad (4)$$

This integration requires the attribute network and the sequence network have the same number of the hidden units (*i.e.*, $d_M = d$).

Since the attributed sequences are unlabeled, we designed the sequence network to predict *the next item in the sequence* as the training strategy. The prediction is carried out by designing an output layer applying a `softmax` function on the hidden states as

$$\mathbf{y}_k^{(t)} = \delta(\mathbf{W}_y \mathbf{h}_k^{(t)} + \mathbf{b}_y) \quad (5)$$

where $\mathbf{y}_k^{(t)} \in \mathbb{R}^r$ is the predicted next item in sequence computed using `softmax` function, \mathbf{W}_y and \mathbf{b}_y are the weights and bias of this output layer. With the `softmax` activation function, the $\mathbf{y}_k^{(t)}$ can be interpreted as the probability distribution over r items.

C. Training

1) *Training Objectives:* We use two different learning objectives for attribute network and sequence network targeting at the unique characteristics of attribute and sequence data.

1) Attribute network aims at minimizing the differences between the input and reconstructed attribute values. The learning objective function of attribute network is defined as

$$L_A = \|\mathbf{x}_k - \widehat{\mathbf{x}}_k\|_2^2 \quad (6)$$

2) Sequence network aims at minimizing log likelihood of incorrect prediction of the next item at each time step. Thus, the sequence network learning objective function can be formulated using categorical cross-entropy as

$$L_S = - \sum_{t=1}^{l_k} \alpha_k^{(t)} \log \mathbf{y}_k^{(t)} \quad (7)$$

2) *Embedding:* After the model is trained, we use the parameters in attribute network and sequence network to embed each attributed sequence. Specifically, the attributed sequences are used as inputs to the *trained* model only with the one forward pass, where the parameters within the model remain unchanged. After the last time step for an attributed sequence \mathbf{S}_k , the cell state of sequence network, namely $\mathbf{c}_k^{(l_k)}$, is used as the embedding of \mathbf{S}_k .

IV. EXPERIMENTAL EVALUATION

In this section, we evaluate NAS framework using real-world application logs from Amadeus and public datasets from Wikipedia [24], [25]. We evaluate the quality of embeddings generated by different methods by measuring the

performance of outlier detection algorithms using different embeddings.

A. Experimental Setup

1) *Data Collection:* We use four datasets in our experiments: two from Amadeus application log files and two from Wikipedia¹. The numbers of attributed sequences in all four datasets vary between $\sim 58k$ and $\sim 106k$.

- AMS-A/B: We extract $\sim 164k$ instances from log files of an Amadeus internal application. Each record is composed of a profile containing information ranging from system configurations to office name, and a sequence of functions invoked by click activities on the web interface. There are 288 distinct functions, 57,270 distinct profiles in this dataset. The average length of the sequences is 18.
- WIKI-A/B: This dataset is sampled from Wikipedia dataset. Wikipedia dataset originated from a human-computation game, called Wikipedia [25]. We use a subset of $\sim 3.5k$ paths from Wikipedia with the average length of the path as 6. We also extract 11 sequence context (*e.g.*, the category of the source page, average time spent on each page) as attributes.

2) *Compared Methods:* To study NAS performance on attributed sequences in real-world applications, we use the following compared methods in our experiments.

- LEN [26]: The attributes are encoded and directly used in the mining algorithm.
- MCC [27]: MCC uses the sequence component of attributed sequence as input and produces log likelihood for each sequence.
- SEQ [9]: Only the sequence inputs are used by an LSTM to generate fixed-length embeddings.
- ATR [16]: A two-layered fully connected neural network is used to generate attribute embeddings.
- EML[28]: Aggregate MCC and LEN scores.
- CSA [29]: The attribute embedding and the sequence embedding are independently generated by ATR and SEQ, then concatenated together.

3) *Network Parameters:* Following the previous work in [30], we initialize weight matrices \mathbf{W}_A and \mathbf{W}_S using the uniform distribution. The recurrent matrix \mathbf{U}_S is initialized using the orthogonal matrix as suggested by [31]. All the bias vectors are initialized with zero vector $\mathbf{0}$. We use stochastic gradient descent as optimizer with the learning rate of 0.01. We use a two-layer encoder-decoder stack as our attribute network.

B. Outlier Detection Tasks

We use outlier detection tasks to evaluate the quality of embeddings produced by NAS. We select k -NN outlier detection algorithm as it has only one important parameter (*i.e.*, the k value). We use ROC AUC as the metric in this set of experiments.

For each of the AMS-A and AMS-B datasets, we ask domain experts to select two users as inlier and outlier. These

¹Personal identity information is not collected for privacy reasons.

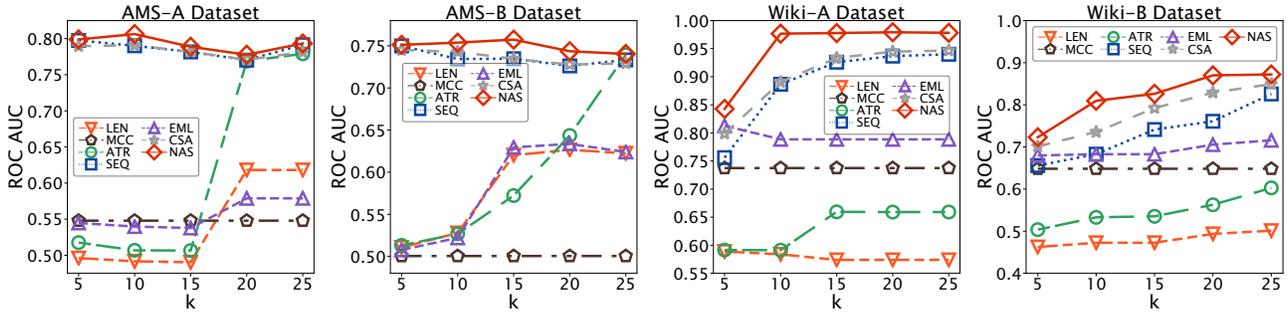


Fig. 4: Parameter sensitivity to different k values. It is shown that the embeddings generated by NAS always have the best performance under different k values.

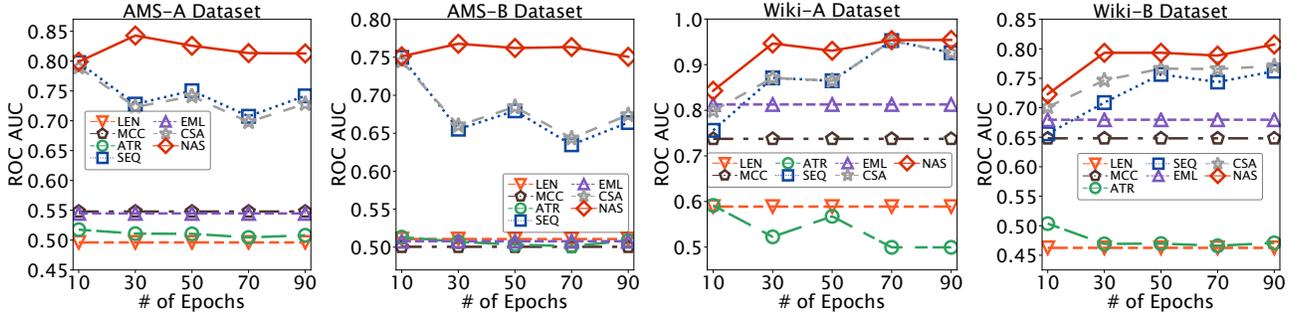


Fig. 5: Performance comparisons using outlier detection tasks. The embeddings generated by NAS can always achieve the best performance compared to baseline methods when the number of training epochs increases.

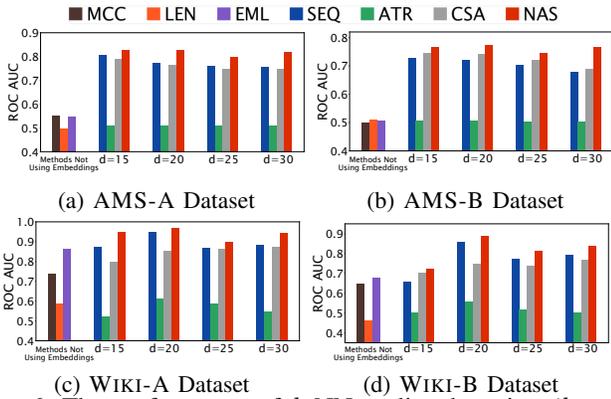


Fig. 6: The performance of k -NN outlier detection ($k = 5$). The *methods not using embeddings* are placed on the left. We vary the number of dimensions on the right. The higher score is better. We observe that the combinations of k -NN and NAS embeddings have the best performance on the four datasets.

two users have completely different behaviors (*i.e.*, sequences) and metadata (*i.e.*, attributes). The percentages of outliers in AMS-A and AMS-B are 1.5% and 2.5% of all attributed sequences, respectively. For the WIKI-A and WIKI-B datasets, each path is labeled based on the category of the source page. Similarly to the previous two datasets, we select paths with two labels as inliers and outliers where the percentage of outlier paths is 2%. The feature used to label paths is excluded from the learning and embedding process.

Performance. The goal of this set of experiments is to demonstrate the performance of outlier detection using all our compared methods. Each method is trained with all the instances. For SEQ, ATR and NAS, the number of learning epochs is set to 10 and we vary the number of embedding dimensions d from 15 to 30. We set $k = 5$ for outlier detection tasks for LEN, SEQ, ATR, CSA and NAS. Choosing the *optimal* k value in the outlier detection tasks is beyond the scope of this work, thus we omit its discussions. We summarize the performance results in Fig. 6.

Analysis. We find that the results based on the embeddings generated by NAS are superior to other methods. That is, NAS maximally outperforms other state-of-the-art algorithms by 32.9%, 27.5%, 44.8% and 48% on AMS-A, AMS-B, WIKI-A and WIKI-B datasets, respectively. It is worth mentioning that NAS outperforms a similar baseline method CSA by incorporating the information of attribute-sequence dependencies.

Parameter Study There are two *key* parameters in our evaluations, *i.e.*, k value for the k -NN algorithm and the learning epochs.

In Fig. 4, we first show that the embeddings (dimension $d = 15$) generated by our NAS assist k -NN outlier detection algorithm to achieve superior performance under a wide range of k values ($k = 5, 10, 15, 20, 25$). We omit the detailed discussions of selecting the optimal k values due to the scope of this work.

Fig. 5 presents the results when we fix $k = 5, d = 15$ and vary the number of epochs in the learning phase. We notice

that compared to its competitor, the embeddings generated by NAS can achieve a higher AUC even with a relatively fewer number of learning epochs. Compared to other neural network-based methods (*i.e.*, SEQ, ATR and CSA), NAS have a more stable performance. The NAS performance gain is not due to the advantage of using both attributes and sequences, but because of taking the various dependencies into account, as the other two competitors (*i.e.*, CSA and EML) also utilize the information from both attributes and sequences.

V. RELATED WORK

Sequence Mining. Many sequence mining work focuses on frequent sequence pattern mining. Recent work in [2] targets finding subsequences of possible non-consecutive actions constrained by a gap within sequences. [32] aims at solving pattern-based sequence classification problems using a parameter-free algorithm from the model space. It defines rule pattern models and a prior distribution on the model space. [33] builds a subsequence interleaving model for mining the most relevant sequential patterns.

Deep Learning. Sequence-to-sequence learning in [9] uses long short-term memory model in machine translation. The hidden representations of sentences in the source language are transferred to a decoder to reconstruct in the target language. The idea is that the hidden representation can be used as a compact representation to transfer sequence similarities between two sequences. Multi-task learning in [11] examines three multi-task learning settings for sequence-to-sequence models that aim at sharing either an encoder or decoder in an encoder-decoder model setting. Although the above work is capable of learning the dependencies within a sequence, none of them focuses on learning the dependencies between attributes and sequences. This new bipartite data type of attributed sequence has posed new challenges of heterogeneous dependencies to sequence models, such as RNN and LSTM. Multimodal deep neural networks [34], [29], [35] is designed for information sharing across multiple neural networks, but none of these work focuses on our attributed sequence embedding problem.

VI. CONCLUSION

In this paper, we study the problem of *unsupervised attributed sequences embedding*. Different from conventional feature learning approaches, which work on either sequences or attributes without considering the *attribute-sequence dependencies*, we identify the three types of dependencies in attributed sequences. We propose a novel framework, called NAS, to learn the heterogeneous dependencies and embed unlabeled attributed sequences. Empirical studies on real-world tasks demonstrate that the proposed NAS effectively boosts the performance of outlier detection tasks compared to baseline methods.

REFERENCES

[1] N. Béchet, P. Cellier *et al.*, “Sequence mining under multiple constraints,” in *SAC*. ACM, 2015.

[2] I. Miliaraki, K. Berberich *et al.*, “Mind the gap: Large-scale frequent sequence mining,” in *SIGMOD*. ACM, 2013.

[3] G. Wang, X. Zhang *et al.*, “Unsupervised clickstream clustering for user behavior analysis,” in *CHI*. ACM, 2016.

[4] W. Wei, J. Li *et al.*, “Effective detection of sophisticated online banking fraud on extremely imbalanced data,” *WWW*, 2013.

[5] T. B. Pedersen and C. S. Jensen, “Multidimensional data modeling for complex data,” in *ICDE*. IEEE, 1999.

[6] J. Yang and W. Wang, “Cluseq: Efficient and effective sequence clustering,” in *ICDE*. IEEE, 2003.

[7] A. Tajer, V. V. Veeravalli, and H. V. Poor, “Outlying sequence detection in large data sets: A data-driven approach,” *IEEE Signal Processing Magazine*, 2014.

[8] K. Cho, B. Van Merriënboer *et al.*, “Learning phrase representations using rnn encoder-decoder for statistical machine translation,” *arXiv preprint arXiv:1406.1078*, 2014.

[9] I. Sutskever, O. Vinyals, and Q. V. Le, “Sequence to sequence learning with neural networks,” in *NIPS*, 2014.

[10] N. Kalchbrenner and P. Blunsom, “Recurrent continuous translation models,” in *EMNLP*, 2013.

[11] M.-T. Luong, Q. V. Le *et al.*, “Multi-task sequence to sequence learning,” *arXiv preprint arXiv:1511.06114*, 2015.

[12] J. Gibert, E. Valveny, and H. Bunke, “Graph embedding in vector spaces by node attribute statistics,” *Pattern Recognition*, 2012.

[13] B. Perozzi, R. Al-Rfou, and S. Skiena, “Deepwalk: Online learning of social representations,” in *SIGKDD*. ACM, 2014.

[14] Z. Zhuang, X. Kong *et al.*, “One-shot learning on attributed sequences,” in *Big Data*. IEEE, 2018, pp. 921–930.

[15] Z. Zhuang, X. Kong, and E. Rundensteiner, “Amas: A attention model for a attributed sequence classification,” in *SDM*, 2019, pp. 109–117.

[16] W. Wang, Y. Huang *et al.*, “Generalized autoencoder: A neural network framework for dimensionality reduction,” in *CVPR Workshops*, 2014.

[17] A. Khaleghi, D. Ryabko *et al.*, “Consistent algorithms for clustering time series,” *JMLR*, 2016.

[18] Amadeus, “Amadeus IT Group,” <http://www.amadeus.com>, 2017, accessed: 2017-09-23.

[19] C.-Y. Liou, W.-C. Cheng *et al.*, “Autoencoder for words,” *Neurocomputing*, 2014.

[20] C.-Y. Liou, J.-C. Huang, and W.-C. Yang, “Modeling word perception using the elman network,” *Neurocomputing*, 2008.

[21] N. Phan, Y. Wang *et al.*, “Differential privacy preservation for deep auto-encoders: an application of human behavior prediction,” in *AAAI*, 2016.

[22] V. Nair and G. E. Hinton, “Rectified linear units improve restricted boltzmann machines,” in *ICML*, 2010.

[23] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, 1997.

[24] R. West and J. Leskovec, “Human wayfinding in information networks,” in *WWW*. ACM, 2012.

[25] R. West, J. Pineau, and D. Precup, “Wikispeedia: An online game for inferring semantic distances between concepts,” in *IJCAI*, 2009.

[26] Z. Akata, F. Perronnin *et al.*, “Label-embedding for attribute-based classification,” in *CVPR*, 2013.

[27] S. D. Bernhard, C. K. Leung *et al.*, “Clickstream prediction using sequential stream mining techniques with markov chains,” in *IDEAS*. ACM, 2016.

[28] R. R. Yager and N. Alajlan, “Probabilistically weighted owa aggregation,” *IEEE Transactions on Fuzzy Systems*, 2014.

[29] J. Ngiam, A. Khosla *et al.*, “Multimodal deep learning,” in *ICML*, 2011.

[30] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *AISTATS*, 2010.

[31] A. M. Saxe, J. L. McClelland, and S. Ganguli, “Exact solutions to the nonlinear dynamics of learning in deep linear neural networks,” *arXiv preprint arXiv:1312.6120*, 2013.

[32] E. Egho, D. Gay *et al.*, “A parameter-free approach for mining robust sequential classification rules,” in *ICDM*. IEEE, 2015.

[33] J. Fowkes and C. Sutton, “A subsequence interleaving model for sequential pattern mining,” *arXiv preprint arXiv:1602.05012*, 2016.

[34] A. Karpathy and L. Fei-Fei, “Deep visual-semantic alignments for generating image descriptions,” in *CVPR*, 2015.

[35] K. Xu, J. Ba *et al.*, “Show, attend and tell: Neural image caption generation with visual attention,” in *ICML*, 2015.