#### Modeling Users' Behavior Sequences with Hierarchical Explainable Network for Cross-domain Fraud Detection

# **Yongchun Zhu<sup>1, 2, 3</sup>**, Dongbo Xi<sup>1, 2, 3</sup>, Bowen Song<sup>3</sup>, Fuzhen Zhuang<sup>1, 2</sup>, Shuai Chen<sup>3</sup>, Xi Gu<sup>3</sup>, Qing He<sup>1, 2</sup>

1. Key Lab of Intelligent Information Processing of Chinese Academy of Sciences (CAS), Institute of Computing

Technology, CAS, Beijing, China

2. University of Chinese Academy of Sciences, Beijing, China

3. Alipay (Hangzhou) Information & Technology Co., Ltd.



### Motivation



Sequence-based Fraud Detection Task: Exploit the users' historical behavior sequences to help differentiate fraudulent payments from regular ones.



# Motivation

Existing Method



Shortcommings:

- The prediction results are difficult to explain for these methods.
- They focus more on the sequential information of the behaviors, but fail to thoroughly exploit the internal information of each behavior, e.g., only the first-order information of fields' embeddings is used to represent events.

In order to tackle the above two problems, we propose a Hierarchical Explainable Network (HEN) to model users' behavior sequences, which could not only improve the performance of fraud detection but also help answer this "why".





### Motivation

New challenge:

Our company is one of the world-leading cross-border e-commerce companies. As its e-commerce business expands to new domains, e.g., new countries or new markets, one major problem for modeling users' behavior in fraud detection systems is the limitation of data collection, e.g., very few data/labels available. (The need of Lazada)

Due to the new challenge, we introduce **cross-domain fraud detection**, which aims to transfer knowledge from existing domains (source domains) with enough and mature data to improve the performance in the new domain (target domain).



## Contributions

- To solve the sequence-based fraud detection problem, we propose hierarchical explainable network (HEN) to model users' behavior sequences, which could not only improve the performance of fraud detection but also give reasonable explanations for the prediction results.
- For cross-domain fraud detection, we propose a general transfer framework that can be applied upon various existing models in the Embedding & MLP paradigm.
- We perform experiments on real-world datasets of four countries to demonstrate the effectiveness of HEN. In addition, we demonstrate our transfer framework is general for various existing models on 90 transfer tasks. Finally, we conduct case study to prove the explainability of HEN and the transfer framework.



5



#### 中國科学院计算技术研究码 INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES



#### Problem Statement:

Event sequence:

$$E = [e_1, e_2, ..., e_T]$$

Event presentation:

$$\boldsymbol{e}_t = [x_1^t, x_2^t, ..., x_n^t]$$

#### Structure:

- Look-up embedding
- Field-level extractor
- Event-level extractor
- □ Wide layer
- □ MLP

1、Look-up embedding

Look-up embedding has been widely adopted to learn dense representations from raw data for prediction.

$$\boldsymbol{v}_i^t = \begin{cases} \Phi_i[x_i^t] & \text{the } i\text{-th field is categorical field} \\ x_i^t \times \Phi_i & \text{the } i\text{-th field is numerical field} \end{cases},$$

- 2、Field-level extractor
- The field-level extractor aims to extract the event representations from the embedding vectors of fields.  $e_t = \left[\sum_{i=1}^n w_i^t v_i^t + \sum_{i=1}^n \sum_{j=i+1}^n v_i^t \odot v_j^t, \qquad w_i^t = \frac{\exp(a_{x_i^t}^i)}{\sum_{j=1}^n \exp(a_{x_i^t}^j)}, \frac{1}{\sum_{j=1}^n \exp(a_{x_i^t}^j)}\right]$



3, Event-level extractor

The event-level extractor aims to extract the sequence embedding from the historical event embedding vectors.  $T_{-1}$ 

$$s = \sum_{t=1}^{T-1} u_t f_1(e_t), \quad \hat{u}_t = \frac{\langle f_2(e_t), f_3(e_t) \rangle}{\sqrt{k}}, u_t = \frac{\exp\left(\hat{u}_t\right)}{\sum_{j=1}^{T-1} \exp\left(\hat{u}_j\right)},$$

4、Wide layer

Wide layer is the linear part just like the "wide" part in Wide & Deep to capture the first-order information.

$$l(e_T) = \sum_{i=1}^{n} c_i(x_i^T) + c_0,$$

5. Prediction and learning

 $\hat{y} = Sigmoid (MLP([s, e_T]) + l(e_T)),$  $L(\theta) = -\frac{1}{N} \sum_{(E,y)\in\mathcal{D}} (y\log\hat{y} + (1-y)\log(1-\hat{y})),$ 

#### Explainability:

- Field-level extractor: Our field-level extractor is able to choose informative fields. w indicates the attention distribution of field embeddings that could explain which field embedding is more important to represent event embedding.
- Event-level extractor: Our event-level extractor is able to score the event importance, and from the score, we could find which event is more important to represent the sequence embedding.
- Wide layer: The wide layer of HEN could help to identify specific value of fields with high-risk/lowrisk, which could be used as whitelist/blacklist.







Figure 3: The structure of the general transfer framework for cross-domain fraud detection.





#### Problem Statement:

There are a source domain and a target domain. The number of labeled samples in target is far less than the source domain.

#### Structure:

- □ The Strategy of Embedding
- Shared and Specific behavior sequence extractor
- Domain Attention
- Aligning Distributions

### Explainability:

 Domain Attention: Domain attention module is capable of learning the importance of domain-shared and domainspecific representations.

- 1、The Strategy of Embedding
- □ We divide the embedding layer into domain-shared and domain-specific.
- 2. Shared and Specific behavior sequence extractor:
- □ We also divide the behavior sequence extractor into domain-shared and domain-specific.
- 3、 Domain Attention
- By combining the domain-shared and domain-specific representations ( $z_{share}$  and  $z_{spe}$ ), it could contain more useful knowledge. (o denotes the domain label)

$$z^{o} = b^{o}_{share} g_{1}(z^{o}_{share}) + b^{o}_{spe} g_{1}(z^{o}_{spe}), \qquad \hat{b}^{o}_{p} = \frac{\langle g_{2}(z^{o}_{p}), g_{3}(z^{o}_{p}) \rangle}{\sqrt{k}}, \\ b^{o}_{p} = \frac{\exp(\hat{b}^{o}_{p})}{\exp(\hat{b}^{o}_{share}) + \exp(\hat{b}^{o}_{spe})}, \qquad \hat{b}^{o}_{p} = \frac{\exp(\hat{b}^{o}_{p})}{\exp(\hat{b}^{o}_{share}) + \exp(\hat{b}^{o}_{spe})},$$





#### 4、 Aligning Distributions

To feed  $z^{src}$  and  $z^{tgt}$  into the same MLP, we need to align the distibutions of them. Most existing transfer approaches aim to align the marginal and conditional distributions. However, in our scenario, the class distribution is extremely unbalanced that the number of non-fraud samples is about 100 times more than fraud samples. Hence, aligning the marginal and conditional distributions would lead to unsatisfying performance. For our scene, we propose Class-aware Euclidean Distance which explicitly takes the class information into account and measures the intra-class and inter-class discrepancy across domains.

Euclidean Distance: 
$$d(\mathcal{D}_{c_1}^{o_1}, \mathcal{D}_{c_2}^{o_2}) = \left\| \frac{1}{N_{c_1}^{o_1}} \sum_{i=1}^{N_{c_1}} z_i^{o_1} - \frac{1}{N_{c_2}^{o_2}} \sum_{i=1}^{N_{c_2}^{o_2}} z_i^{o_2} \right\|,$$

Class-aware Euclidean Distance:  $\Delta(\mathcal{D}^{src}, \mathcal{D}^{tgt}) = \frac{\sum_{c \in \{0,1\}} d(\mathcal{D}_{c}^{src}, \mathcal{D}_{c}^{tgt})}{\sum_{o_{1} \in \{src,tgt\}} \sum_{o_{2} \in \{src,tgt\}} d(\mathcal{D}_{c=0}^{o_{1}}, \mathcal{D}_{c=1}^{o_{2}})}$ 







The general transfer framework could be applied upon various existing models in the Embedding & MLP paradigm. The most important thing to apply the transfer framework is to define the behavior sequence extractor. For our HEN, the behavior sequence extractor contains a field-level extractor and an eventlevel extractor as shown in the left.

Dataset:

| Dataset    | #pos | #neg  | #pos ratio | #fields | #events | #avglen |
|------------|------|-------|------------|---------|---------|---------|
| C1         | 10K  | 1.93M | 5.2‰       | 56      | 3.57M   | 4.94    |
| C2         | 4.2K | 1.48M | 2.8‰       | 56      | 3.91M   | 4.73    |
| C3         | 15K  | 1.37M | 10.8‰      | 56      | 4.28M   | 14.97   |
| <b>C</b> 4 | 5.7K | 174K  | 31.7‰      | 56      | 353K    | 4.91    |

The fraud detection dataset is collected from one of the world-leading cross-border ecommerce company, which utilizes the risk management system to detect the transaction frauds.



### Experiments Baselines

- □ W & D[1]:In real industrial applications, Wide & Deep model has been widely accepted.
- □ NFM[2]: It is a recent state-of-the-art simple and efficient neural factorization machine model.
- □ LSTM4FD[3]: Some work has applied LSTM for the sequence-based fraud detection tasks, and we called these methods as LSTM4FD.
- □ M3R[4]: It is a most recent hierarchical sequence-based model (M3R and M3C) which deals with both short-term and long-term dependencies with mixture models. We choose the better hierarchical model M3R as the baseline.

Cheng, Heng-Tze, et al. "2016. Wide & deep learning for recommender systems." In DLRS, 2016.
He, Xiangnan, et al. "Neural factorization machines for sparse predictive analytics." In SIGIR, 2017.
Wang S, et al. "Session-based fraud detection in online e-commerce transactions using recurrent neural networks"//ECML&EKDD, 2017.
Tang, Jiaxi, et al. "Towards neural mixture recommender for long range dependent user sequences." In WWW, 2019.



In binary prediction tasks, AUC (Area Under ROC) is a widely used metric. However, in our real card-stolen fraud detection scenario, we should increase the recall rate, while avoid disturbing the normal users as few as possible. In other words, the task is improving the True Positive Rate (TPR) on the basis of low False Positive Rate (FPR). Therefore, we adopt the standardized partial AUC (SPAUC<sub>FPR≤maxfpr</sub>)

$$\begin{aligned} \text{SPAUC}_{\text{FPR} \leq \text{maxfpr}} &= \frac{1}{2} \left( 1 + \frac{\text{AUC}_{\text{FPR} \leq \text{maxfpr}} - \text{minarea}}{\text{maxarea} - \text{minarea}} \right), \\ \text{where} \quad \text{maxarea} &= \text{maxfpr}, \\ \text{minarea} &= \frac{1}{2} \times \text{maxfpr}^2. \end{aligned}$$





#### Standard supervised prediction tasks.



3

4

5

Transfer tasks



We observe that the transfer framework is able to improve the performance of base models with different sizes of training sets, which proves that the transfer framework is compatible with many models in the Embedding & MLP paradigm.

#### Abalation study:

#### Table 2: Results (SPAUC<sub>FPR $\leq 1\%$ </sub>) of ablation study on transfer tasks from dataset C1 to dataset C4 based with HEN.

| Methods       | 3days                 | 1week               | 2weeks              | 3weeks                | 4weeks                | 5weeks                | avg    |
|---------------|-----------------------|---------------------|---------------------|-----------------------|-----------------------|-----------------------|--------|
| target        | $0.8534 {\pm} 0.0173$ | $0.8811 \pm 0.0232$ | $0.8971 \pm 0.0111$ | $0.9009 \pm 0.0117$   | $0.9105 \pm 0.0099$   | $0.9050 \pm 0.0171$   | 0.8913 |
| source        | $0.5724 {\pm} 0.0230$ | $0.5732 \pm 0.0272$ | $0.5909 \pm 0.0576$ | $0.5877 {\pm} 0.0334$ | $0.5474 {\pm} 0.0393$ | $0.5705 \pm 0.0419$   | 0.5737 |
| pretrain      | $0.8639 \pm 0.0047$   | $0.8920 \pm 0.0066$ | $0.8983 \pm 0.0115$ | $0.9120 \pm 0.0051$   | $0.9172 \pm 0.0064$   | $0.9191 \pm 0.0054$   | 0.9004 |
| domain-shared | $0.8811 \pm 0.0110$   | $0.8995 \pm 0.0119$ | $0.9025 \pm 0.0153$ | $0.9157 \pm 0.0109$   | $0.9204 \pm 0.0068$   | $0.9194 \pm 0.0067$   | 0.9065 |
| our structure | $0.8801 {\pm} 0.0101$ | $0.9071 \pm 0.0042$ | $0.9038 \pm 0.0085$ | $0.9152 {\pm} 0.0094$ | $0.9226 \pm 0.0081$   | $0.9233 \pm 0.0113$   | 0.9087 |
| Coral         | $0.8566 \pm 0.0156$   | $0.8938 \pm 0.0096$ | $0.9028 \pm 0.0112$ | $0.9163 \pm 0.0128$   | $0.9236 \pm 0.0080$   | $0.9221 \pm 0.0081$   | 0.9025 |
| Adversarial   | $0.8763 \pm 0.0113$   | $0.8963 \pm 0.0135$ | $0.8951 \pm 0.0109$ | $0.9116 \pm 0.0099$   | $0.9289 \pm 0.0053$   | $0.9260 \pm 0.0069$   | 0.9057 |
| MMD           | $0.8744 {\pm} 0.0107$ | $0.8964 \pm 0.0105$ | $0.9049 \pm 0.0091$ | $0.9185 \pm 0.0058$   | $0.9249 \pm 0.0075$   | $0.9213 \pm 0.0057$   | 0.9067 |
| CMMD          | $0.8797 \pm 0.0136$   | $0.8897 \pm 0.0132$ | $0.8963 \pm 0.0145$ | 0.9121±0.0135         | $0.9219 \pm 0.0069$   | $0.9146 \pm 0.0122$   | 0.9024 |
| ED            | $0.8767 {\pm} 0.0148$ | $0.9010 \pm 0.0088$ | $0.9180 \pm 0.0063$ | $0.9186 \pm 0.0073$   | $0.9258 \pm 0.0076$   | $0.9251 \pm 0.0046$   | 0.9109 |
| CED           | $0.8866 \pm 0.0091$   | $0.9076 \pm 0.0073$ | $0.9140 \pm 0.0059$ | $0.9220 \pm 0.0070$   | $0.9305 \pm 0.0076$   | $0.9298 \pm 0.0059$   | 0.9151 |
| ours          | $0.9038 \pm 0.0052$   | $0.9188 \pm 0.0052$ | 0.9242±0.0063       | 0.9283±0.0039         | $0.9367 \pm 0.0041$   | <b>0.9340</b> ±0.0059 | 0.9243 |





Case Study:

#### Explanation of Wide Layer:

|     |               | Email   | Card   | Issuer  |  |
|-----|---------------|---|--|---|--|
| er: | High-<br>Risk | email1(123/131)<br>email2(27/27)<br>email3(54/59) | card1(33/33)<br>card2(42/54)<br>card3(77/78)   | issuer1(602/607)<br>issuer2(76/108)<br>issuer3(77/90) |  |
|     | Low-<br>Risk  | email4(0/382)<br>email5(0/298)<br>email6(0/471)   | card4(0/2365)<br>card5(0/245)<br>card6(0/5972) | issuer4(27/12491)<br>issuer5(0/789)<br>issuer6(1/725) |  |







True label: 1(fraud) Pred:0.97

Target-specife: 0.44







#### Name: Yongchun Zhu | Email: zhuyongchun18s@ict.ac.cn



