



# Mining Precise-positioning Episode Rules from Event Sequences

Xiang Ao<sup>1</sup>, Ping Luo<sup>1</sup>, Jin Wang<sup>2</sup>, Fuzhen Zhuang<sup>1</sup>, Qing He<sup>1</sup>

中国科学院  
INSTITUTE OF COMPUTING TECHNOLOGY

Institute of Computing Technology, CAS, China<sup>1</sup>

University of California at Los Angeles, USA<sup>2</sup>



## MOTIVATION

### Traditional Episode Rule

Given a frequent episode  $\alpha$ , a **traditional episode rule** in the form of  $lhs \rightarrow rhs$  is generated straightforwardly: The antecedent  $lhs$  is the prefix of  $\alpha$  and the consequent  $rhs$  is the last event in  $\alpha$ , if its confidence is larger than a user-specified threshold.

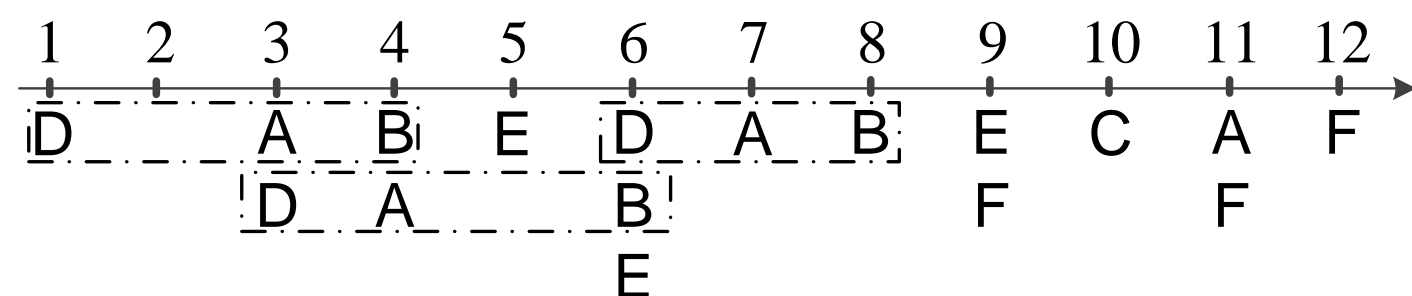


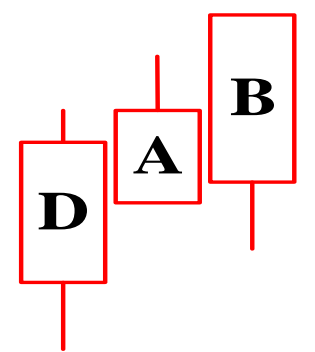
Fig.1 The running example event sequence.

From Fig.1,  $\langle D, A \rangle \rightarrow \langle B \rangle$  is a **traditional episode rule** which indicates it is **within 2 time intervals** after the occurrence of  $\langle D, A \rangle$  that B will occur (with 100% confidence).

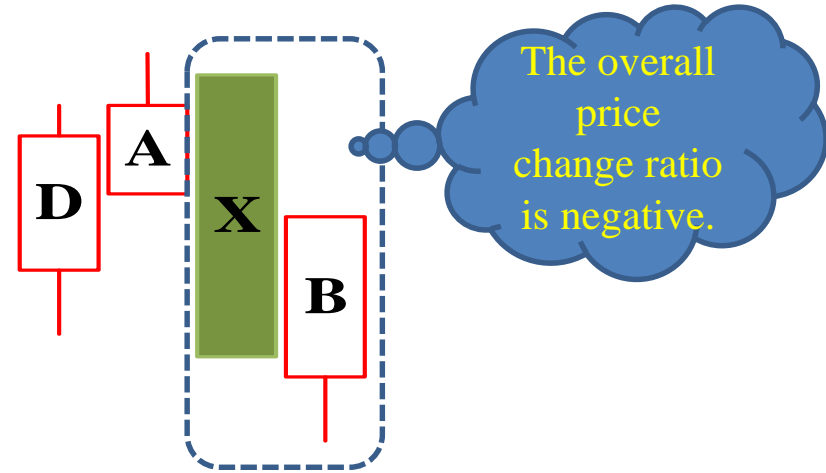
### Limitation of Traditional Episode Rule

**Example:** In stock investment application, we can map price change ratios to events and use candlestick charts to represent events. **Red bars denote price increase** of a stock, and **green bars denote prices decrease**.

- The episode rule  $\langle D, A \rangle \rightarrow \langle B \rangle$  **predicts correct** in the following two cases, however **we will lose money in Case 2** if we long the stock after we observed the antecedent of the rule.



Case 1: B occurs right behind  $\langle D, A \rangle$ .



Case 2: B occurs after  $\langle D, A \rangle$  appears within two days but right behind a significant decrease.

### Precise-positioning Episode Rule (PER)

We define **precise-positioning episode rule** in the form of:

$$\Gamma = \alpha \xrightarrow{\Delta t} \beta$$

$\alpha$ : a **traditional episode**, as the antecedent;  
 $\beta$ : a **fixed-gap episode**, as the consequent;  
 $\Delta t$ : the **time constraint** between the antecedent and the consequent.

**Fixed-gap episode:**

$$\beta = (\langle e_{\beta_1}, \dots, e_{\beta_k} \rangle, \langle \Delta t_1, \dots, \Delta t_{k-1} \rangle)$$

time constraints between two consecutive events

- The traditional episode rule  $\langle D, A \rangle \rightarrow \langle B \rangle$  in Fig.1 becomes two PERs:  $\langle D, A \rangle \xrightarrow{1} \langle B \rangle$  and  $\langle D, A \rangle \xrightarrow{2} \langle B \rangle$ .

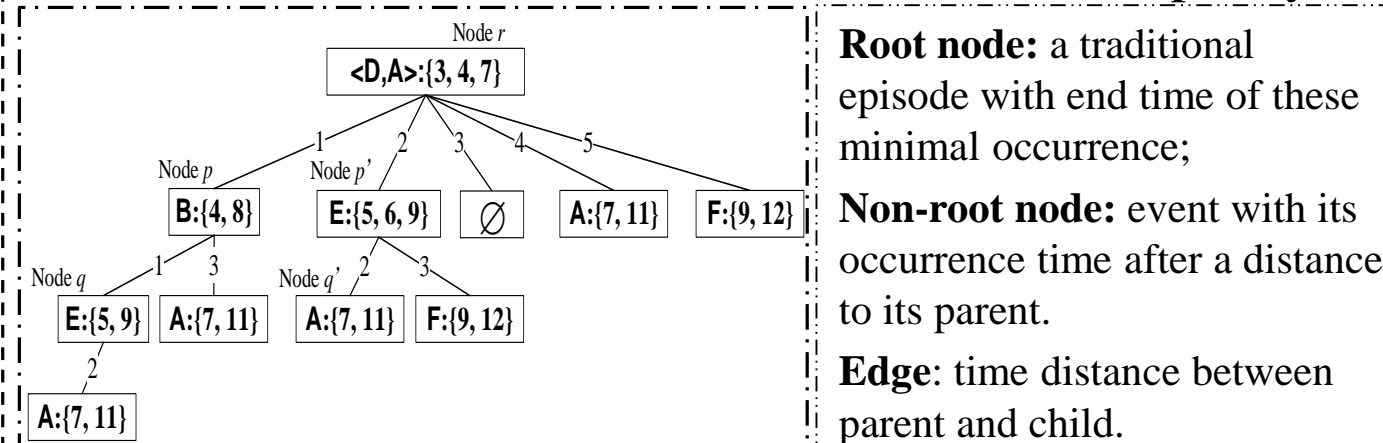
## Mining ALGORITHM & EFFICIENCY

### 1. MIP-ENUM Algorithm

The basic idea of MIP-ENUM is to enumerate PER candidates by concatenating discovered traditional episode with fixed-gap episode and subsequently filter the invalid ones according their confidence values.

### 2. MIP-TRIE Algorithm

**Data structure: PER-trie** stores valid PER compactly.



**Root node:** a traditional episode with end time of these minimal occurrence;  
**Non-root node:** event with its occurrence time after a distance to its parent.  
**Edge:** time distance between parent and child.

**Algorithm: MIP-TRIE(DFS) and MIP-TRIE(PRU).**

We use PER-trie to store all valid PER given an antecedent  $\alpha$  and propose two algorithms to build complete PER-trie.

- **MIP-TRIE(DFS)** expands the PER-trie by a recursively depth first search manner.
- **MIP-TRIE(PRU)** adopts an improved traverse strategy with pruning technique.

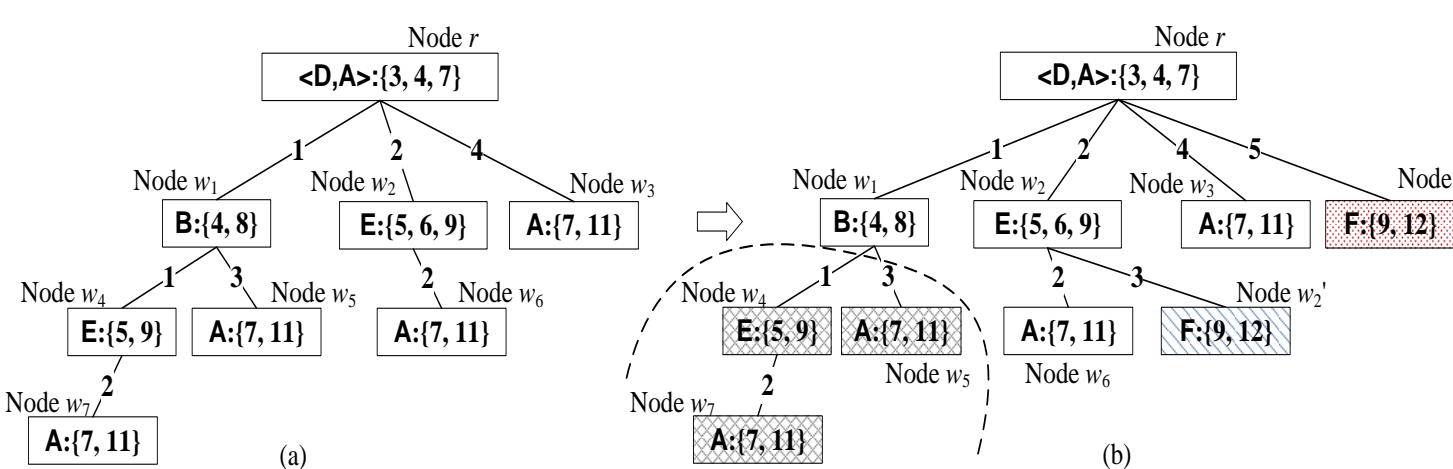
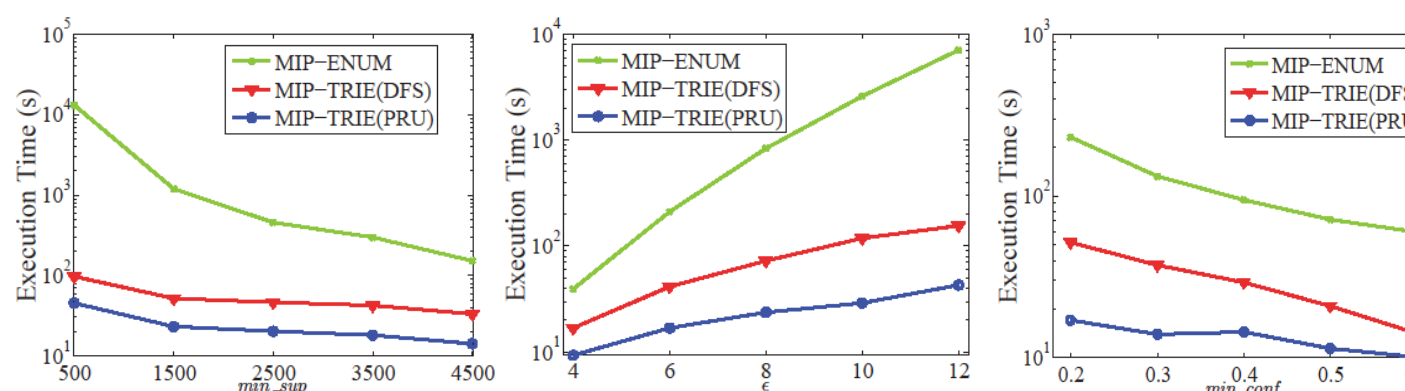


Fig.2 The expansion process of MIP-TRIE(PRU).

- $q'$  is expanded first as child of  $r$ , and we traverse  $w_1 - w_3$  and pruning  $w_4, w_5$  and  $w_7$  and finally traverse  $w_6$  with  $q'$ .

### Efficiency Comparisons



Dataset: Retail -- <http://fimi.cs.helsinki.fi/data/>

Observations: 1. MIP-TRIE(PRU) outperforms MIP-TRIE(DFS) and MIP-ENUM algorithm; 2. MIP-TRIE algorithms significantly outperform MIP-ENUM.

## EFFECTIVENESS of PER

**DATASET:** 150 related industry sector pairs of China stock market from Jan. 1, 2010 to Aug. 29, 2014.

**EVT SEQ. CONSTRUCTION:** **UP** (if the price increases) and **DN** (otherwise) for each industry sector.

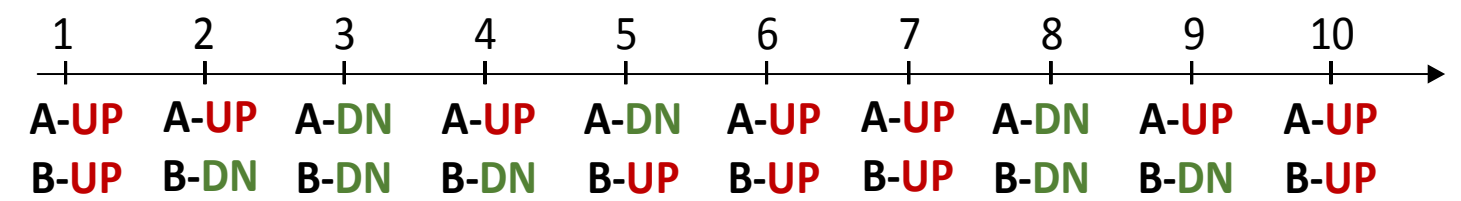


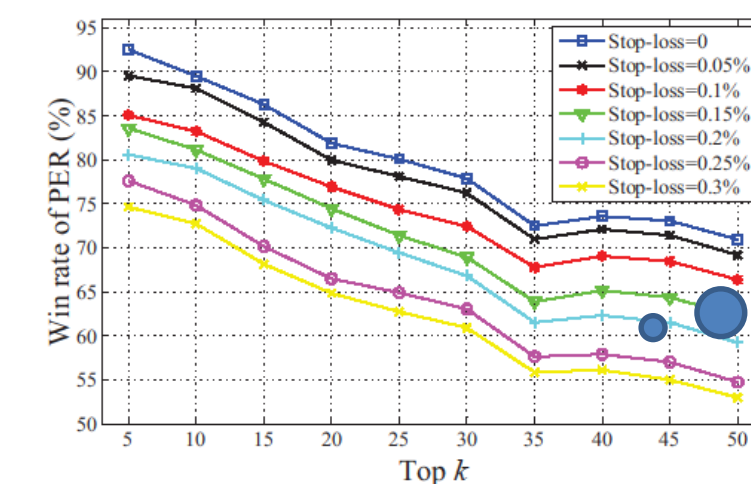
Fig.3 The example stock industry sector event sequence.

- A and B denote stock industry sectors.

**SETTINGS:** We use first 4-year sequence as the training set to mine PER on each sequence and degrade PER whose  $\Delta t = 5$  to traditional episode rule (denoted as TDR), then test prediction ability of them on the rest.

**COMPARISON:** For PER, we trade strictly according to the rule; for TDR, we trade after antecedent occurs and close out either consequent appears or the maximal occurrence window for consequent reaches.

**MEASURE:** We close out when the float loss exceeds a stop-loss threshold during the holdings by TDR. We compute the return of holdings and visualize the winning rate of PER under different stop-loss threshold.



The winning rate of PER > 50% indicates PER is more effective than TDR.

## VENUE & CONTACT INFORMATION

The 33<sup>rd</sup> IEEE International Conference on Data Engineering, San Diego, California, USA, April 19-22, 2017.

Email: {aoxiang, luop}@ict.ac.cn, jinwang@cs.ucla.edu, {zhuangfz, heq}@ics.ict.ac.cn

Homepage of MLDM Group, ICT, CAS: <http://mldm.ict.ac.cn>

Xiang's personal homepage: <http://mldm.ict.ac.cn/MLDM/~ao>

Public account on Wechat:

