

# Predicting Rare Events In Temporal Domains

Ricardo Vilalta and Sheng Ma

ICDM 2002

Presented by: Fahad Arshad

School of Electrical and Computer Engineering  
Purdue University



Slide 1/20

**PURDUE**  
UNIVERSITY

## Outline

- Introduction
- Event Prediction Problem
- Searching for Eventsets
- Building a Rule-base Model for Prediction
- Experiments and Results
- Conclusion



Slide 2/20

**PURDUE**  
UNIVERSITY

## Introduction

- Learning to predict rare events (*target events*) is a difficult problem
  - Attack on a network, fraudulent transactions at bank etc
- Prediction strategy
  1. Characterize target events by finding the types of events frequently preceding target events within fixed time window  $W$
  2. Validate that these event types uniquely characterize target events and do not occur often far from target events
  3. Combine validated event types using association rule mining to build a rule-based system for prediction



Slide 3/20



## Introduction

- Association Rule Mining
  - Let  $I = \{i_1, i_2, \dots, i_m\}$  be a set of  $m$  binary attributes called *items* and  $D = \{t_1, t_2, \dots, t_n\}$  be a set of transactions called the *database*.
  - A *rule* is defined as an implication of the form  $X \rightarrow Y$  where  $X, Y$  (*itemsets*) are subsets of  $I$  and  $X \cap Y = \text{empty}$
  - In the real world example, a *rule*:  $\{\text{milk, bread}\} \rightarrow \{\text{butter}\}$
  - $\text{Support}(X)$  is the proportion of transactions in the dataset which contain the itemset  $X$
  - Confidence of a rule is:

$$\text{conf}(X \Rightarrow Y) = \text{supp}(X \cup Y) / \text{supp}(X)$$

transaction ID	milk	bread	butter	beer
1	1	1	0	0
2	0	1	1	0
3	0	0	0	1
4	1	1	1	0
5	0	1	0	0



Slide 4/20



## Event Prediction Problem

- **Definition 1:**
  - A sequence of events is an ordered collection of events  $D = \langle d_1, d_2, \dots, d_n \rangle$ , where each event  $d_i$  is a pair  $d_i = (e_i, t_i)$ .  $e_i$  indicates the event type and  $t_i$  its occurrence time.
- $D_{target}$  (target events) a subset of  $D$  has size  $m$  and  $m \ll n$
- **Definition 2:**
  - An eventset  $Z$  is a set of event types  $\{e_i\}$ . Event set  $Z$  matches the set of events in window  $W$ , if every event type  $e_i \in Z$  is found in  $W$

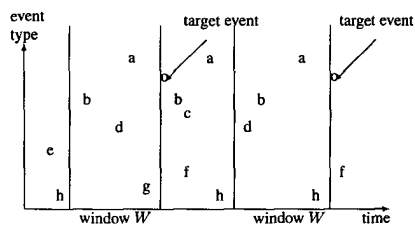


Slide 5/20



## Event Prediction Problem

- **Definition 3:**
  - Eventset  $Z$  has *support*  $s$  in  $D$  if  $s\%$  of all windows of size  $W$  preceding target events are matched by  $Z$ . Eventset  $Z$  is *frequent* if  $s$  is above a minimum user defined threshold
- **Definition 4:**
  - Eventset  $Z$  has *confidence*  $c$  in  $D$  if  $c\%$  of all time windows of size  $W$  matched by  $Z$  precede a target event. Eventset  $Z$  is *accurate* if  $c$  is above a minimum user defined threshold



Slide 6/20



## Searching for Eventsets

- **Frequent Eventsets**
  - Maintain in memory events within a sliding window of size  $W$
  - With each new target event, all event types within  $W$  are stored in a new transaction
  - A-priori algorithm is used to find all eventsets above a minimum user defined threshold.
  - Order of events and inter-arrival times between events within each  $W$  are not relevant.

### Algorithm 1: Finding Frequent Eventsets

**Input:** event sequence  $D$ , window size  $W$ , minimum support  $s\%$ , target-event type  $e^*$

**Output:** frequent eventsets  $\mathcal{F}$

FREQUENTEVENTSETS( $D, W, s, e^*$ )

```

(1)   $B = \emptyset; T = \emptyset$ 
(2)  foreach event  $d_i = (e_i, t_i) \in D$ 
(3)    currentTime =  $t_i$ 
(4)    foreach event  $d_j = (e_j, t_j) \in T$ 
(5)      if (currentTime -  $t_j$ ) >  $W$ 
(6)        Remove  $d_j$  from  $T$ 
(7)    end
(8)    if  $d_i$  is a target event (i.e.,  $e_i = e^*$ )
(9)       $B = B \cup \{e_j \mid (e_j, \cdot) \in T\}$ 
(10)    $T = T \cup d_i$ 
(11)  end
(12)  Use A-priori over  $B$  to find all frequent
(13)  eventsets with minimum support  $s\%$ .
(14)  Let  $\mathcal{F}$  be the set of all frequent eventsets.
(15)  return  $\mathcal{F}$ 
  
```



Slide 7/20



## Searching for Eventsets

- **Accurate Eventsets**
  - Find the number of times that each of the eventsets occur outside  $W$  preceding target events
  - Compute the confidence of each frequent eventset and eliminate those below the threshold
  - $B'$  a database of all eventsets not preceding target events is found.
  - If  $x_1$  and  $x_2$  are the number of transaction in  $B$  and  $B'$  respectively
  - confidence( $Z, B, B'$ ) =  $x_1 / (x_1 + x_2)$

### Algorithm 2: Finding Confident Eventsets

**Input:** event sequence  $D$ , minimum confidence  $c\%$ , time intervals  $I$ , frequent eventsets  $\mathcal{F}$ , database eventsets  $B$

**Output:** confident eventsets  $\mathcal{F}'$

CONFIDENTEVENTSETS( $D, c, I, \mathcal{F}, B$ )

```

(1)   $T = \emptyset; [a, b] = \text{next interval from } I$ 
(2)  foreach  $d_i = (e_i, t_i) \in D$ 
(3)    if  $t_i \in [a, b]$ 
(4)       $T = T \cup d_i$ 
(5)    if  $t_i > b$ 
(6)       $B' = B' \cup \{e_j \mid (e_j, \cdot) \in T\}$ 
(7)       $T = \emptyset; [a, b] = \text{next interval from } I$ 
(8)      Add  $d_i$  to  $D$ 
(9)  end
(10)  $\mathcal{F}' = \emptyset$ 
(11) foreach eventset  $Z$  in  $\mathcal{F}$ 
(12)   if confidence( $Z, B, B'$ ) >  $c$  AND
(13)    $P(Z|B) > P(Z|B')$ 
(14)      $\mathcal{F}' = \mathcal{F}' \cup Z$ 
(15)  end
(16)  return  $\mathcal{F}'$ 
  
```



Slide 8/20



## Searching for Eventsets

- **Validation step for Accurate Eventset:**

- Let  $P(Z|B)$  be the probability of  $Z$  occurring within database  $B$  and  $P(Z|B')$  the corresponding probability within  $B'$
- Event  $Z$  is validated if we reject the null hypothesis  $H_0$  with high confidence

$$H_0 : P(Z|B) \leq P(Z|B')$$

- If number of events is large, one assume a Gaussian distribution and reject the null hypothesis in favor of the alternate hypothesis  $H_1$

$$H_1 : P(Z|B) > P(Z|B')$$

- For a given confidence level  $\alpha$ , if the difference between the two probability is significant, we reject  $H_0$ . By choosing a small  $\alpha$  we can almost be certain about the relation of  $Z$  with target events.



Slide 9/20



## Building a Rule-base Model

- **Definition 5:**

- Eventset  $Z_i$  is said to be more specific than eventset  $Z_j$ , if  $Z_j$  is a subset of  $Z_i$ 
  - e.g.  $\{a, b, c\}$  is more specific than  $\{a, b\}$

- **Definition 6:**

- Eventset  $Z_i$  is said to have higher rank over eventset  $Z_j$ , represented as  $Z_i > Z_j$  if any of the following conditions is true
  1. The confidence of  $Z_i$  is greater than that of  $Z_j$
  2. The confidence of  $Z_i$  equals that of  $Z_j$ , but the support of  $Z_i$  is greater than that of  $Z_j$
  3. The confidence and support of  $Z_i$  equal that  $Z_j$ , but  $Z_i$  is more specific than  $Z_j$



Slide 10/20



## Building a Rule-base Model

- Find the most accurate and specific rule first
- Rule-based system  $R$  can be used for prediction by checking the occurrence of any of the events in  $R$  along the event sequence used for testing.
- The model predicts finding a target event within a time window of size  $W$  after any such eventset is detected.

**Algorithm 3:** Building Rule-Based Model

**Input:** eventsets  $\mathcal{F}'$

**Output:** Set of rules  $\mathcal{R}$

**RULE-BASED-EVENTSETS**( $\mathcal{F}'$ )

```
(1)  $\mathcal{R} = \emptyset$ 
(2) Sort  $\mathcal{F}'$  in decreasing order by rank
(3) while  $\mathcal{F}'$  is not empty
(4)   Let  $Z_i$  be the first eventset in  $\mathcal{F}'$ 
(5)   if  $Z_j \subset Z_i, i \neq j$ , remove  $Z_j$  from  $\mathcal{F}'$ 
(6)   Make a new rule  $r : Z_i \rightarrow \text{targetevent}$ 
(7)    $\mathcal{R} = \mathcal{R} \cup r$ 
(8)   Remove  $Z_i$  from  $\mathcal{F}'$ 
(9) end
(10) return  $\mathcal{R}$ 
```



Slide 11/20



## Experiments and Results

- **Artificial Data**
  - Data generator outputs uniformly distributed sequence of events over a fixed time interval
  - Time-interval= 1 week, number of event types = 50,  $W = 5$  min, target events = 50,  $s = 0.1$ , significance level for hypothesis testing  $\alpha = 0.01$ .
  - For each event sequence first 50% is for training and the rest for testing. Results are averaged over 30 runs.



Slide 12/20



## Experiments and Results

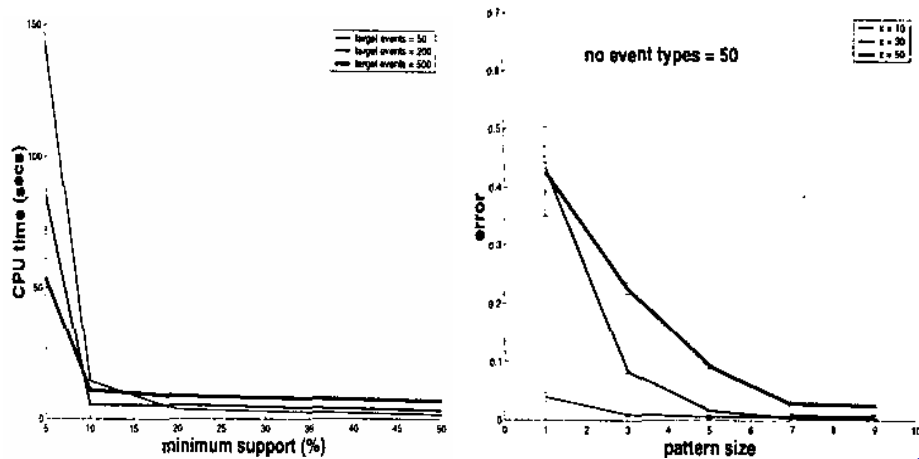
- Non-overlapping windows of size  $W$  that do not intersect the set of time windows preceding target events are negative examples.
- All time windows preceding target events are considered positive examples
- *Error* is the fraction of examples incorrectly classified by the rule-based system.



Slide 13/20



## Experiments and Results



Slide 14/20



## Experiments and Results

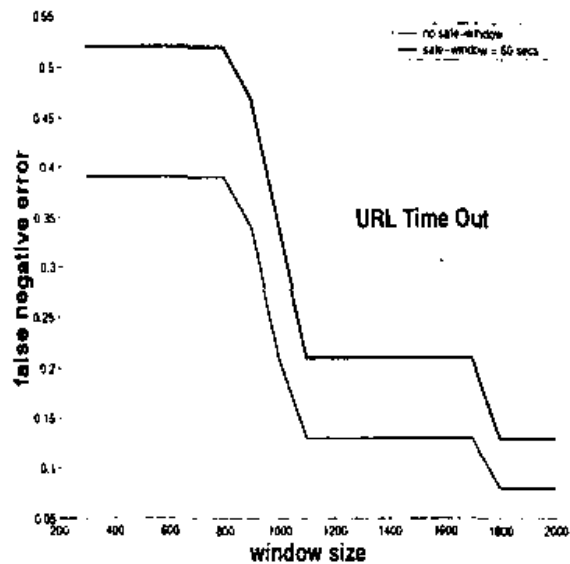
- **Real Production Data:**
  - Data collected over 1 month on a network with 750 hosts
  - 26000 events, 165 different types of events
- **Two types of target events considered**
  - EPP event indicates that end-to-end response time for a host is above a critical threshold
  - URL timeout indicates that a website is inaccessible.
- **Event is characterized by time, event type, and host**
  - Merge type and host to identify the nature of the event



Slide 15/20



## Experiments and Results

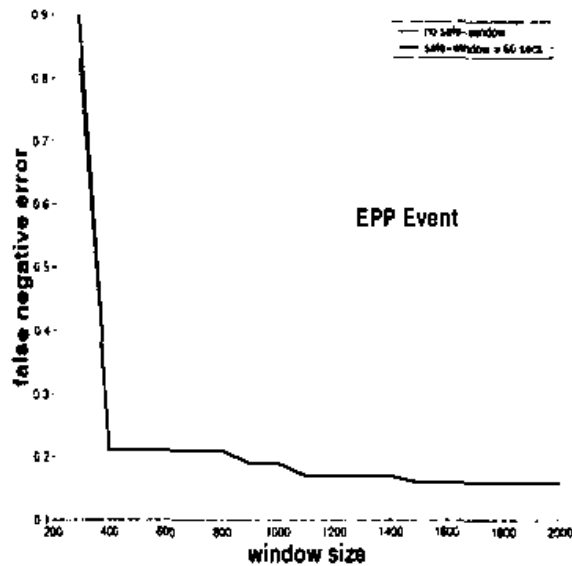


Slide 16/20





## Experiments and Results



Slide 17/20



## Conclusions

- An approach to detecting patterns in events sequences before a target event.
- Size of the time window preceding target events is crucial to this approach
- Two different combination of event-types and host of interest show how the false negative rate decreases significantly with increase in window size  $W$
- The approach is contingent on sets of events frequently occurring before target events.



Slide 18/20

