



Data Classification Using Sequential Pattern Mining

1

Under the Guidance of:-
Mrs. Dhanshree Tayade

¹SSBT's College of Engineering And Technology, Bambhori Jalgaon - 425001, Maharashtra, India



Outline of Topics

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- 1 Introduction
- 2 Literature Survey
- 3 Interesting Pattern
- 4 Methods
- 5 UML Diagram
- 6 Conclusion



Outline

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- 1 Introduction
- 2 Literature Survey
- 3 Interesting Pattern
- 4 Methods
- 5 UML Diagram
- 6 Conclusion



Introduction

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

Sequence classification system classifies dataset into sequential pattern. User inputs dataset which is to be classified. The classification involves the measure of interesting pattern in a class of sequence. This involves use of techniques such as support and cohesion. The interesting patterns are then used to build a sequence classifier.



Introduction

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

The pattern can both, itemset and subsequences. The support measures in how many sequence the pattern appears. Cohesion measures how close the items are to each other on an average. The classifier is built using the classification rules. There may be additional rules which occupy a lot of space. These rules also take time for processing. There is a need to prune such rules.



Introduction

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

This pruning is done by using a technique known as Lift. Lift is the ratio of observed support to the expected, if x and y were independent. A rule that has independent events are pointless to process. Such rules are eliminated using Lift. The pruned subset classifies the interesting patterns of the dataset into a sequential pattern.



Outline

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- 1 Introduction
- 2 Literature Survey**
- 3 Interesting Pattern
- 4 Methods
- 5 UML Diagram
- 6 Conclusion



Literature Survey

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

Existing system classifies data using interesting patterns. The patterns are connects into rules to further refine the rules. The system uses confidence Method. The confidence value of a rule $X = Y$ with respect to a set of transactions T , is the proportion of the transaction that contains X as well as Y .



Literature Survey

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

The proposed system classifies data using interesting patterns. Later, rules are drawn from these interesting patterns. By using confidence, we may find a rule to be true for a particular instances, but the same rule won't be true for all instances. Hence, Lift method is used in the proposed system. 'Lift' prunes only those rules in which positive dependency are stored in the databases. These rules can be used for classification the antecedent and consequent are independent of each other. The rules which possess dependency are stored in the database. These rules can be used for classification depending on the different instances.



System Architecture

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

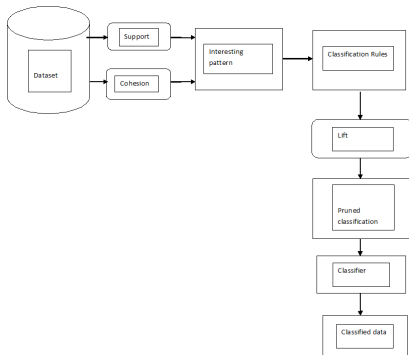


Figure: System Architecture



Outline

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- 1 Introduction
- 2 Literature Survey
- 3 Interesting Pattern**
- 4 Methods
- 5 UML Diagram
- 6 Conclusion



Interesting Pattern

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- Interesting Patterns in sequences, a pattern is typically evaluated based on how often it occurs (support) along with the proximity of the items (cohesion).
- Utilise interesting patterns to build classifiers.
- Use cohesion and support methods to define interesting patterns in a sequence dataset.

$$I(P) = f(P) * c(P) \quad (1)$$



Outline

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- 1 Introduction
- 2 Literature Survey
- 3 Interesting Pattern
- 4 Methods**
- 5 UML Diagram
- 6 Conclusion



Support

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- Support is an indication of how frequently the item-set appears in the database.
- The support count of a pattern is defined as the number of different sequences in which the pattern occurs regardless of how many times the pattern occurs in any single sequence.

$$F(P) = N(P)/S \quad (2)$$

where N =set of sequences containing all items of X .



Cohesion

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- Cohesion measures how close the items making up the pattern are to each other on average, using the lengths of the shortest intervals containing the pattern in different sequences.
- The cohesion of P in a single sequence s is defined as

$$C(P, s) = |P|/W(P, s) \quad (3)$$



'Lift' Method

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- Lift interestingness measure defines the number of transaction that contain the item used to find interesting patterns.
- The Lift measure is denote by $Lift(X = Y)$ as shown in

$$Lift(X \rightarrow Y) = \frac{supp(X + Y)}{supp(X) * supp(Y)} \quad (4)$$



Example

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

Example database with 5 transactions and 5 items

ID	Milk	Bread	Butter	Books	Pens
1	1	1	0	0	0
2	0	0	1	0	0
3	0	0	0	1	1
4	1	1	1	0	0
5	0	1	0	0	0



- $Support = F(P) = N(P)/S$

$$F(P) = 1/5$$

$$F(P) = 0.2$$

- $Lift(X \implies Y) = supp(X \cup Y) / supp(X) * supp(Y)$
 $Lift(Milk, Bread \implies Butter) = 0.2 / 0.4 * 0.4 = 1.25$



Outline

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- 1 Introduction
- 2 Literature Survey
- 3 Interesting Pattern
- 4 Methods
- 5 UML Diagram**
- 6 Conclusion



UseCase Diagram

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

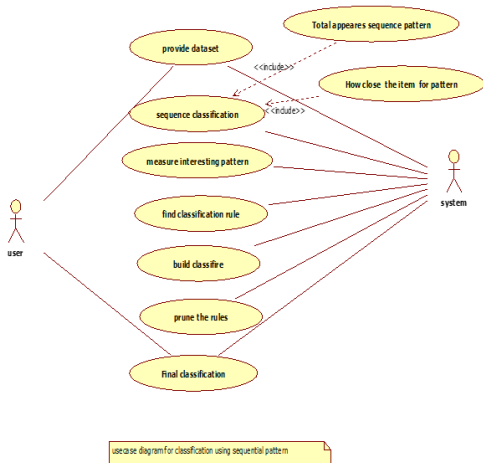
Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion





Outline

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

- 1 Introduction
- 2 Literature Survey
- 3 Interesting Pattern
- 4 Methods
- 5 UML Diagram
- 6 Conclusion**



Conclusion

Data
Classification
Using
Sequential
Pattern
Mining

Outline

Introduction

Literature
Survey

Interesting
Pattern

Methods

UML Diagram

Conclusion

The system classifies data using rules which are drawn from Interesting Pattern. These rules are pruned using Lift instead of Confidence. The proposed algorithm depends mainly on support compared to the algorithm proposed in the paper which uses both support and confidence. This results in simplification of the process as well as the rules that possesses dependency are not pruned. Thus, making the system more versatile.