

Pattern Structures

Pattern Structures

- Models describe whole or a large part of the data
- Pattern characterizes some local aspect of the data
- Pattern is a predicate that returns “true” for those objects or parts of objects in the data for which the pattern occurs and “false” otherwise

Pattern Specification

- To specify a pattern, need to specify
 - Syntax of the patterns (language specifying how they are defined)
 - Semantics of the patterns (interpretation of what they tell us about the data)
- Patterns can be considered in two different types of discrete-valued data
 1. Data in standard matrix form
 2. Data described as strings

Patterns in Data Matrices

- Start from primitive patterns and combine using logical connectives
- Data Matrix notation:
 - p variables X_1, \dots, X_p
 - $\mathbf{x} = \{x_1, \dots, x_p\}$ is a p -dimensional vector of measurements

Primitive Patterns

- Subset of all possible observations over variables X_1, \dots, X_p
- If c is a possible value of X_k then $X_k = c$ is a primitive pattern
- If values of X_k are ordered then $X_k \leq c$ is a primitive condition
- Multivariate conditions: $X_k X_j > 2$
 $X_k = X_j$

Complex Patterns

- Given a set of primitive patterns we can form more complex patterns by using logical connectives such as AND and OR
- Example: $(\text{age} < 40) \wedge (\text{income} < 10)$
- $(\text{chips} = 1) \wedge (\text{beer} = 1) \vee (\text{soft-drink} = 1)$
is a subset of a market-basket database

Pattern Class

- Pattern class is a set of legal patterns
- Defined by specifying a collection of primitive patterns and the legal ways of combining primitive patterns
- Example: If variables X_1, \dots, X_p all range over $\{0, 1\}$ we can define a class of patterns C consisting of all possible conjunctions of the form $(X_{j_1}=1) \wedge (X_{j_2}=1) \wedge \dots \wedge (X_{j_k}=1)$
- Conjunctive patterns such as frequent sets are relatively easy to discover

Frequency of a Pattern Class

- Given a Pattern class and a Data Set D , an important property of a pattern is its frequency
- Frequency $fr(\rho)$ of a pattern ρ is defined as
The relative number of observations in the dataset about which ρ is true

Importance of Frequency of a Pattern

- Patterns that occur reasonably often are of interest in data mining
- Frequency of a pattern close to 0 can also be informative
 - Rare and unusual phenomenon
- Other properties of relevance:
 - Semantic simplicity, understandability, novelty and surprise
- Example of uninteresting pattern
 - Disjunction of all conjunctive patterns in the data set forms a pattern of frequency 1 – which is uninteresting

Pattern Discovery Task

- Find all patterns from that class that satisfy certain conditions with respect to the data sets
- Example: Find all the frequent set patterns whose frequency is at least 0.1 and where variable X_7 occurs in the pattern
- Might include conditions on the informativeness, novelty and understandability of the pattern
- Challenge is to find the right balance between
 - expressivity of the patterns,
 - comprehensibility and
 - computational complexity of solving the discovery task

Rule

- A rule is an expression of the form $\rho \rightarrow \phi$

- **Accuracy** of the rule

$$p(\phi | \rho) = \frac{fr(\rho \wedge \phi)}{fr(\rho)}$$

- **Support** of the rule

$fr(\rho \rightarrow \phi)$ of the rule $\rho \rightarrow \phi$

is defined either as

$fr(\rho)$: fraction of objects to which the rule applies

Or

$fr(\rho \wedge \phi)$: fraction of objects for which both the left hand and right hand sides apply

Association Rule

- A rule would have the form

$$\{A_1, \dots, A_k\} \rightarrow \{B_1, \dots, B_h\}$$

where each of the A_k s and B_j s are binary variables

- Which when written out in full has the form

$$(A_1 = 1) \wedge \dots \wedge (A_k = 1) \rightarrow (B_1 = 1) \wedge \dots \wedge (B_h = 1)$$

Functional Dependency

- Previously each pattern referred to a single observation
- Patterns can be defined by referring to several variables
- Example: identify all points in a geographical database that form the vertices in an equilateral triangle

Formal Functional Dependency

- Expression of the form

$$A_{i_1}A_{i_2}\dots A_{i_k} \rightarrow A_{i_{k+1}}$$

where $1 \leq i_j \leq p$ for $i = 1, \dots, k+1$

- A dataset has this property if for all pairs of observations x and y in the dataset, if x and y agree on all the variables A_{i_j} for $j = 1, \dots, k$ then x and y agree also on $A_{i_{k+1}}$

Patterns that Specify a Set of Records

- Previous specifications of patterns refer to only a single record in the database
- Describing patterns that refer to several records, e.g.,

$$\{x_k \mid \text{age} \leq 40 \wedge \text{income} \leq 10\}$$

Criteria for Interestingness

- Given a rule $\rho \rightarrow \phi$, its interestingness can be defined in many ways
- Background knowledge about variables referred to in the patterns ρ and ϕ have an influence on the interestingness of the rule
- Examples:
 - In credit scoring data set decide beforehand that rules connecting month of birth and credit score are not interesting
 - In market-basket case, interest in a rule is directly proportional to the frequency of the rules multiplied by the prices of the items mentioned, i.e., more interested in rules of high frequency that connect expensive items

Statistical Criteria for Interestingness

- Purely statistical criteria are easier to use in an application-independent way
- Construct a 2 x 2 contingency table using presence or absence of ρ and ϕ as the variables and having as the counts the frequencies of the four different combinations

	ϕ	$\sim\phi$
ρ	$fr(\rho \wedge \phi)$	$fr(\rho \wedge \sim\phi)$
$\sim\rho$	$fr(\sim\rho \wedge \phi)$	$fr(\sim\rho \wedge \sim\phi)$

Cross-Entropy Measure of Interestingness

	ϕ	$\sim\phi$
ρ	$fr(\rho \wedge \phi)$	$fr(\rho \wedge \sim\phi)$
$\sim\rho$	$fr(\sim\rho \wedge \phi)$	$fr(\sim\rho \wedge \sim\phi)$

Cross entropy between the binary variable ϕ with and without conditioning on the event ρ

$$J(\rho \rightarrow \phi) = p(\rho) \left(p(\phi | \rho) \log \frac{p(\phi | \rho)}{p(\phi)} + (1 - p(\phi | \rho)) \log \frac{1 - p(\phi | \rho)}{1 - p(\phi)} \right)$$

Empirically observed accuracy of the rule

How widely is the rule applicable?

Empirically observed marginal probabilities

How dissimilar is our knowledge about ϕ is from only knowing about marginal $p(\phi)$ compared with knowing that ρ holds

Patterns for Strings

- Different Types of Patterns are required for data in the form of strings
- String over an alphabet S is a sequence a_1, \dots, a_n of elements (letters) of S
- Examples of alphabets:
 - Binary $\{0,1\}$
 - Set of ASCII codes
 - DNA alphabet $\{A,C,G,T\}$
 - Set of all words consisting of ASCII characters
- Set of all strings built from letters from S is denoted by S^*

String Data

- No fixed set of variables
- For notions of probability we consider each of the letters of the string to be a random variable
- Interested in finding how many times a certain pattern occurs in strings
- Example: no of exact occurrences of a certain DNA sequence in a large collection of sequences
- Simplest string pattern is a substring: the pattern $b_1 \dots b_k$ occurs in the string $a_1 \dots a_n$ at position i
- Examples:
 - For DNA subsequences we need to find occurrences of ATTATTAA
 - For strings over ASCII alphabet whether the pattern “data mining” occurs

Specifying a larger Class of Patterns: Regular Expressions

- Regular Expression E defines a set $L(E)$ of strings
- Expression E is one of:
 - A string s ; then $L(s) = \{s\}$
 - A concatenation E_1E_2 ; the set $L(E_1E_2)$ consists of all strings that are a concatenation of a string in $L(E_1)$ and a string in $L(E_2)$
 - A choice $E_1|E_2$; then $L(E_1|E_2) = L(E_1) \cup L(E_2)$
 - An iteration E^* ; then $L(E^*)$ that can be written as a concatenation of 0 or more strings from $L(E)$
- $10(00|11)^*01$ is a regular expression that describes all strings that start with 10 and end with 10 and inbetween contain a sequence of pairs 00 and 11
- Many complicated phenomena can be captured, but not balanced sequences of parentheses

Episodes

- Regular Expressions are not sufficiently expressive for expressing variations in the occurrence times of events
- Episodes can do this
- Partially ordered collection of events occurring together
 - Events may be of different types and may refer to different variables
- Example from biostatistics: event is a headache followed by a sense of disorientation occurring within a given period of time
- Be insensitive to intervening events, e.g., alarms in telecom network, logs of user interface actions