Data Mining Part 1. Introduction

1.4 Input

Spring 2010

Instructor: Dr. Masoud Yaghini

Outline

- Instances
- Attributes
- References

Instances

Instances

• Instance:

- Individual, independent example of the concept to be learned.
- Characterized by a predetermined set of attributes
- Input to learning process: set of instances/dataset
- Each dataset is represented as a matrix of instances versus attributes
 - Represented as a single table or flat file
- Rather restricted form of input
 - No relationships between objects
 - Problems often involve relationships between objects rather than separate, independent instances.

An example: A family tree

• Example:

- a family tree is given, and we want to learn the concept sister.
- This tree is the input to the learning process, along with a list of pairs of people and an indication of whether they are sisters or not.



Two ways of expressing the sister-of relation

| first person | second person | sister of? | |
|-----------------|------------------|---------------|--|
| Peter | Peggy | no | |
| Peter | Steven | no | |
| | | | |
| Steven | Peter | no | |
| Steven | Graham | no | |
| Steven | Pam | yes | |
| Steven | Grace | no | |
| | | | |
| lan | Pippa | yes | |
| | | | |
| Anna | Nikki | yes | |
| | | | |
| Nikki | Anna | yes | |
| | | | |

| person Pam | of? | |
|---------------|---------------------------------|--|
| Pam | | |
| | yes | |
| Pam | yes | |
| Pippa | yes | |
| Pippa | yes | |
| Nikki | yes | |
| Anna | yes | |
| All the rest | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | Pippa Pippa Nikki Anna | |

• Neither table is of any use without the family tree itself.

Family tree represented as a table

| Name | Gender | Parent1 | Parent2 | |
|--------|--------|---------|---------|--|
| Peter | male | ? | ? | |
| Peggy | female | ? | ? | |
| Steven | male | Peter | Peggy | |
| Graham | male | Peter | Peggy | |
| Pam | female | Peter | Peggy | |
| lan | male | Grace | Ray | |
| | | | | |

 These tables do not contain independent sets of instances because values in the Name, Parent1, and Parent2 columns refer to rows of the family tree relation.

The sister-of relation represented in a table

| First person | | | Second person | | | | | |
|--------------|--------|---------|---------------|-------|--------|---------|---------|------------|
| Name | Gender | Parent1 | Parent2 | Name | Gender | Parent1 | Parent2 | Sister of? |
| Steven | male | Peter | Peggy | Pam | female | Peter | Peggy | yes |
| Graham | male | Peter | Peggy | Pam | female | Peter | Peggy | yes |
| lan | male | Grace | Ray | Pippa | female | Grace | Ray | yes |
| Brian | male | Grace | Ray | Pippa | female | Grace | Ray | yes |
| Anna | female | Pam | lan | Nikki | female | Pam | lan | yes |
| Nikki | female | Pam | lan | Anna | female | Pam | lan | yes |
| | | | all the | rest | | | | no |

• Each of instance is an individual, independent example of the concept that is to be learned.

A simple rule for the sister-of relation

- A simple rule for the sister-of relation is as follows:
 - If second person's gender = female
 and first person's parent1 = second person's parent1
 then sister-of = yes

Denormalization

Denormalization or flattening:

- Several relations are joined together to make one
- to recast data into a set of independent instances
- Possible with any finite set of finite relations
- Problem:
 - Denormalization may produce false regularities that reflect structure of database
 - Example: "supplier" predicts "supplier address"

Instances

- The input to a data mining scheme is generally expressed as a table of independent instances of the concept to be learned.
- The instances are the rows of the tables the attributes are the columns.

Attributes

Attributes

 Each instance is described by a fixed predefined set of features or attributes

- Problem: Number of attributes may vary in different instances
 - Example: the instances were transportation vehicles
 - Possible solution: to make each possible feature an attribute and to use a special flag value to indicate that a particular attribute is not available for a particular case.

Attributes

- Another problem: existence of an attribute may depend of value of another one
 - Spouse's name depends on the value of married or single attribute

Attributes Types

• Possible attribute types (levels of measurement):

- Numeric (Interval-scaled) quantities
- Binary (Boolean) quantities
- Categorical (Nominal) quantities
- Ordinal quantities

Numeric Quantities

- Numeric (Interval-scaled) quantities are continuous measurements of a roughly linear scale.
- Numeric quantities are not only ordered but measured in fixed and equal units
- Example 1: attribute "temperature" expressed in degrees Fahrenheit
- **Example 2**: attribute "date" (year)
- Difference of two values makes sense

Binary Quantities

- A **Binary (Boolean) quantity** can take has only two states: 0 or 1, where 0 means that the variable is absent, and 1 means that it is present.
- Given the variable **smoker** describing a patient,
 - 1 indicates that the patient smokes
 - 0 indicates that the patient does not.

Categorical Quantities

- Categorical (Nominal) attributes take on values in a prespecified, finite set of possibilities.
- Categorical quantities values are distinct symbols
 - Values themselves serve only as labels or names
- Example: attribute "outlook" from weather data
 - Values: "sunny", "overcast", and "rainy"
- No relation is implied among categorical values (no ordering or distance measure)
- Special case: binary attribute
 - Example: true/false or yes / no

Categorical Quantities

- Note: addition, subtraction, and comparing don't make sense
- Only equality tests can be performed
 - Example:

| outlook: | sunny | \rightarrow | no |
|----------|----------|---------------|-----|
| | overcast | \rightarrow | yes |
| | rainy | \rightarrow | yes |

Ordinal quantities

- Ordinal quantities are ones that make it possible to rank order the categories.
- But: no distance between values defined
- Example: attribute "temperature" in weather data
 - Or: "hot" > "mild" > "cool"
- Note: it makes sense to compare two values, but addition and subtraction don't make sense
- Example rule:
 - temperature < hot => play = yes
- Distinction between nominal and ordinal not always clear (e.g. attribute "outlook")

References

References

 Ian H. Witten and Eibe Frank, Data Mining: Practical Machine Learning Tools and Techniques, 2nd Edition, Elsevier Inc., 2005. (Chapter 2)

The end