

Nature of Data Sets

- Data Set: set of measurements from an environment or process
- Simple case: n objects with p measurements each: $n \times p$ matrix
- p columns are called *variables*, *features*, *attributes* or *fields*

US Census Bureau: Public Use Microdata

Sample data sets (PUMS)

| ID | Age Quantitative Continuous | Sex Categorical Nominal | Marital Status | Education Categorical Ordinal | Income |
|-----|--------------------------------|----------------------------|-------------------|----------------------------------|----------------------------------|
| 248 | 54 | Male | Married | High School grad | 100000 Noisy data A guess? |
| 249 | ?? Missing data | Female | Married | HS grad | 12000 |
| 250 | 29 | Male | Married | Some College | 23000 |
| 251 | 9 | Male | Not Married | Child | 0 |
| | | | | | |

PUMS: Public use data that has identifying information removed.
Available in 5% and 1% sample sizes. 1% sample has 2.7 million records.

Medical Patient Data

- BP at different times of day
- Image data (x-ray or MRI)
- Specialist's comments (text)
- Hierarchy of relationships between patients, doctors, hospitals

$N \times p$ data matrix is often an oversimplification of what occurs in practice

Data Matrix Names

- Data set
- Training data
- Sample
- database

Mining Text Documents

- **Tasks:** Text categorization, Clustering Similar Documents, Finding documents that match a query

A document is viewed as a sequence of words and punctuations

Reuters collection is at <http://www.research.att.com/~lewis>

- Can be represented as matrix (d,w) where document d and word w has value 1 or 0 (sparse matrix)

Transaction Data

List of store purchases: date, customer ID, list of items and prices

Web transaction log -sequence of triples: (user id, web page, time)

Can be transformed
into
binary-valued
matrix

Individuals

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | | 1 | 1 | 1 | | | | 1 |
| 1 | | | | | 1 | | | | |
| 1 | 1 | 1 | | | 1 | 1 | 1 | 1 | |
| 1 | | | | | | | 1 | | |
| 1 | | 1 | | 1 | | 1 | | | 1 |
| 1 | | | 1 | | | | | | 1 |

Web Page Visited

Representations in Data Mining: Models and Patterns

- *Global Model*
 - *Make a statement about any point in p-space*
 - *E.g., assign a point to a cluster*
 - *Even when some values are missing*
 - *Simple model: $Y = aX + c$*
 - *Functional model is linear*
 - *Linear in variables rather than parameters*

Representations in Data Mining: Models and Patterns

- Local Patterns
 - *Make a statement about restricted regions of space spanned by variables*
 - *E.g. 1: if $X > \text{thresh1}$ then $\text{Prob}(Y > \text{thresh2}) = p$*
 - *E.g. 2: certain classes of transactions do not show peaks and troughs (bank discovers dead people's open accounts)*

Data Mining Tasks

(different objectives for mining the data)

1. Exploratory Data Analysis
Visualizing Data (histograms)
2. Descriptive Modeling
Density Estimation, Cluster Analysis
3. Predictive Modeling
 - a) Classification (response variable is a class)
 - b) Regression (response variable is numeric)
4. Discovering Patterns and Rules
Association Rules
5. Retrieval by Content
Text – based (Google) and Image-based (QBIC)

Data Mining Algorithms

(Four basic components in each algorithm*)

1. Model or Pattern Structure

Determining underlying structure or functional form we seek from data

2. Score Function

Judging the quality of the fitted model

3. Optimization and Search Method

Searching over different model and pattern structures

4. Data Management Strategy

Handling data access efficiently

***Illustrated earlier in Regression example**