



UNIT I

Data Reduction

GALGOTIAS
UNIVERSITY

Data Reduction Strategies

Data reduction: Obtain a reduced representation of the data set that is much smaller in volume but yet produces the same (or almost the same) analytical results

Why data reduction? — A database/data warehouse may store terabytes of data. Complex data analysis may take a very long time to run on the complete data set.

Data reduction strategies

- Dimensionality reduction, e.g., remove unimportant attributes

 - Wavelet transforms

 - Principal Components Analysis (PCA)

 - Feature subset selection, feature creation

- Numerosity reduction (some simply call it: Data Reduction)

 - Regression and Log-Linear Models

 - Histograms, clustering, sampling

 - Data cube aggregation

- Data compression

 -

 -

Data Reduction 1: Dimensionality Reduction

- **Curse of dimensionality**
 - When dimensionality increases, data becomes increasingly sparse
 - Density and distance between points, which is critical to clustering, outlier analysis, becomes less meaningful
 - The possible combinations of subspaces will grow exponentially
- **Dimensionality reduction**
 - Avoid the curse of dimensionality
 - Help eliminate irrelevant features and reduce noise
 - Reduce time and space required in data mining
 - Allow easier visualization
- **Dimensionality reduction techniques**
 - Wavelet transforms
 - Principal Component Analysis
 - Supervised and nonlinear techniques (e.g., feature selection)

Data Reduction 2: Numerosity Reduction

- Reduce data volume by choosing alternative, *smaller forms* of data representation
- **Parametric methods** (e.g., regression)
 - Assume the data fits some model, estimate model parameters, store only the parameters, and discard the data (except possible outliers)
 - Ex.: Log-linear models—obtain value at a point in m -D space as the product on appropriate marginal subspaces
- **Non-parametric methods**
 - Do not assume models
 - Major families: histograms, clustering, sampling, ...

Parametric Data Reduction: Regression and Log-Linear Models

- **Linear regression**
 - Data modeled to fit a straight line
 - Often uses the least-square method to fit the line
- **Multiple regression**
 - Allows a response variable Y to be modeled as a linear function of multidimensional feature vector
- **Log-linear model**
 - Approximates discrete multidimensional probability distributions

- **Linear regression**

- Data modeled to fit a straight line
- Often uses the least-square method to fit the line

- **Multiple regression**

- Allows a response variable Y to be modeled as a linear function of multidimensional feature vector

- **Log-linear model**

- Approximates discrete multidimensional probability distributions

GALGOTIAS
UNIVERSITY

Linear regression: $Y = w X + b$

Two regression coefficients, w and b , specify the line and are to be estimated by using the data at hand

Using the least squares criterion to the known values of $Y_1, Y_2, \dots, X_1, X_2, \dots$

Multiple regression: $Y = b_0 + b_1 X_1 + b_2 X_2$

Many nonlinear functions can be transformed into the above

Log-linear models:

Approximate discrete multidimensional probability distributions

Estimate the probability of each point (tuple) in a multi-dimensional space for a set of discretized attributes, based on a smaller subset of dimensional combinations

Useful for dimensionality reduction and data smoothing

DSDM framework

References Jiawei Han, Micheline Kamber and Jian Pei Data Mining: Concepts and Techniques, 3rd ed. The Morgan Kaufmann Series in Data Management Systems Morgan Kaufmann Publishers, July 2011. ISBN 978-0123814791

The logo of Galgotias University is a circular emblem with a stylized 'G' shape in the center. The 'G' is composed of several overlapping, curved segments in shades of yellow, orange, and blue. The background of the emblem is a light, swirling pattern.

GALGOTIAS
UNIVERSITY



Thank You