

The logo of Galgotias University is a stylized 'G' composed of three curved, overlapping bands in shades of yellow, blue, and red. It is centered in the background of the slide.

UNIT I INTRODUCTION

GALGOTIAS
UNIVERSITY

Data Integration

Data integration:

Combines data from multiple sources into a coherent store

Schema integration: e.g., A.cust-id \equiv B.cust-#

Integrate metadata from different sources

Entity identification problem:

Identify real world entities from multiple data sources, e.g., Bill Clinton = William Clinton

Detecting and resolving data value conflicts

For the same real world entity, attribute values from different sources are different

Possible reasons: different representations, different scales, e.g., metric vs. British units



Handling Redundancy in Data Integration

- Redundant data occur often when integration of multiple databases
 - *Object identification*: The same attribute or object may have different names in different databases
 - *Derivable data*: One attribute may be a “derived” attribute in another table, e.g., annual revenue
- Redundant attributes may be able to be detected by *correlation analysis* and *covariance analysis*
- Careful integration of the data from multiple sources may help reduce/avoid redundancies and inconsistencies and improve mining speed and quality

Correlation Analysis (Nominal Data)

- **X^2 (chi-square) test**
- The larger the X^2 value, the more likely the variables are related
- The cells that contribute the most to the X^2 value are those whose actual count is very different from the expected count
- Correlation does not imply causality
 - # of hospitals and # of car-theft in a city are correlated
 - Both are causally linked to the third variable: population

Chi-Square Calculation: An Example

	Play chess	Not play chess	Sum (row)
Like science fiction	250(90)	200(360)	450
Not like science fiction	50(210)	1000(840)	1050
Sum(col.)	300	1200	1500

- χ^2 (chi-square) calculation (numbers in parenthesis are expected counts calculated based on the data distribution in the two categories)

$$\chi^2 = \frac{(250 - 90)^2}{90} + \frac{(50 - 210)^2}{210} + \frac{(200 - 360)^2}{360} + \frac{(1000 - 840)^2}{840} = 507.93$$

- It shows that like_science_fiction and play_chess are correlated in the group

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 stylized, circular emblem. It features a central white swirl that transitions into a yellow and orange arc at the top, and a blue and white arc at the bottom. The overall shape is reminiscent of a 'G' or a globe.

GALGOTIAS
UNIVERSITY



Thank You