# Data Preprocessing

Universitas Airlangga, 18 September 2018

# Why Preprocess the Data?

*Data in real world is dirty*

- Incomplete

  *lacking attribute values or certain attributes of interest, or containing only aggregate data.* e.g., sales=" "

- Inaccurate or Noisy

  *containing errors, or values that deviate from the expected.*

  *e.g., age="-10"*

- Inconsistent

  *containing discrepancies in the codes or names.*

  *e.g., Was rating "1,2,3", now rating "A, B, C"*

- Timeliness

# Why Is Data Preprocessing Important?

- No quality data, no quality mining results!
  - Quality decisions must be based on quality data
    - e.g., duplicate or missing data may cause incorrect or even misleading statistics.
  - Data warehouse needs consistent integration of quality data
- Data extraction, cleaning, and transformation comprises the majority of the work of building a data warehouse. —Bill Inmon

# Major Tasks in Data Preprocessing

- Data cleaning
  - Fill in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies

- Data integration
  - Integration of multiple databases, data cubes, or files

- Data transformation
  - Normalization and aggregation

- Data reduction
  - Obtains reduced representation in volume but produces the same or similar analytical results

- Data discretization
  - Part of data reduction but with particular importance, especially for numerical data

# Data Cleaning

# Data Cleaning

- Importance
  - "Data cleaning is one of the three biggest problems in data warehousing"—Ralph Kimball
  - "Data cleaning is the number one problem in data warehousing"—DCI survey

- Data cleaning tasks

  - Fill in missing values

  - Identify outliers and smooth out noisy data

  - Correct inconsistent data

  - Resolve redundancy caused by data integration

# Missing Data

- Data is not always available
  - E.g., many tuples have no recorded value for several attributes, such as customer income in sales data

- Missing data may be due to
  - equipment malfunction
  - inconsistent with other recorded data and thus deleted
  - data not entered due to misunderstanding
  - certain data may not be considered important at the time of entry
  - not register history or changes of the data

- Missing data may need to be inferred.

# How to Handle Missing Data?

- Ignore the tuple: usually done when class label is missing (assuming the tasks in classification—not effective when the percentage of missing values per attribute varies considerably.

- Fill in the missing value manually: tedious + infeasible?

- Fill in it automatically with
  - a global constant : e.g., "unknown", a new class?!
  - the attribute mean
  - the attribute mean for all samples belonging to the same class: smarter
  - the most probable value: inference-based such as Bayesian formula or decision tree

# Noisy Data

- Noise: random error or variance in a measured variable
- Incorrect attribute values may due to
  - faulty data collection instruments
  - data entry problems
  - data transmission problems
  - technology limitation
  - inconsistency in naming convention
- Other data problems which requires data cleaning
  - duplicate records
  - incomplete data
  - inconsistent data

# How to Handle Noisy Data?

- Binning method:
  - first sort data and partition into (equi-depth) bins
  - then one can smooth by bin means,  smooth by bin median, smooth by bin boundaries, etc.
- Clustering
  - detect and remove outliers
- Combined computer and human inspection
  - detect suspicious values and check by human (e.g., deal with possible outliers)
- Regression
  - smooth by fitting the data into regression functions

# Data Transformation

- Smoothing: remove noise from data

- Aggregation: summarization

- Generalization: concept hierarchy climbing

- Normalization: scaled to fall within a small, specified range
  - min-max normalization
  - z-score normalization
  - normalization by decimal scaling

- Attribute/feature construction
  - New attributes constructed from the given ones

# Data Transformation: Normalization

- min-max normalization

$$v' = \frac{v - min_A}{max_A - min_A}(new\_max_A - new\_min_A) + new\_min_A$$

- z-score normalization

$$v' = \frac{v - mean_A}{stand\_dev_A}$$

- normalization by decimal scaling

$$v' = \frac{v}{10^j}$$  Where $j$ is the smallest integer such that $Max(|v'|) < 1$

# Data Reduction

# Data Reduction Strategies

- A data warehouse may store terabytes of data

  *Complex data analysis/mining may take a very long time to run on the complete data set*

- Data reduction

  *Obtain a reduced representation of the data set that is much smaller in volume but yet produce the same (or almost the same) analytical results*

# Dimensionality Reduction

- Feature selection (i.e., attribute subset selection):
  - Select a minimum set of features such that the probability distribution of different classes given the values for those features is as close as possible to the original distribution given the values of all features
  - reduce # of patterns in the patterns, easier to understand

- Heuristic methods (due to exponential # of choices):
  - step-wise forward selection
  - step-wise backward elimination
  - combining forward selection and backward elimination
  - decision-tree induction

# Handling Inconsistent Data

- Manual correction (expensive and tedious)

- Use routines designed to detect inconsistencies and manually correct them. E.g., the routine may use the check global constraints (age>10) or functional dependencies

- Other inconsistencies (e.g., between names of the same attribute) can be corrected during the data integration process

# Reference



THIRD EDITION

DATA MINING
Concepts and Techniques

Jiawei Han | Micheline Kamber | Jian Pei