



## 3.1 Introduction

# Types of Data Preprocessing

Data cleaning

• Fill in missing values, smooth noisy data, identify or remove outliers, resolve inconsistencies

Data integration

• Integration of multiple databases, data cubes, or files

Data transformation

• Normalization and aggregation

Data reduction

• Reduce number of records, attributes or attribute values

SFU, CMPT 740, 03-3, Martin Ester



# 3.2 Data Cleaning

## Handling Missing Data

- Ignore the record: usually done when class label is missing
- Fill in missing value manually: tedious + infeasible?
- Use a default to fill in the missing value:

e.g., "unknown", a new class, . . .

- Use the attribute mean to fill in the missing value for classification: mean for all records of the same class
- Use the most probable value to fill in the missing value: inference-based such as Bayesian formula or decision tree

SFU, CMPT 740, 03-3, Martin Ester



# 3.2 Data Cleaning

## Handling Noisy Data

#### Binning

- sort data and partition into (equi-depth) bins
- smooth by bin means, bin median, bin boundaries, etc.

#### Regression

• smooth by fitting a regression function

#### Clustering

• detect and remove outliers

Combined computer and human inspection

• detect suspicious values and check by human

SFU, CMPT 740, 03-3, Martin Ester









3.3 Data Integration	
Approach	
Identification	
• Detect corresponding tables from different sources manual	
<ul> <li>Detect corresponding attributes from different sources may use correlation analysis</li> <li>e.g., A.cust-id ≡ B.cust-#</li> </ul>	
<ul> <li>Detect duplicate records from different sources involves approximate matching of attribute values e.g. 3.14283 ≡ 3.1, Schwartz ≡ Schwarz</li> </ul>	
Treatment	
Merge corresponding tables	
• Use attribute values as synonyms	
Remove duplicate records	
$\mathbb{Q}$ data warehouses are already integrated	
SFU, CMPT 740, 03-3, Martin Ester 95	





# 3.4 Data Transformation

## Discretization

Three types of attributes

- Nominal (categorical) values from an unordered set
- Ordinal values from an ordered set
- Continuous (numerical) real numbers

Motivation for discretization

- Some data mining algorithms only accept categorical attributes
- May improve understandability of patterns

SFU, CMPT 740, 03-3, Martin Ester

98









## 3.5 Data Reduction

## Feature Selection

Feature selection methods

- Feature independence assumption: choose features independently by their significance
- Greedy bottom-up feature selection:
  - The best single-feature is picked first
  - Then next best feature condition to the first, ...
- Greedy top-down feature elimination:
  - Repeatedly eliminate the worst feature
- Branch and bound
  - Returns optimal set of features
  - Requires monotone structure of the feature space

SFU, CMPT 740, 03-3, Martin Ester











