

# Introduction to Data Mining

Key components

Collection  
Preprocessing  
Mining  
Evaluation  
Presentation

Why?

- explosive data growth
- Drawing in data
- Need automatical

4 views

- Data
  - The (3, 4, 5) vs Relational, transaction data  
Sequential, Temporal, streaming data
- Knowledge
  - Freq pattern, association, correlation  
Categorization  
Anomaly, outliers
- Application
  - Market Analysis, target advertisement  
Healthcare, medical research  
Security

Technical

- Freq pattern analysis
  - Itemset structure
  - Assoc.
  - Correlation analysis
- Classification
  - pre-defined classes
  - Need training data
  - Build model to distinguish classes
- Prediction
  - Numerical
  - No predefined classes
- Clustering
  - Intra-cluster similarity
  - Inter-cluster dissimilarity
- Trend and evolution analysis — changes overtime
- Anomaly detection — Anomaly / outlier

What?

Extraction  
Interesting  
Huge

Data Understanding

Data Objects & Attribute

Descriptive Statistics

Data Visualization

## Introduction to Data Understanding

Data Similarity

Key Measures

Application

Harnessing data understanding and similarity for effective data analysis and decision making

Handling missing data

Data Integration techniques

Normalization and scaling

Dimensionality Reduction

## Introduction to Data Preprocessing

### Data Preprocessing

Cleaning

missing value  
handling duplicate records  
correcting inaccuracies  
Effective data cleaning

Integration

challenge  
Data format  
Data redundancy

Benefits

Comprehensive analysis  
Improved decision-making  
Data integration  
Streamlining data utilization

Transformation

Normalization  
Attribute Construction

Reduction

Dimensionality Reduction  
Numerosity Reduction

Benefit of Data preprocessing

Feature Selection

# Introduction to Data Warehousing

