

TIM 245 Lecture 2 (4/5/17)

Agenda

- 1) Review of Lecture 1
- 2) Data
- 3) Data pre processing
- 4) Project kick off and Project Phase I assignment
- 5) work on project

② Data

Dataset: collection of instances (rows)
across a set of attributes
(columns)

Example: class grades

Attributes

Instances

Assignment ID	Name	Date	Level	Score	Grade
01	John	6/10/17	Soph	95	A

Types of attributes

Nominal = Value are distinct ($=, \neq$)

Example: name

Ordinal: Values are distinct and can be ordered ($=, \neq, <, >$)

Example: Level

Interval: Values are distinct, can be ordered, and have a meaningful difference ($=, \neq, <, >, +, -$)

Ratio: Value are distinct, can be ordered, and have a meaningful difference and ratio ($=, \neq, <, >, +, -, *, /$)

Example: score

Types of data-sets

Record Data : flat files (CSV)
 (Tables) hierarchical files (XML, JSON)
 Relational databases (SQL)

Time series Data : flat files (CSV)
 (ordered data)

Text Data : hierarchical files (XML, JSON)

Graph Data : hierarchical files (XML, JSON)
 Graph databases (Neo4J, MongoDB)

Where to get data :

Data mining : Kaggle
 Competitions : Data Driven

Website API : Facebook
 Twitter
 Google
 Yelp
 Zillow

Website : Forum
 Scoping : product listings

Open Source : UCI Machine Learning
 Datasets : Data.gov
 AWS Public Datasets

Create your own data : Wearable devices,
 research experiments

③ Data Preprocessing

Why do we need to do pre-processing?

Real-World data is often

- 1) Incomplete : lacks certain attributes of interest
- 2) Noisy : contains errors or outlier values
- 3) Inconsistent : uses different values for the same thing (lingo problem)

Addressing these issues will improve the quality of the data mining results:

Descriptive Analysis : patterns that are more useful and relevant to the problem under consideration

Predictive Analysis : models that have better predictive accuracy and generalization

There are four basic steps in data pre-processing :

1) Data Cleaning : filling in missing values,
(Instances / Rows) Smoothing noisy data
removing outliers,
removing duplicates

2) Data Integration : joining multiple datasets
(attributes / columns) entity recognition
repeat data cleaning

3) Data Transformation : normalization
(instances / rows) binning

4) Data Reduction : attribute selection
(attributes / columns, Dimensionality Reduction
Instances / rows) Sampling