Week 2 Engineering Data (Part I)

Seokho Chi

Assistant Professor I Ph.D. SNU Construction Innovation Lab



Source: Tan, Kumar, Steinback (2006)



What is Data?

Objects

Collection of data objects and their attributes

- An attribute is a property or characteristic of an object
 - Examples: eye color of a person, temperature, etc.
 - Attribute is also known as variable, field, characteristic, or feature
- A collection of attributes describe an object
 - Object is also known as record, point, case, sample, entity, or instance

Attributes



Data set for predicting borrowers who will default on loan payments

Types of Attributes

	Attribute Type	Description	Examples	Operations	
	Nominal	The values of a nominal attribute are just different names, i.e., nominal attributes provide only enough information to distinguish one object from another. (DISTINCTNESS =, \neq)	zip codes, employee ID numbers, eye color, sex: { <i>male, female</i> }	mode, entropy, contingency correlation, χ^2 test	
)	Ordinal	The values of an ordinal attribute provide enough information to order objects. (ORDER <, >)	hardness of minerals, {good, better, best}, grades, street numbers	median, percentiles, rank correlation, run tests, sign tests	
	Interval	For interval attributes, the differences between values are meaningful, i.e., a unit of measurement exists. (ADDITION +, -)	calendar dates, temperature in Celsius or Fahrenheit (differ in the location of their zero value)	mean, standard deviation, Pearson's correlation, <i>t</i> and <i>F</i> tests	
	Ratio	For ratio variables, both differences and ratios are meaningful. (MULTIPLICATION *, /)	monetary quantities, counts, age, mass, length, electrical current	geometric mean, harmonic mean, percent variation	

Discrete and Continuous Attributes

- Discrete Attribute
 - Has only a finite or countable infinite set of values
 - Examples: zip codes, counts, or the set of words in a collection of documents
 - Often represented as integer variables.
 - Note: binary attributes are a special case of discrete attributes
- Continuous Attribute
 - Has real numbers as attribute values
 - Examples: temperature, height, or weight.
 - Practically, real values can only be measured and represented using a finite number of digits.
 - Continuous attributes are typically represented as floating-point variables.

Types of data sets

- Record
 - Data Matrix
 - Document Data
 - Transaction Data
- Graph
 - World Wide Web
 - Molecular Structures
- Ordered
 - Spatial Data
 - Temporal Data
 - Sequential Data
 - Genetic Sequence Data

Record Data

 Data that consists of a collection of records, each of which consists of a fixed set of attributes

Tid	House Owner	Marital Status	Taxable Income	Defau Ited Borro wer
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Data Matrix

- If data objects have the same fixed set of numeric attributes, then the data objects can be thought of as points in a multi-dimensional space, where each dimension represents a distinct attribute → Possible 3D Plotting
- Such data set can be represented by an m by n matrix, where there are m rows, one for each object, and n columns, one for each attribute

Projection of x Load	Projection of y load	Distance	Load	Thickness
10.23	5.27	15.22	2.7	1.2
12.65	6.25	16.22	2.2	1.1

Document Data

- Each document becomes a `term' vector,
 - each term is a component (attribute) of the vector,
 - the value of each component may be the number of times the corresponding term occurs in the document.

	team	coach	pla y	ball	score	game	n <u>V</u> .	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

Transaction Data

- A special type of record data, where
 - each record (transaction) involves a set of items.
 - For example, consider a grocery store. The set of products purchased by a customer during one shopping trip constitute a transaction, while the individual products that were purchased are the items.

TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

Graph Data

Examples: Generic graph and HTML Links



 Data Mining

 Graph Partitioning

<|i>

Parallel Solution of Sparse Linear System of Equations

N-Body Computation and Dense Linear System Solvers



Chemical Data

Benzene Molecule: C₆H₆



Ordered Data

Sequences of transactions



Ordered Data

DNA sequencing: the process of determining the exact order of the 3 billion chemical building blocks (called bases and abbreviated A, T, C, and G) that make up the DNA of the 24 different human chromosomes

GGTTCCGCCTTCAGCCCGCGCGCC CGCAGGGCCCGCCCGCGCGCGTC GAGAAGGGCCCGCCTGGCGGGGCG GGGGGAGGCGGGGGCCGCCCGAGC CCAACCGAGTCCGACCAGGTGCC CCCTCTGCTCGGCCTAGACCTGA GCTCATTAGGCGGCAGCGGACAG GCCAAGTAGAACACGCGAAGCGC

Ordered Data

Spatio-Temporal Data

Jan

Average Monthly Temperature of land and ocean

Data Quality

- What kinds of data quality problems?
- How can we detect problems with the data?
- What can we do about these problems?
- Examples of data quality problems:
 - Noise and outliers
 - missing values
 - duplicate data

Noise

- Noise refers to modification of original values
 - Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen



Outliers

 Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set



Missing Values

- Reasons for missing values (data is not available)
 - Information is not collected

(e.g., people decline to give their age and weight)

Attributes may not be applicable to all cases

(e.g., annual income is not applicable to children)

Duplicate Data

- Data set may include data objects that are duplicates, or almost duplicates of one another
 - Major issue when merging data from heterogeneous sources: where assign? or different objects?
- Examples:
 - Same person with multiple email addresses

Why Data Preprocessing?

- Data in the real world is dirty
 - incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
 - noisy: containing errors or outliers
 - inconsistent: containing discrepancies in codes or names
- Less quality data, less quality mining results!
 - Quality decisions must be based on quality data

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

Forms of data preprocessing



Data Cleaning

- Data cleaning tasks:
 - Fill in missing values
 - Identify outliers and smooth out noisy data
 - Correct inconsistent data

Missing Data

- Data is not always available
 - E.g., many records 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 record: 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?
- Use a global constant to fill in the missing value: e.g., "unknown", a new class?!
- Use the attribute mean to fill in the missing value
- Use the attribute mean for all samples belonging to the same class to fill in the missing value: smarter
- Use the most probable value to fill in the missing 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

How to Handle Noisy Data?

- Binning method:
 - first sort data and partition into 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
- Regression
 - smooth by fitting the data into regression functions

Data Integration

Metadata: DB 시스템에서 데이터 관리상 필요한 작성자, 목적, 저장 장소 등 속성에 관한 데이터

- Data integration:
 - combines data from multiple sources into a coherent store
- Schema integration
 - integrate metadata from different sources
 - Entity identification problem: identify real world entities from multiple data sources

e.g., A.cust-id = B.cust-# | (02)111-1111 = 111-1111 = 02 111 1111

- Detecting and resolving data value conflicts
 - for the same real world entity, attribute values from different sources are different e.g., work phone = phone = tel.
 - possible reasons: different representations, different scales, e.g., metric vs. British units

Handling Redundant Data in Data Integration

- Redundant data occur often when integration of multiple databases
 - The same attribute may have different names in different databases
 - One attribute may be a "derived" attribute in another table, e.g., annual revenue
- Careful integration of the data from multiple sources may help reduce/avoid redundancies and inconsistencies and improve mining speed and quality

Data Transformation

- Smoothing: remove noise from data
- Aggregation: summarization, data cube construction
- 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 Preprocessing: Generic Approaches

- Aggregation
- Sampling
- Dimensionality Reduction
- Feature subset selection
- Feature creation
- Discretization and Binarization
- Attribute Transformation

Aggregation

 Combining two or more attributes (or objects) into a single attribute (or object)

Purpose

- Data reduction
 - Reduce the number of attributes or objects
- Change of scale
 - Cities aggregated into regions, states, countries, etc
- More "stable" data
 - Aggregated data tends to have less variability

Aggregation

Variation of Precipitation in Australia

Sampling

- Sampling is the main technique employed for data selection.
 - It is often used for both the preliminary investigation of the data and the final data analysis.
- Statisticians sample because obtaining the entire set of data (population) of interest is too expensive or time consuming.
- Sampling is used in data mining because processing the entire set of data of interest is too expensive or time consuming.

Sampling (2)

- The key principle for effective sampling is the following:
 - using a sample will work almost as well as using the entire data sets, if the sample is representative

 A sample is representative if it has approximately the same property (of interest) as the original set of data

Types of Sampling

- Simple Random Sampling
 - There is an equal probability of selecting any particular item
- Sampling without replacement
 - As each item is selected, it is removed from the population
- Sampling with replacement
 - Objects are not removed from the population as they are selected for the sample.
 - In sampling with replacement, the same object can be picked up more than once
- Stratified sampling
 - Split the data into several partitions; then draw random samples from each partition

Sample Size

8000 points

2000 Points

500 Points

Curse of Dimensionality

Data analysis becomes significantly harder as the dimensionality of the data increases. *Dimension = Attribute

- When dimensionality increases, data becomes increasingly sparse in the space that it occupies
- Definitions of density and distance between points, which is critical for clustering and outlier detection, become less meaningful

- Randomly generate 500 points
- Compute difference between max and min distance between any pair of points

Dimensionality Reduction

- Purpose:
 - Avoid curse of dimensionality
 - Reduce amount of time and memory required by data mining algorithms
 - Allow data to be more easily visualized
 - May help to eliminate irrelevant features or reduce noise
- Techniques: Linear Algebra Techniques
 - Principle Component Analysis
 - Singular Value Decomposition
 - Others: supervised and non-linear techniques

Principal Component Analysis

- Linear algebra technique for continuous attributes that finds new attributes (principal components) that
 - Are linear combinations of the original attributes;
 - Orthogonal (perpendicular) to each other; and,
 - Capture the maximum amount of variation in the data
- For instance,
 - The first two principal components capture: as much of the variation in the data as is possible with two orthogonal attributes that are linear combinations of the original attributes

Feature Subset Selection

- Redundant features
 - duplicate much or all of the information contained in one or more other attributes
 - Example: purchase price of a product and the amount of sales tax paid
- Irrelevant features
 - contain no information that is useful for the data mining task at hand
 - Example: students' ID is often irrelevant to the task of predicting students' GPA

Feature Subset Selection

Techniques:

Brute-force approach (ideal approach)

Try all possible feature subsets as input to data mining algorithm

Embedded approaches

- Feature selection occurs naturally as part of the data mining algorithm
- The data mining algorithm itself decides which attributes to use and which to ignore

Filter approaches

- Features are selected before data mining algorithm is run
- Independent of data mining task
- We might select sets of attributes whose pairwise correlation is as low as possible

– Wrapper approaches

- Use the data mining algorithm as a black box to find best subset of attributes
- Not the all possible subsets

Feature Creation

- Create new attributes that can capture the important information in a data set much more efficiently than the original attributes
- Three general methodologies:
 - Feature Extraction
 - − For classifying pics that contain human face, a set of pixels (raw data) \rightarrow higher level new feature such as certain types of edges and areas
 - Mapping Data to New Space
 - Feature Construction
 - combining features

Mapping Data to a New Space

- Fourier transform (linear, signal processing)
- Wavelet transform (discrete vs continuous)

Discretization

- Discretization of continuous attributes:
 - Transform continuous attribute to categorical attribute
 e.g., 1, 2, 3, 4, 5, 6, 7, 8, 9 → low, medium, high
 - Divide the range of a continuous attribute into intervals in order to reduce data size
 - Subtasks:
 - Decide how many categories to have
 - Determine how to map the values of the continuous values to these categories

Discretization Using Class Labels (supervised)

Entropy based approach

3 categories for both x and y

5 categories for both x and y

Discretization Without Using Class Labels (unsupervised)

Unknown Classes

Attribute Transformation

- A function that maps the entire set of values of a given attribute to a new set of replacement values such that each old value can be identified with one of the new values
 - Simple functions: x^k , log(x), e^x , |x|
 - Standardization and Normalization

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