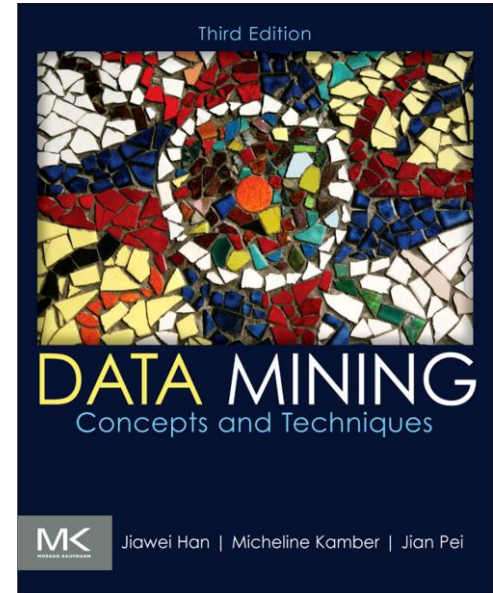


CISI612 March 16, 2024



3. Data Preprocessing

Dr. Kadan Aljoumaa

Kadan.aljoumaa@hiast.edu.sy

Chapter 3: Data Preprocessing

1. Data Preprocessing: An Overview

- a. Data Quality
- b. Major Tasks in Data Preprocessing

2. Data Cleaning

3. Data Integration

4. Data Reduction

5. Data Transformation and Data Discretization

6. Summary

Data Quality: Why Preprocess the Data?

❑ Measures for data quality: A multidimensional view

- Accuracy: correct or wrong, accurate or not
- Completeness: not recorded, unavailable, ...
- Consistency: some modified but some not, dangling, ...
- Timeliness: timely update?
- Believability: how trustable the data are correct?
- Interpretability: how easily the data can be understood?

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 reduction

- Dimensionality reduction
- Numerosity reduction
- Data compression

❑ Data transformation and data discretization

- Normalization
- Concept hierarchy generation

2. Data Cleaning

- ❑ Data in the Real World Is Dirty: Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, transmission error
 - incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
 - ▶ e.g., *Occupation*=“ ” (missing data)
 - noisy: containing noise, errors, or outliers
 - ▶ e.g., *Salary*=“-10” (an error)
 - inconsistent: containing discrepancies in codes or names, e.g.,
 - ▶ *Age*=“42”, *Birthday*=“03/07/2010”
 - ▶ Was rating “1, 2, 3”, now rating “A, B, C”
 - ▶ discrepancy between duplicate records
 - Intentional (e.g., *disguised missing data*)
 - ▶ Jan. 1 as everyone’s birthday?

Incomplete (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 (when doing classification)—not effective when the % 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 be due to
 - faulty data collection instruments
 - data entry problems
 - data transmission problems
 - technology limitation
 - inconsistency in naming convention
- ❑ Other data problems which require data cleaning
 - duplicate records
 - incomplete data
 - inconsistent data

How to Handle Noisy Data?

❑ Binning

- first sort data and partition into (equal-frequency) bins
- then one can smooth by bin means, smooth by bin median, smooth by bin boundaries, etc.

❑ Regression

- smooth by fitting the data into regression functions

❑ Clustering

- detect and remove outliers

❑ Combined computer and human inspection

- detect suspicious values and check by human (e.g., deal with possible outliers)

Data Cleaning as a Process

❑ Data discrepancy detection

- Use metadata (e.g., domain, range, dependency, distribution)
- Check field overloading
- Check uniqueness rule, consecutive rule and null rule
- Use commercial tools
 - ▶ Data scrubbing: use simple domain knowledge (e.g., postal code, spell-check) to detect errors and make corrections
 - ▶ Data auditing: by analyzing data to discover rules and relationship to detect violators (e.g., correlation and clustering to find outliers)

❑ Data migration and integration

- Data migration tools: allow transformations to be specified
- ETL (Extraction/Transformation/Loading) tools: allow users to specify transformations through a graphical user interface

❑ Integration of the two processes

- Iterative and interactive (e.g., Potter's Wheels)

3. Data Integration

❑ Data integration:

- Combines data from multiple sources into a coherent store

❑ Schema integration: e.g., $A.cust-id \equiv B.cust-#$

- Integrate metadata from different sources

❑ Entity identification problem:

- Identify real world entities from multiple data sources, e.g., Bill Clinton = William Clinton

❑ Detecting and resolving data value conflicts

- For the same real world entity, attribute values from different sources are different
- Possible reasons: different representations, different scales, e.g., metric vs. British units

Handling Redundancy in Data Integration

- ❑ Redundant data occur often when integration of multiple databases
 - *Object identification*: The same attribute or object may have different names in different databases
 - *Derivable data*: One attribute may be a “derived” attribute in another table, e.g., annual revenue
- ❑ Redundant attributes may be able to be detected by *correlation analysis* and *covariance analysis*
- ❑ Careful integration of the data from multiple sources may help reduce/avoid redundancies and inconsistencies and improve mining speed and quality

Correlation Analysis (Nominal Data)

▶ χ^2 (chi-square) test

$$\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$$

- ▶ The larger the χ^2 value, the more likely the variables are related
- ▶ The cells that contribute the most to the χ^2 value are those whose actual count is very different from the expected count
- ▶ Correlation does not imply causality
 - ▶ # of hospitals and # of car-theft in a city are correlated
 - ▶ Both are causally linked to the third variable: population

Chi-Square Calculation: An Example

	Play chess	Not play chess	Sum (row)
Like science fiction	250(90)	200(360)	450
Not like science fiction	50(210)	1000(840)	1050
Sum(col.)	300	1200	1500

- ▶ χ^2 (chi-square) calculation (numbers in parenthesis are expected counts calculated based on the data distribution in the two categories)

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} = 507.93$$

- ▶ It shows that like_science_fiction and play_chess are correlated in the group

Correlation Analysis (Numeric Data)

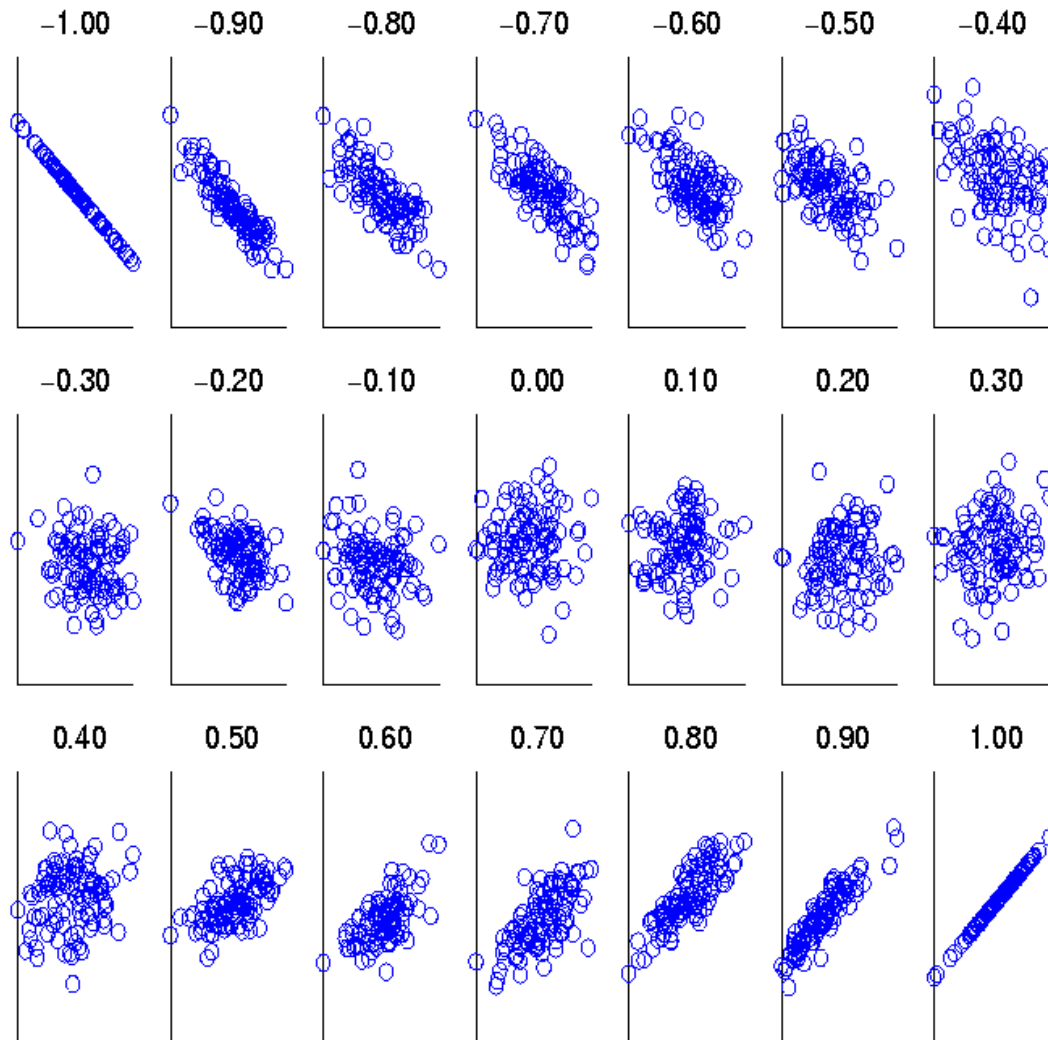
- ▶ Correlation coefficient (also called Pearson's product moment coefficient)

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{(n-1)\sigma_A\sigma_B} = \frac{\sum_{i=1}^n (a_i b_i) - n\bar{A}\bar{B}}{(n-1)\sigma_A\sigma_B}$$

where n is the number of tuples, \bar{A} and \bar{B} are the respective means of A and B , σ_A and σ_B are the respective standard deviation of A and B , and $\sum(a_i b_i)$ is the sum of the AB cross-product.

- ▶ If $r_{A,B} > 0$, A and B are positively correlated (A 's values increase as B 's). The higher, the stronger correlation.
- ▶ $r_{A,B} = 0$: independent; $r_{AB} < 0$: negatively correlated

Visually Evaluating Correlation



**Scatter plots
showing the
similarity from
-1 to 1.**

Correlation (viewed as linear relationship)

- ❑ Correlation measures the linear relationship between objects
- ❑ To compute correlation, we standardize data objects, A and B, and then take their dot product

$$a'_k = (a_k - \text{mean}(A)) / \text{std}(A)$$

$$b'_k = (b_k - \text{mean}(B)) / \text{std}(B)$$

$$\text{correlation}(A, B) = A' \bullet B'$$

Covariance (Numeric Data)

- ❑ **Covariance is similar to correlation**

$$Cov(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

- ❑ **Correlation coefficient:** $r_{A,B} = \frac{Cov(A, B)}{\sigma_A \sigma_B}$

where n is the number of tuples, \bar{A} and \bar{B} are the respective mean or **expected values** of A and B , σ_A and σ_B are the respective standard deviation of A and B .

- ❑ **Positive covariance:** If $Cov_{A,B} > 0$, then A and B both tend to be larger than their expected values.
- ❑ **Negative covariance:** If $Cov_{A,B} < 0$ then if A is larger than its expected value, B is likely to be smaller than its expected value.
- ❑ **Independence:** $Cov_{A,B} = 0$ but the converse is not true:
 - Some pairs of random variables may have a covariance of 0 but are not independent. Only under some additional assumptions (e.g., the data follow multivariate normal distributions) does a covariance of 0 imply independence

Co-Variance: An Example

$$Cov(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

- It can be simplified in computation as

$$Cov(A, B) = E(A \cdot B) - \bar{A}\bar{B}$$

- Suppose two stocks A and B have the following values in one week: (2, 5), (3, 8), (5, 10), (4, 11), (6, 14).
- Question: If the stocks are affected by the same industry trends, will their prices rise or fall together?
 - $E(A) = (2 + 3 + 5 + 4 + 6) / 5 = 20 / 5 = 4$
 - $E(B) = (5 + 8 + 10 + 11 + 14) / 5 = 48 / 5 = 9.6$
 - $Cov(A, B) = (2 \times 5 + 3 \times 8 + 5 \times 10 + 4 \times 11 + 6 \times 14) / 5 - 4 \times 9.6 = 4$
- Thus, A and B rise together since $Cov(A, B) > 0$.

4. Data Reduction Strategies

- ❑ **Data reduction:** Obtain a reduced representation of the data set that is much smaller in volume but yet produces the same (or almost the same) analytical results
- ❑ **Why data reduction?** — A database/data warehouse may store terabytes of data. Complex data analysis may take a very long time to run on the complete data set.
- ❑ **Data reduction strategies**
 - **Dimensionality reduction**, e.g., remove unimportant attributes
 - Wavelet transforms
 - Principal Components Analysis (PCA)
 - Feature subset selection, feature creation
 - **Numerosity reduction** (some simply call it: Data Reduction)
 - Regression and Log-Linear Models
 - Histograms, clustering, sampling
 - Data cube aggregation
 - **Data compression**

Data Reduction 1: Dimensionality Reduction

❑ **Curse of dimensionality**

- When dimensionality increases, data becomes increasingly sparse
- Density and distance between points, which is critical to clustering, outlier analysis, becomes less meaningful
- The possible combinations of subspaces will grow exponentially

❑ **Dimensionality reduction**

- Avoid the curse of dimensionality
- Help eliminate irrelevant features and reduce noise
- Reduce time and space required in data mining
- Allow easier visualization

❑ **Dimensionality reduction techniques**

- Wavelet transforms
- Principal Component Analysis
- Supervised and nonlinear techniques (e.g., feature selection)

Attribute Subset Selection

- ❑ Another way to reduce dimensionality of data
- ❑ Redundant attributes
 - Duplicate much or all of the information contained in one or more other attributes
 - E.g., purchase price of a product and the amount of sales tax paid
- ❑ Irrelevant attributes
 - Contain no information that is useful for the data mining task at hand
 - E.g., students' ID is often irrelevant to the task of predicting students' GPA

Data Reduction 2: Numerosity Reduction

- ❑ Reduce data volume by choosing alternative, *smaller forms* of data representation
- ❑ **Parametric methods** (e.g., regression)
 - Assume the data fits some model, estimate model parameters, store only the parameters, and discard the data (except possible outliers)
 - Ex.: Log-linear models—obtain value at a point in m -D space as the product on appropriate marginal subspaces
- ❑ **Non-parametric methods**
 - Do not assume models
 - Major families: histograms, clustering, sampling, ...

Parametric Data Reduction: Regression and Log-Linear Models

❑ Linear regression

- Data modeled to fit a straight line
- Often uses the least-square method to fit the line

❑ Multiple regression

- Allows a response variable Y to be modeled as a linear function of multidimensional feature vector

❑ Log-linear model

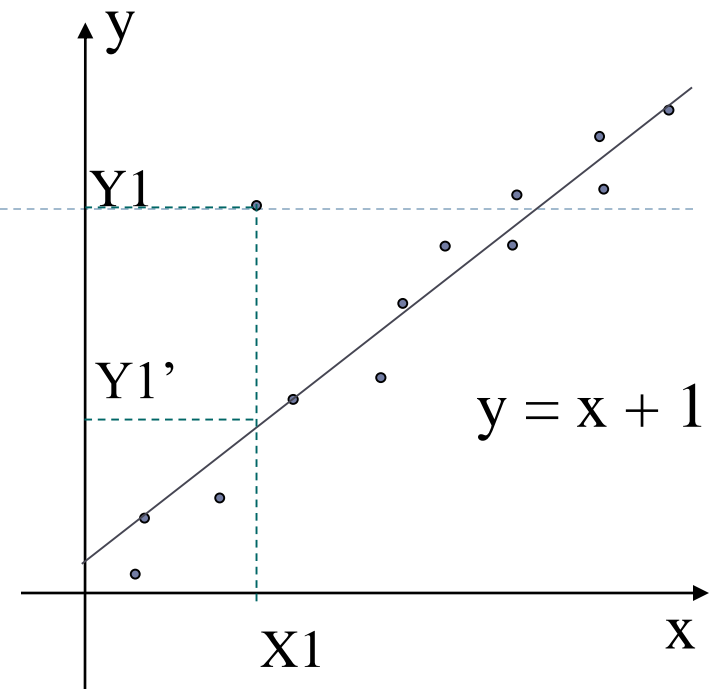
- Approximates discrete multidimensional probability distributions

Regression Analysis

❑ Regression analysis: A collective name for techniques for the modeling and analysis of numerical data consisting of values of a **dependent variable** (also called **response variable** or *measurement*) and of one or more *independent variables* (aka. **explanatory variables** or **predictors**)

❑ The parameters are estimated so as to give a **"best fit"** of the data

❑ Most commonly the best fit is evaluated by using the **least squares method**, but other criteria have also been used



Used for prediction (including forecasting of time-series data), inference, hypothesis testing, and modeling of causal relationships

Clustering

- ❑ Partition data set into clusters based on similarity, and store cluster representation (e.g., centroid and diameter) only
- ❑ Can be very effective if data is clustered but not if data is “smeared”
- ❑ Can have hierarchical clustering and be stored in multi-dimensional index tree structures
- ❑ There are many choices of clustering definitions and clustering algorithms
- ❑ Cluster analysis will be studied in depth in Chapter 10

Sampling

- ❑ Sampling: obtaining a small sample s to represent the whole data set N
- ❑ Allow a mining algorithm to run in complexity that is potentially sub-linear to the size of the data
- ❑ Key principle: Choose a **representative** subset of the data
 - Simple random sampling may have very poor performance in the presence of skew
 - Develop adaptive sampling methods, e.g., stratified sampling:
- ❑ Note: Sampling may not reduce database I/Os (page at a time)

Data Cube Aggregation

- ❑ The lowest level of a data cube (base cuboid)
 - The aggregated data for an individual entity of interest
 - E.g., a customer in a phone calling data warehouse
- ❑ Multiple levels of aggregation in data cubes
 - Further reduce the size of data to deal with
- ❑ Reference appropriate levels
 - Use the smallest representation which is enough to solve the task
- ❑ Queries regarding aggregated information should be answered using data cube, when possible

Data Reduction 3: Data Compression

❑ String compression

- There are extensive theories and well-tuned algorithms
- Typically lossless, but only limited manipulation is possible without expansion

❑ Audio/video compression

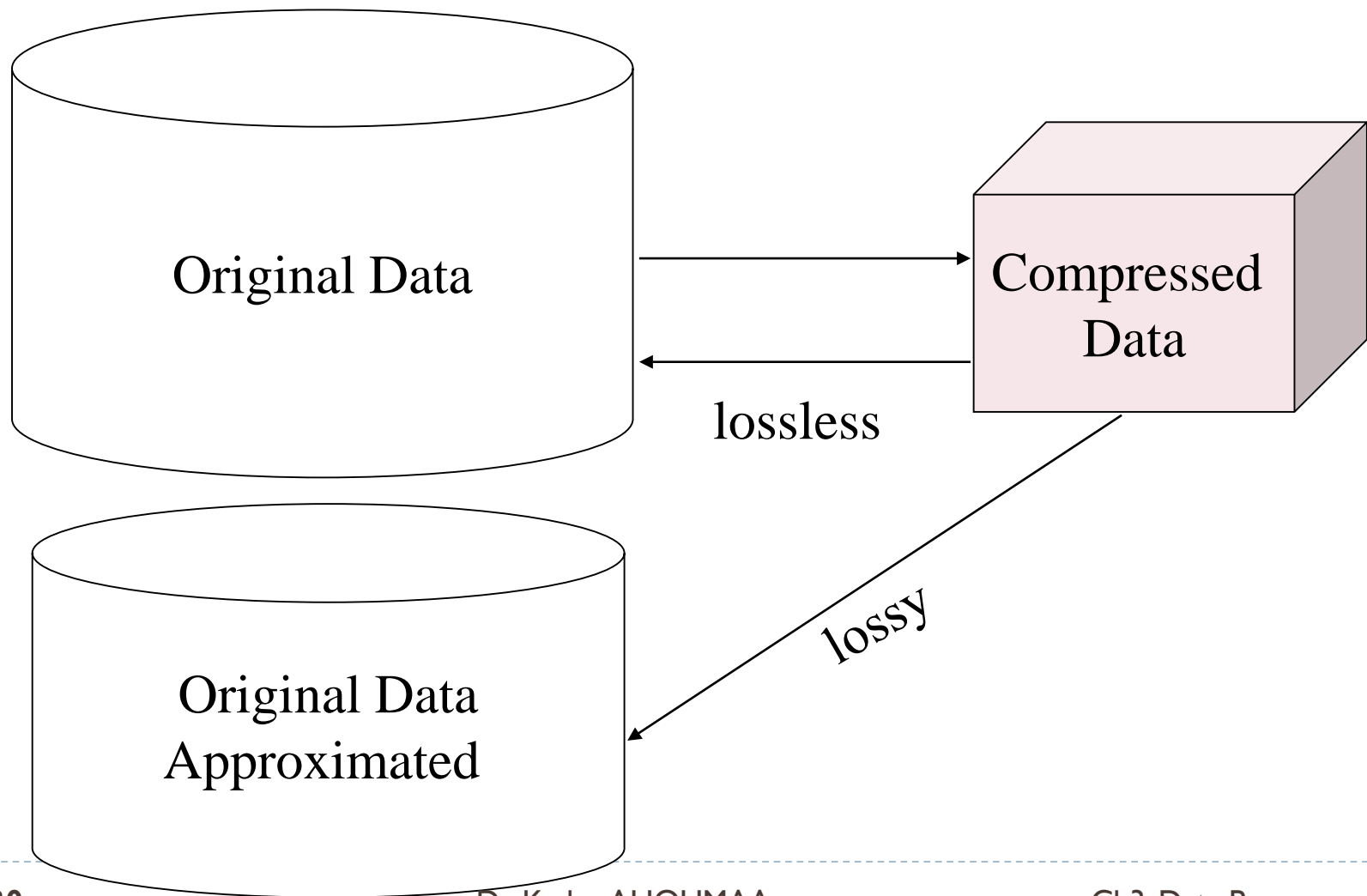
- Typically lossy compression, with progressive refinement
- Sometimes small fragments of signal can be reconstructed without reconstructing the whole

❑ Time sequence is not audio

- Typically short and vary slowly with time

❑ Dimensionality and numerosity reduction may also be considered as forms of data compression

Data Compression



5. Data Transformation

- ❑ A function that maps the entire set of values of a given attribute to a new set of replacement values s.t. each old value can be identified with one of the new values
- ❑ Methods
 - Smoothing: Remove noise from data
 - Attribute/feature construction
 - New attributes constructed from the given ones
 - Aggregation: Summarization, data cube construction
 - Normalization: Scaled to fall within a smaller, specified range
 - min-max normalization
 - z-score normalization
 - normalization by decimal scaling
 - Discretization: Concept hierarchy climbing

Normalization

- ▶ **Min-max normalization:** to $[new_min_A, new_max_A]$

$$v' = \frac{v - min_A}{max_A - min_A} (new_max_A - new_min_A) + new_min_A$$

- ▶ Ex. Let income range \$12,000 to \$98,000 normalized to $[0.0, 1.0]$. Then \$73,600 is mapped to $\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$

- ▶ **Z-score normalization** (μ : mean, σ : standard deviation):

$$v' = \frac{v - \mu_A}{\sigma_A}$$

$$\frac{73,600 - 54,000}{16,000} = 1.225$$

- ▶ Ex. Let $\mu = 54,000$, $\sigma = 16,000$. Then

- ▶ **Normalization by decimal scaling**

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

Discretization

❑ Three types of attributes

- Nominal—values from an unordered set, e.g., color, profession
- Ordinal—values from an ordered set, e.g., military or academic rank
- Numeric—real numbers, e.g., integer or real numbers

❑ Discretization: Divide the range of a continuous attribute into intervals

- Interval labels can then be used to replace actual data values
- Reduce data size by discretization
- Supervised vs. unsupervised
- Split (top-down) vs. merge (bottom-up)
- Discretization can be performed recursively on an attribute
- Prepare for further analysis, e.g., classification

Data Discretization Methods

- ❑ Typical methods: All the methods can be applied recursively
 - Binning
 - Top-down split, unsupervised
 - Histogram analysis
 - Top-down split, unsupervised
 - Clustering analysis (unsupervised, top-down split or bottom-up merge)
 - Decision-tree analysis (supervised, top-down split)
 - Correlation (e.g., χ^2) analysis (unsupervised, bottom-up merge)

6. Summary

- ❑ **Data quality:** accuracy, completeness, consistency, timeliness, believability, interpretability
- ❑ **Data cleaning:** e.g. missing/noisy values, outliers
- ❑ **Data integration** from multiple sources:
 - Entity identification problem
 - Remove redundancies
 - Detect inconsistencies
- ❑ **Data reduction**
 - Dimensionality reduction
 - Numerosity reduction
 - Data compression
- ❑ **Data transformation and data discretization**
 - Normalization
 - Concept hierarchy generation

References

- ❑ D. P. Ballou and G. K. Tayi. Enhancing data quality in data warehouse environments. *Comm. of ACM*, 42:73-78, 1999
- ❑ A. Bruce, D. Donoho, and H.-Y. Gao. Wavelet analysis. *IEEE Spectrum*, Oct 1996
- ❑ T. Dasu and T. Johnson. *Exploratory Data Mining and Data Cleaning*. John Wiley, 2003
- ❑ J. Devore and R. Peck. *Statistics: The Exploration and Analysis of Data*. Duxbury Press, 1997.
- ❑ H. Galhardas, D. Florescu, D. Shasha, E. Simon, and C.-A. Saita. Declarative data cleaning: Language, model, and algorithms. *VLDB'01*
- ❑ M. Hua and J. Pei. Cleaning disguised missing data: A heuristic approach. *KDD'07*
- ❑ H. V. Jagadish, et al., Special Issue on Data Reduction Techniques. *Bulletin of the Technical Committee on Data Engineering*, 20(4), Dec. 1997
- ❑ H. Liu and H. Motoda (eds.). *Feature Extraction, Construction, and Selection: A Data Mining Perspective*. Kluwer Academic, 1998
- ❑ J. E. Olson. *Data Quality: The Accuracy Dimension*. Morgan Kaufmann, 2003
- ❑ D. Pyle. *Data Preparation for Data Mining*. Morgan Kaufmann, 1999
- ❑ V. Raman and J. Hellerstein. Potters Wheel: An Interactive Framework for Data Cleaning and Transformation, *VLDB'2001*
- ❑ T. Redman. *Data Quality: The Field Guide*. Digital Press (Elsevier), 2001
- ❑ R. Wang, V. Storey, and C. Firth. A framework for analysis of data quality research. *IEEE Trans. Knowledge and Data Engineering*, 7:623-640, 1995