



Dr. Vishwanath Karad

**MIT WORLD PEACE  
UNIVERSITY** | PUNE

TECHNOLOGY, RESEARCH, SOCIAL INNOVATION & PARTNERSHIPS

## **Data Engineering Concepts UNIT II**

**SCHOOL OF COMPUTER ENGINEERING & TECHNOLOGY**

# Syllabus

- Data Preprocessing: An Overview, Methods: Data Cleaning, Data Integration, Data Reduction, Data Transformation, Data Discretization.
- **Data Cleaning:** Handling Missing Values, Noisy Data, Data Cleaning as a Process, **Data Integration:** Entity Identification Problem, Redundancy and Correlation Analysis, Tuple Duplication, Data Value Conflict Detection and Resolution, **Data Reduction:** Attribute Subset
- Selection, Histograms, Sampling, **Data discretization** – binning, histogram analysis, decision tree and correlation analysis, concept hierarchy for nominal data.

# Data Pre-processing: An Overview

- Data preprocessing is the process of transforming raw data into an understandable format.
- It is also an important step in data mining as we cannot work with raw data.
- The quality of the data should be checked before applying machine learning or data mining algorithms.
- The goal of data preprocessing is to improve the quality of the data and to make it more suitable for the specific data mining task.
- Data preprocessing is a process of preparing the raw data and making it suitable for a machine learning model and it is the first and crucial step while creating a machine learning model.

# Why is Data Preprocessing Important?

- A real-world data generally contains noises, missing values, and maybe in an unusable format which cannot be directly used for machine learning models.
- Data preprocessing is required tasks for cleaning the data and making it suitable for a machine learning model which also increases the accuracy and efficiency of a machine learning model.
- *“How can the data be preprocessed in order to help improve the quality of the data and, consequently, of the mining results? How can the data be preprocessed so as to improve the efficiency and ease of the mining process?”*

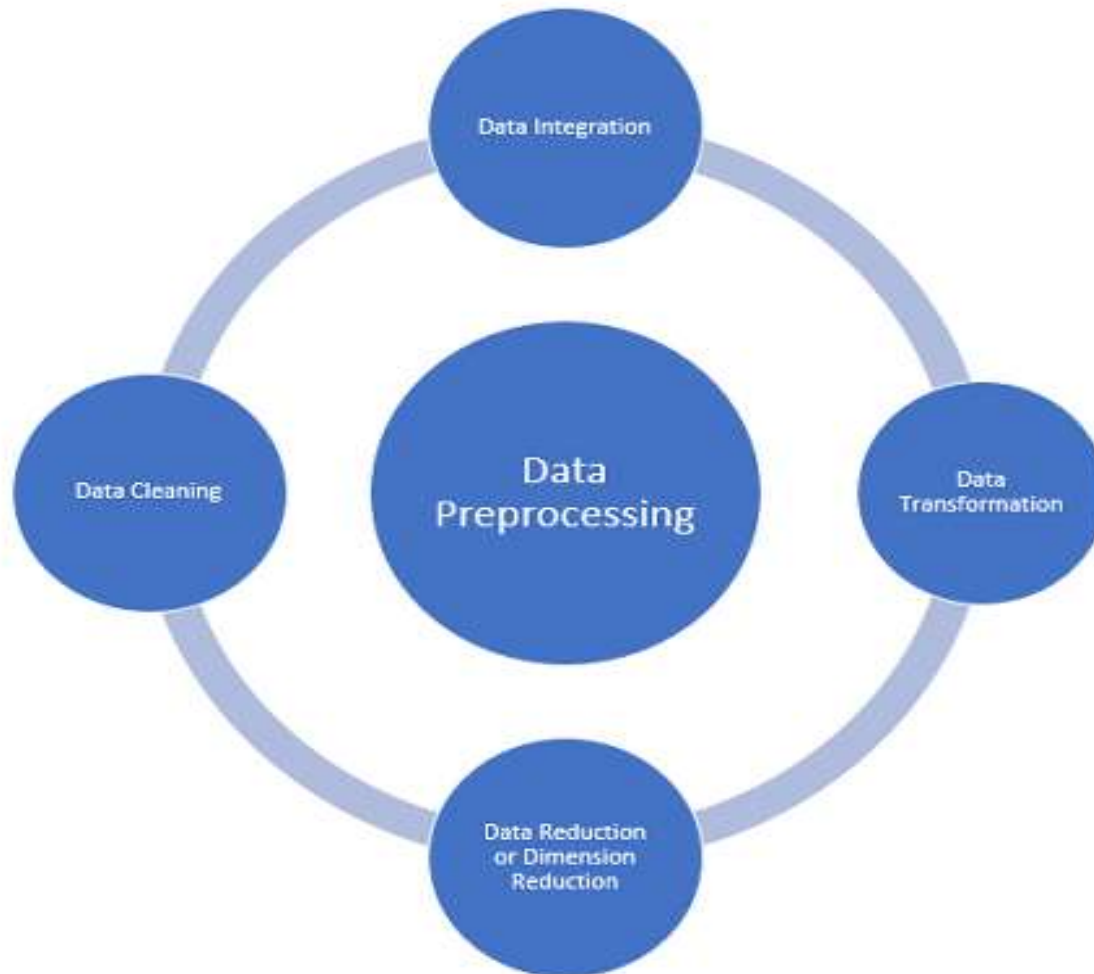
# Why is Data Preprocessing Important?

- Preprocessing of data is mainly to check the data quality and it can be checked with the following parameters:
- **Accuracy:** To check whether the data entered is correct or not.
- **Completeness:** To check whether the data is available or not recorded.
- **Consistency:** To check whether the same data is kept in all the places that do or do not match.
- **Timeliness:** The data should be updated correctly.
- **Believability:** The data should be trustable.
- **Interpretability:** The understandability of the data.

# Why is Data Preprocessing Important?

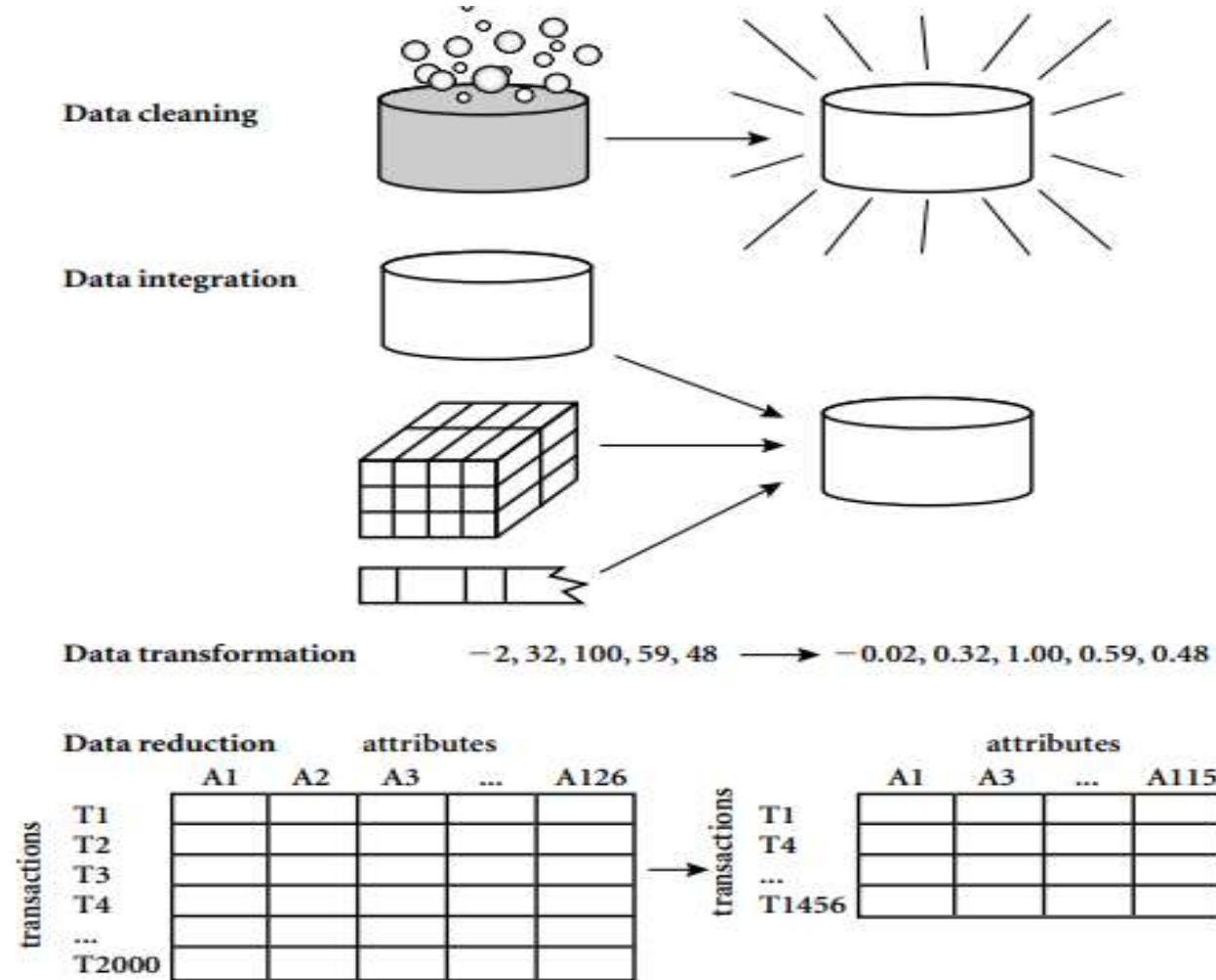
- *There are a number of data preprocessing techniques: **Data cleaning** can be applied to remove noise and correct inconsistencies in the data. **Data integration** merges data from multiple sources into a coherent data store, such as a data warehouse. **Data reduction** can reduce the data size by aggregating, eliminating redundant features, or clustering, for instance. **Data transformations, such as normalization**, may be applied, where data are scaled to fall within a smaller range like 0.0 to 1.0. This can improve the accuracy and efficiency of mining algorithms involving distance measurements*

# Major Tasks in Data Preprocessing



There are 4 major tasks in data preprocessing – Data cleaning, Data integration, Data reduction, and Data transformation.

# Major Tasks in Data Preprocessing





# Data Cleaning

- Real-world data tend to be incomplete, noisy, and inconsistent. Data cleaning (or data cleansing) routines attempt to fill in missing values, smooth out noise while identifying outliers, and correct inconsistencies in the data.
- Data cleaning is the process of removing incorrect data, incomplete data, and inaccurate data from the datasets, and it also replaces the missing values.
- Topics under Data Cleaning are: Handling missing values, Noisy data, Data cleaning as a process

# Data Cleaning

- **Missing or incomplete records:** Missing data sometimes appears as empty cells, values (e.g., NULL or N/A), or a particular character, such as a question mark

age	weight
[50-60)	?
[20-30)	[50-75)
[80-90)	?
[50-60)	?
[50-60)	?
[70-80)	?

# Data Cleaning

- 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.
- It is important to note that, a missing value may not always imply an error. (for example, Null-allow attri. )

# Data Cleaning

- How can you go about filling in the missing values for this attribute? Let's look at the following methods:
- **1. Ignore the tuple:** This is usually done when the class label is missing (assuming the mining task involves classification).
- This method is not very effective, unless the tuple contains several attributes with missing values.
- It is especially poor when the percentage of missing values per attribute varies considerably.
- By ignoring the tuple, we do not make use of the remaining attributes values in the tuple.
- **2. Fill in the missing value manually:** In general, this approach is time consuming and may not be feasible given a large data set with many missing values.

# Data Cleaning

- **3. Use a global constant to fill in the missing value:** Replace all missing attribute values by the same constant, such as a label like “Unknown” or  $-\infty$ .
- If missing values are replaced by, say, “Unknown,” then the mining program may mistakenly think that they form an interesting concept, since they all have a value in common—that of “Unknown.”
- **4. Use a measure of central tendency for the attribute (such as the mean or median) to fill in the missing value:** measures of central tendency, which indicate the “middle” value of a data distribution.
- For normal (symmetric) data distributions, the mean can be used, while skewed data distribution should employ the median.

# Data Cleaning

- **5. Use the attribute mean or median for all samples belonging to the same class as the given tuple:** For example, if classifying customers according to credit risk, we may replace the missing value with the average income value for customers in the same credit risk category as that of the given tuple.
- If the data distribution for a given class is skewed, the median value is a better choice.
- **6. Use the most probable value to fill in the missing value:** This may be determined with regression, inference-based tools using a Bayesian formalism, or decision tree induction.
- For example, using the other customer attributes in your data set, you may construct a decision tree to predict the missing values for income.

# Data Cleaning: Noisy Data

Binning: Binning methods smooth a sorted data value by consulting its “neighborhood,” that is, the values around it. The sorted values are distributed into a number of “buckets,” or bins

- Sorted data for price (in dollars): 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34

- \* Partition into (equi-depth) bins:

  - Bin 1: 4, 8, 9, 15

  - Bin 2: 21, 21, 24, 25

  - Bin 3: 26, 28, 29, 34

- \* Smoothing by bin means:

  - Bin 1: 9, 9, 9, 9  $(4+8+9+15/4) = 9$

  - Bin 2: 23, 23, 23, 23  $(21+21+24+25/4) = 23$

  - Bin 3: 29, 29, 29, 29  $(26+28+29+34/4) = 29$

- \* Smoothing by bin boundaries:

  - Bin 1: 4, 4, 4, 15

  - Bin 2: 21, 21, 25, 25

  - Bin 3: 26, 26, 26, 34

# Data Cleaning

- In this example, the data for price are first sorted and then partitioned into equal-frequency bins of size 3 (i.e., each bin contains three values).
- In smoothing by bin means, each value in a bin is replaced by the mean value of the bin. For example, the mean of the values 4, 8, and 15 in Bin 1 is 9.
- Therefore, each original value in this bin is replaced by the value 9.
- Similarly, smoothing by bin medians can be employed, in which each bin value is replaced by the bin median.
- In smoothing by bin boundaries, the minimum and maximum values in a given bin are identified as the bin boundaries.
- Each bin value is then replaced by the closest boundary value.

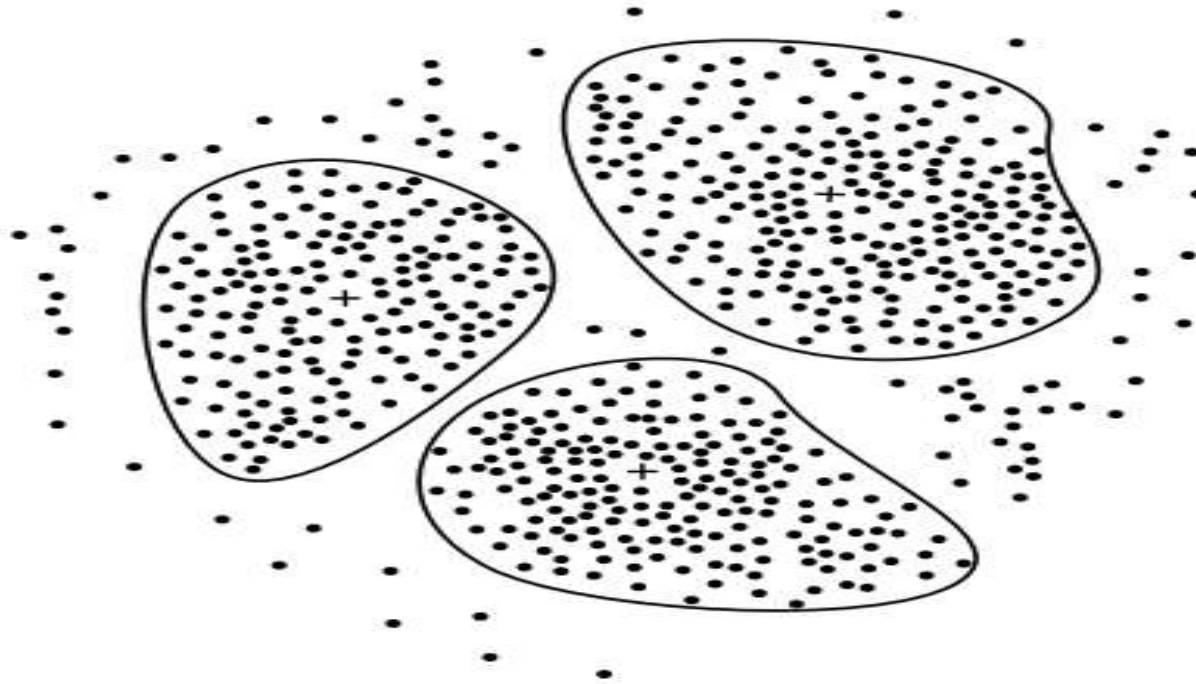


# Data Cleaning

- **Regression:** Data smoothing can also be done by conforming data values to a function, a technique known as regression.
- Linear regression involves finding the “best” line to fit two attributes (or variables), so that one attribute can be used to predict the other.
- Multiple linear regression is an extension of linear regression, where more than two attributes are involved and the data are fit to a multidimensional surface.

# Data Cleaning

- **Outlier analysis:** Outliers may be detected by clustering, for example, where similar values are organized into groups, or “clusters.” Intuitively, values that fall outside of the set of clusters may be considered outliers



A 2-D plot of customer data with respect to customer locations in a city, showing three data clusters. Each cluster centroid is marked with a “+”, representing the average point in space for that cluster. Outliers may be detected as values that fall outside of the sets of clusters.

# Data Cleaning as a Process

- *Missing values, noise, and inconsistencies contribute to inaccurate data. So far, we have looked at techniques for handling missing data and for smoothing data. “But data cleaning is a big job. What about data cleaning as a process? How exactly does one proceed in tackling this task? Are there any tools out there to help?”*
- The first step in data cleaning as a process is discrepancy (inconsistency) detection.
- Discrepancies can be caused by several factors, including poorly designed data entry forms that have many optional fields, human error in data entry, deliberate errors (e.g., respondents not wanting to disclose information about themselves), and data decay (e.g., outdated addresses).
- Discrepancies may also arise from inconsistent data representations and the inconsistent use of codes.
- Errors in instrumentation devices that record data, and system errors, are another source of discrepancies.
- Errors can also occur when the data are (inadequately) used for purposes other than originally intended and there may also be inconsistencies due to data integration (e.g., where a given attribute can have different names in different databases).

# Data Cleaning as a Process

- *“So, how to proceed with discrepancy detection?” As a starting point, use any knowledge you may already have regarding properties of the data. Such knowledge or “data about data” is referred to as metadata.*
- For example, what are the data type and domain of each attribute? What are the acceptable values for each attribute? The basic statistical data descriptions those are useful here to grasp data trends and identify anomalies.
- For example, find the mean, median, and mode values.
- Are the data symmetric or skewed? What is the range of values? Do all values fall within the expected range? What is the standard deviation of each attribute?
- Values that are more than two standard deviations away from the mean for a given attribute may be flagged as potential outliers.
- Are there any known dependencies between attributes?

# Data Cleaning as a Process

- There are a number of different commercial tools that can aid in the step of discrepancy detection.
- Data scrubbing tools use simple domain knowledge (e.g., knowledge of postal addresses, and spell-checking) to detect errors and make corrections in the data.
- These tools rely on parsing and fuzzy matching techniques when cleaning data from multiple sources.
- Data auditing tools find discrepancies by analyzing the data to discover rules and relationships, and detecting data that violate such conditions.
- They are variants of data mining tools. For example, they may employ statistical analysis to find correlations, or clustering to identify outliers.

# Data Cleaning as a Process

- ETL (extraction/transformation/loading) tools allow users to specify transforms through a graphical user interface (GUI).
- These tools typically support only a restricted set of transforms so that, often, we may also choose to write custom scripts for this step of the data cleaning process.
- The two-step process of discrepancy detection and data transformation (to correct discrepancies) iterates.
- This process, however, is error-prone and time consuming and some transformations may introduce more discrepancies while some nested discrepancies may only be detected after others have been fixed.

# Data Cleaning as a Process

- New approaches to data cleaning emphasize increased interactivity.
- Potter's Wheel, for example, is a publicly available data cleaning tool that integrates discrepancy detection and transformation.
- Users gradually build a series of transformations by composing and debugging individual transformations, one step at a time, on a spreadsheet-like interface.
- The transformations can be specified graphically or by providing examples.
- Results are shown immediately on the records that are visible on the screen.
- The user can choose to undo the transformations, so that transformations that introduced additional errors can be “erased.”
- The tool performs discrepancy checking automatically in the background on the latest transformed view of the data.
- Users can gradually develop and refine transformations as discrepancies are found, leading to more effective and efficient data cleaning.
- Another approach to increased interactivity in data cleaning is the development of declarative languages for the specification of data transformation operators.
- Such work focuses on defining powerful extensions to SQL and algorithms that enable users to express data cleaning specifications efficiently.



# Exercise

- Suppose a group of 12 sales price records has been sorted as follows:  
5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215 Partition them into three bins solve it by each of the following methods:
  - (a) equi-depth partitioning
  - (b) Smoothing by bin boundaries
- Solution:
  - (a) equi-depth partitioning -Bin-1: 5, 10, 11, 13, Bin-2: 15, 35, 50, 55, Bin-03: 72, 92, 204, 215
  - (b) Smoothing by bin boundaries-Smoothing by bin boundaries: Bin-1: 5,13,13,13 , Bin-2: 15,15,55,55, Bin-3:72,72,215,215



# Data Integration

- The merging of data from multiple data stores.
- Careful integration can help reduce and avoid redundancies and inconsistencies in the resulting data set.
- This can help improve the accuracy and speed of the subsequent mining process.
- *The semantic heterogeneity and structure of data pose great challenges in data integration. How can we match schema and objects from different sources?*

# Data Integration

Data integration techniques:

- Schema matching
- Instance conflict resolution
- Source selection
- Result merging
- Quality composition

# Data Integration

- Combines data from multiple sources into a coherent store
- Careful integration can help reduce & avoid redundancies and inconsistencies
- This helps to improve accuracy & speed of subsequent data mining
- Heterogeneity & structure of data pose great challenges
- Issues that need to be addressed:
  1. How to match schema & objects from different sources? (**Entity identification problem**)
  2. Are any attributes correlated?
  3. Tuple duplication
  4. Detection & resolution of data value conflicts

# Data Integration

Data integration techniques:

- Schema matching
  - Instance conflict resolution
  - Source selection
  - Result merging
  - Quality composition
- 
- Entity Identification Problem, Redundancy and Correlation Analysis, Tuple Duplication, Data Value Conflict Detection and Resolution

# Data Integration

- **Data integration:**
  - Combines data from multiple sources into a coherent store
  - Careful integration can help reduce & avoid redundancies and inconsistencies
  - This helps to improve accuracy & speed of subsequent data mining
  - Heterogeneity & structure of data pose great challenges
  - Issues that need to be addressed:
    1. How to match schema & objects from different sources? (**Entity identification problem**)
    2. Are any attributes correlated?
    3. Tuple duplication
    4. Detection & resolution of data value conflicts

# Issues during Data Integration(Contd..)

## 1. Schema integration:

Integrate metadata from different sources

- e.g. customer\_id in one database and cust\_number in another

Entity identification problem: Identify real world entities from multiple data sources,

- e.g., Bill Clinton = William Clinton

## 2. Redundancy:

- An attribute may be redundant if it can be derived or obtaining from another attribute or set of attribute.
- Inconsistencies in attribute can also cause redundancies in the resulting data set.
- Some redundancies can be detected by correlation analysis.

## 3. 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

# Data Integration(Contd..)

- **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
- **Redundancy:**
  - Inconsistencies in attribute or dimensions naming can cause redundancy
- **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

# Data Integration

## Entity identification problem

- How can equivalent real-world entities from multiple data sources be matched up?
- Schema integration: e.g., A.cust-id  $\equiv$  B.cust-number
- Integrate using **metadata**(meaning, data type, range of values,etc.) from different sources
- Ensure that any attribute functional dependencies & referential constraints in source system match those in target system



# Handling Redundancy in Data Integration

## Redundancy & Correlation Analysis:

- 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

# Detection of Data Redundancy-Correlation

- Redundancies can be detected using following methods:
  - $\chi^2$ Test (Used for nominal Data or categorical or qualitative data)
  - Correlation coefficient and covariance (Used for numeric Data or quantitative data)
- The term "correlation" refers to a mutual relationship or association between quantities.
- In almost any business, it is useful to express one quantity in terms of its relationship with others.
  - For example, sales might increase when the marketing department spends more on TV advertisements,
  - Customer's average purchase amount on an e-commerce website might depend on a number of factors related to that customer.
- Often, correlation is the first step to understanding these relationships and subsequently building better business and statistical models.

# Why is correlation a useful metric?

- Correlation can help in **predicting** one quantity from another
- Correlation can (but often does not, as we will see in some examples below) **indicate the presence of a causal relationship**
- Correlation is used as a **basic quantity** and foundation for many other modeling techniques
- More formally, correlation is a **statistical measure** that describes the association between **random variables**.
- There are several methods for calculating the correlation coefficient, each measuring different types of strength of association

# Correlation Analysis (Numeric data )

- Evaluate correlation between 2 attributes, A & B, by computing **correlation coefficient(Pearson's product moment coefficient)**

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

Where  $n$  is the number of tuples,  $a_i$  and  $b_i$  are the respective values of  $A$  and  $B$  in tuple  $i$ ,

$\bar{A}$  and  $\bar{B}$  are the respective mean values of  $A$  and  $B$ ,

$\sigma_A$  and  $\sigma_B$  are the respective standard deviations of  $A$  and  $B$

$\Sigma(a_i b_i)$  is the sum of the  $AB$  cross-product (i.e., for each tuple, the value for  $A$  is multiplied by the value for  $B$  in that tuple)

# Correlation- Pearson Correlation Coefficient

- Pearson is the most widely used correlation coefficient.
- Pearson correlation measures the linear association between continuous variables.
- Pearson correlation coefficient

The **covariance** between  $A$  and  $B$  is defined as

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

$$E(A) = \bar{A} = \frac{\sum_{i=1}^n a_i}{n}$$

$$E(B) = \bar{B} = \frac{\sum_{i=1}^n b_i}{n}.$$

- The Correlations coefficient is a statistic and it can range between +1 and -1

# Correlation Analysis (Numeric data )

- Note that  $-1 \leq r_{A,B} \leq +1$
- If resulting value is greater than 0, then  $A$  and  $B$  are *positively correlated*, meaning that the values of  $A$  increase as the values of  $B$  increase
- Higher the value, the stronger the correlation
- Higher value may indicate that  $A$  (or  $B$ ) may be removed as a redundancy
- If the resulting value is equal to 0, then  $A$  and  $B$  are *independent* and there is no
- correlation between them
- If the resulting value is less than 0, then  $A$  and  $B$  are *negatively correlated*, where the values of one attribute increase as the values of the other attribute decrease

# Covariance Analysis(Numeric Data)

- Correlation and covariance are two similar measures for assessing how much two attributes change together
- Consider two numeric attributes  $A$  and  $B$ , and a set of  $n$  observations

$$\{(a_1, b_1), \dots, (a_n, b_n)\}$$

- Mean values of  $A$  and  $B$ , respectively, are also known as the **expected values** on  $A$  and  $B$ , that is

$$E(A) = \bar{A} = \frac{\sum_{i=1}^n a_i}{n}$$

$$E(B) = \bar{B} = \frac{\sum_{i=1}^n b_i}{n}.$$

- **covariance** between  $A$  and  $B$  is defined as  $Cov(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}.$

# Covariance Analysis(Numeric Data)

$$r_{A,B} = \frac{Cov(A,B)}{\sigma_A \sigma_B}$$

- Also, for simplified calculations

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

- For two attributes  $A$  and  $B$  that tend to change together,
  - if  $A$  is larger than  $\bar{A}$ , then  $B$  is likely to be larger than  $\bar{B}$ . Hence, the covariance between  $A$  and  $B$  is **positive**
  - if one of the attributes tends to be above its expected value when the other attribute is below its expected value, then the covariance of  $A$  and  $B$  is **negative**
  - covariance value is 0 means attributes are independent



# Correlation

- **Positive covariance:** If  $\text{Cov}(A, B) > 0$ , then A and B both tend to be larger than **their expected values**
- **Negative covariance:** If  $\text{Cov}(A, B) < 0$  then if A is larger than its expected value, B is likely to be smaller than its **expected value**
- **Independence:**  $\text{Cov}(A, B) = 0$  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

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

# Covariance Analysis(Numeric Data)

Stock Prices for *AllElectronics* and *HighTech*

<i>Time point</i>	<i>AllElectronics</i>	<i>HighTech</i>
t1	6	20
t2	5	10
t3	4	14
t4	3	5
t5	2	5

- Table shows stock prices of two companies at five time points. If the stocks are affected by same industry trends, determine whether their prices rise or fall together?

# Covariance Analysis(Numeric Data)

$$E(AllElectronics) = \frac{6 + 5 + 4 + 3 + 2}{5} = \frac{20}{5} = \$4$$

and

$$E(HighTech) = \frac{20 + 10 + 14 + 5 + 5}{5} = \frac{54}{5} = \$10.80.$$

we compute

$$\begin{aligned} Cov(AllElectronics, HighTech) &= \frac{6 \times 20 + 5 \times 10 + 4 \times 14 + 3 \times 5 + 2 \times 5}{5} - 4 \times 10.80 \\ &= 50.2 - 43.2 = 7. \end{aligned}$$

Therefore, given the positive covariance we can say that stock prices for both companies rise together. ■

# Correlation: example

An example of stock prices observed at five time points for *AllElectronics* and *HighTech*, a high-tech company. If the stocks are affected by the same industry trends, will their prices rise or fall together?

Time Point	All Electronics	High Tech
t1	2	5
t2	3	8
t3	5	10
t4	4	11
t5	6	14

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  $Cov(A, B) > 0$ .

# Data Value Conflict Detection and Resolution

- For the same real-world entity, attribute values from different sources may differ
  - Eg. Prices of rooms in different cities may involve different currencies
- Attributes may also differ on the abstraction level, where an attribute in one system is recorded at, say, a lower abstraction level than the “same” attribute in another.
  - Eg. total sales in one database may refer to one branch of All\_Electronics, while an attribute of the same name in another database may refer to the total sales for All\_Electronics stores in a given region.
- To resolve, data values have to be converted into consistent form

# Data Transformation

- **Smoothing**: remove noise from data (binning, clustering, regression)
- **Aggregation**: summarization, data cube construction
- **Generalization**: concept hierarchy climbing
- **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
  - E.g. add attribute *area* based on the attributes *height* and *width*

# Data Transformation: Normalization

- Min-max normalization

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new\_max}_A - \text{new\_min}_A) + \text{new\_min}_A$$

- Z-score normalization

$$v' = \frac{v - \text{mean}_A}{\text{stand\_dev}_A}$$

- Normalization by decimal scaling

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

# Data Transformation: Min Max-Normalization

- Min-max normalization: to  $[new\_minA, new\_maxA]$ 
  - Performs a linear transformation on the original data.
  - Suppose that  $min_a$  and  $max_a$  are the minimum and maximum values of an attribute,  $a$ .
  - Min-max normalization maps a value,  $v_i$ , of  $a$  to the range  $[new\_min_a, new\_max_a]$  by computing

- Let **income** range \$12,000 to \$98,000 normalized to  $[0.0, 1.0]$ .

- Then \$73,000 is mapped to

$$\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$$

- Eg. Suppose that the minimum and maximum values for the attribute income are \$12,000 and \$98,000, respectively. We would like to map income to the range  $[0.0, 1.0]$ . By min-max normalization, a value of \$73,600 for income is transformed to **0.716**



# Data Transformation: Z-score Normalization

- Z-score normalization (uses mean &  $\sigma$ : standard deviation):
- Values for an attribute, A, are normalized based on the mean (i.e., average) and standard deviation of A
- Eg. Suppose that the mean and standard deviation of the values for the attribute income are \$54,000 and \$16,000, respectively.

$$v'_i = \frac{v_i - \bar{A}}{\sigma_A}$$

- Let  $\bar{A} = 54,000$ ,  $\sigma_A = 16,000$ , for the attribute *income*
- With z-score normalization, a value of \$73,600 for *income* is transformed to:

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

# Data Transformation: Decimal scaling Normalization

- Suppose that the recorded values of A ranges from -986 to 917
- The maximum absolute value of A is 986
- To normalize by decimal scaling, we therefore divide each value by 1000
- So that -986 normalize to -0.986 and
- 917 normalize to -0.917

# Data Warehouse and Data Mining

- Introduction to KDD
- Data Preprocessing: An Overview
  - Data Quality
  - Major Tasks in Data Preprocessing
- Data Cleaning
- Data Integration
- **Data Reduction**
- Data Transformation and Data Discretization
- Summary

# Exercise

- Define Normalization.
- What is the value range of min-max. Use min-max normalization to normalize the following group of data: 8,10,15,20.

- Solution:

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new\_max}_A - \text{new\_min}_A) + \text{new\_min}_A$$

Marks	Marks after Min-Max normalization
8	0
10	0.16
15	0.58
20	1

# Data Reduction Strategies

- Data is too big to work with
- Data reduction
  - Obtain a reduced representation of the data set that is much smaller in volume but yet produce 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 — remove unimportant attributes
  - Aggregation and clustering
  - Sampling

# Data Reduction Strategies

- Obtains a reduced representation of the data set that is much smaller in volume but yet produces the same (or almost the same) analytical results
- Data reduction strategies
  - Data cube aggregation
  - Dimensionality reduction
  - Data compression
  - Numerosity reduction
  - Discretization and concept hierarchy generation

# Data Cube Aggregation

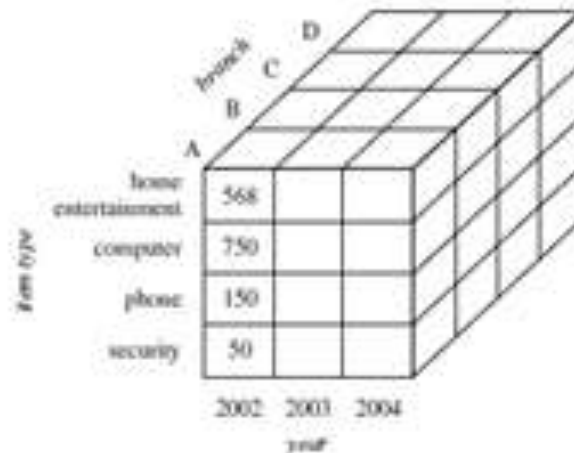
- Multiple levels of aggregation in data cubes
  - Further reduce the size of data to deal with
- Reference appropriate levels
  - Use the smallest representation capable to solve the task
- Queries regarding aggregated information should be answered using data cube, when possible

# Data Cube Aggregation

Year 2002	
Quarter	Sales
Q1	\$224,000
Q2	\$408,000
Q3	\$350,000
Q4	\$586,000

Year	Sales
2002	\$1,568,000
2003	\$2,356,000
2004	\$3,594,000

**Figure 2.13** Sales data for a given branch of *AllElectronics* for the years 2002 to 2004. On the left, the sales are shown per quarter. On the right, the data are aggregated to provide the annual sales.



**Figure 2.14** A data cube for sales at *AllElectronics*.



# Data Reduction Strategies

- 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

- **Problem:** Feature selection (i.e., attribute subset selection):
  - Select a minimum set of features such that the probability distribution of different classes given the values for those features is as close as possible to the original distribution given the values of all features
  - Nice side-effect: reduces # of attributes in the discovered patterns (which are now easier to understand)
- **Solution:** Heuristic methods (due to exponential # of choices) usually greedy:
  - step-wise forward selection
  - step-wise backward elimination
  - combining forward selection and backward elimination
  - decision-tree induction
  - E.g. Decision Tree

# Dimensionality reduction :Attribute Subset Selection

- Reduces the data set size by removing irrelevant or redundant attributes (or dimensions)
- Goal of attribute subset selection is to find a minimum set of attributes
- Improves speed of mining as dataset size is reduced
- Mining on a reduced data set also makes the discovered pattern easier to understand
- Redundant attributes
  - Duplicate information contained in one or more 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' telephone number is often irrelevant to the task of predicting students' CGPA

# Attribute Subset Selection

Best and worst attributes are determined using **tests of statistical significance**

- A test of significance is a formal procedure for comparing observed data with a claim (also called a hypothesis), the truth of which is being assessed
- Claim is a statement about a parameter
- Results of a significance test are expressed in terms of a probability that measures how well the data and the claim agree

# Heuristic (Greedy) methods for attribute subset selection

## 1. Stepwise Forward Selection:

- Starts with an empty set of attributes as the reduced set
- Best of the relevant attributes is determined and added to the reduced set
- In each iteration, best of remaining attributes is added to the set

## 2. Stepwise Backward Elimination:

- Here all the attributes are considered in the initial set of attributes
- In each iteration, worst attribute remaining in the set is removed

## 3. Combination of Forward Selection and Backward Elimination:

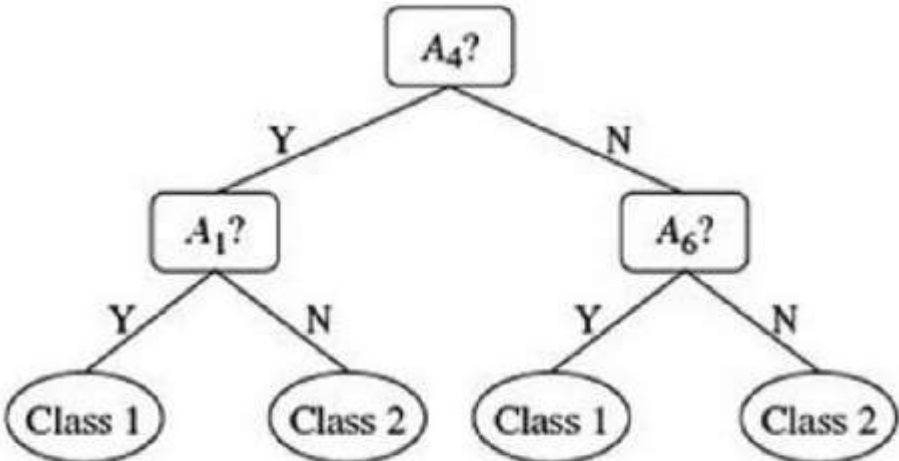
- Stepwise forward selection and backward elimination are combined
- At each step, the procedure selects the best attribute and removes the worst from among the remaining attributes

# Heuristic (Greedy) methods for attribute subset selection(cont)

## 4. Decision Tree Induction:

- This approach uses decision tree for attribute selection.
- It constructs a flow chart like structure having nodes denoting a test on an attribute.
- Each branch corresponds to the outcome of test and leaf nodes is a class prediction.
- The attribute that is not the part of tree is considered irrelevant and hence discarded

# Example of Decision Tree Induction

Forward selection	Backward elimination	Decision tree induction
<p>Initial attribute set: <math>\{A_1, A_2, A_3, A_4, A_5, A_6\}</math></p> <p>Initial reduced set: <math>\{\}</math> <math>\Rightarrow \{A_1\}</math> <math>\Rightarrow \{A_1, A_4\}</math> <math>\Rightarrow</math> Reduced attribute set: <math>\{A_1, A_4, A_6\}</math></p>	<p>Initial attribute set: <math>\{A_1, A_2, A_3, A_4, A_5, A_6\}</math></p> <p><math>\Rightarrow \{A_1, A_3, A_4, A_5, A_6\}</math> <math>\Rightarrow \{A_1, A_4, A_5, A_6\}</math> <math>\Rightarrow</math> Reduced attribute set: <math>\{A_1, A_4, A_6\}</math></p>	<p>Initial attribute set: <math>\{A_1, A_2, A_3, A_4, A_5, A_6\}</math></p>  <pre>graph TD; A4["A4?"] -- Y --&gt; A1["A1?"]; A4 -- N --&gt; A6["A6?"]; A1 -- Y --&gt; C1_1((Class 1)); A1 -- N --&gt; C2_1((Class 2)); A6 -- Y --&gt; C1_2((Class 1)); A6 -- N --&gt; C2_2((Class 2));</pre> <p><math>\Rightarrow</math> Reduced attribute set: <math>\{A_1, A_4, A_6\}</math></p>

# Data Reduction 2: Numerosity Reduction

- **Parametric methods**

- Assume the data fits some model, estimate model parameters, store only the parameters, and discard the data (except possible outliers)
- E.g.: 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



# Regression and Log-Linear Models

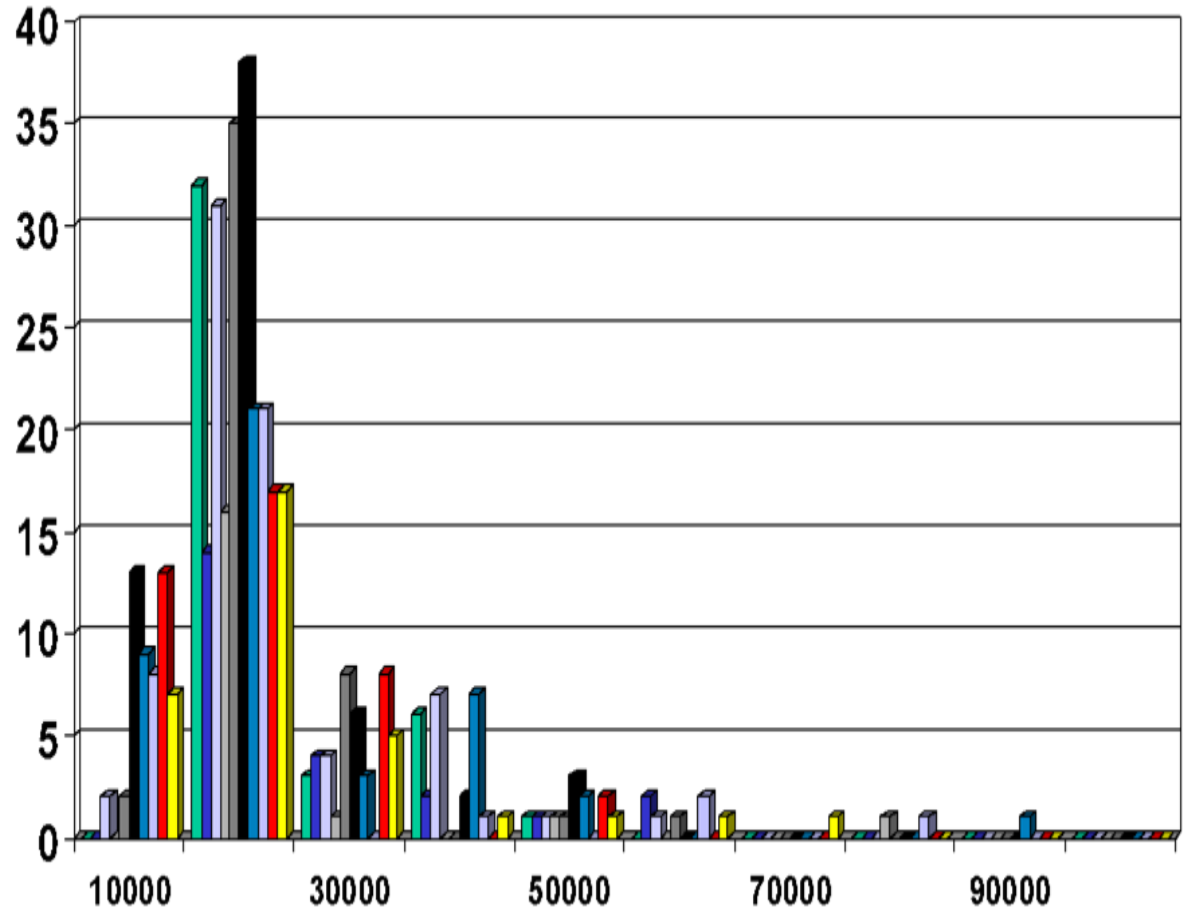
- Linear regression: Data are modeled to fit a straight line:  $Y = \alpha + \beta X$ 
  - Often uses the least-square method to fit the line
- Multiple regression: allows a response variable  $y$  to be modeled as a linear function of multidimensional feature vector (predictor variables)  $Y = b_0 + b_1 X_1 + b_2 X_2$ .
- Log-linear model: approximates discrete multidimensional joint probability distributions  $p(a, b, c, d) = \alpha_{ab} \beta_{ac} \gamma_{ad} \delta_{bcd}$

# Histograms

- Histograms (or frequency histograms) are at least a century old and are widely used.
- Plotting histograms is a graphical method **for summarizing the distribution of a given attribute, X.**
- Height of the bar indicates the frequency (i.e., count) of that X value
- Range of values for X is partitioned into disjoint consecutive subranges.
- Subranges, referred to as **buckets or bins**, are disjoint subsets of the data distribution for X.
- Range of a bucket is known as the width
- Typically, the buckets are of equal width.
- Eg. a price attribute with a value range of \$1 to \$200 can be partitioned into subranges 1 to 20, 21 to 40, 41 to 60, and so on.
- For each subrange, a bar is drawn with a height that represents the total count of items observed within the subrange

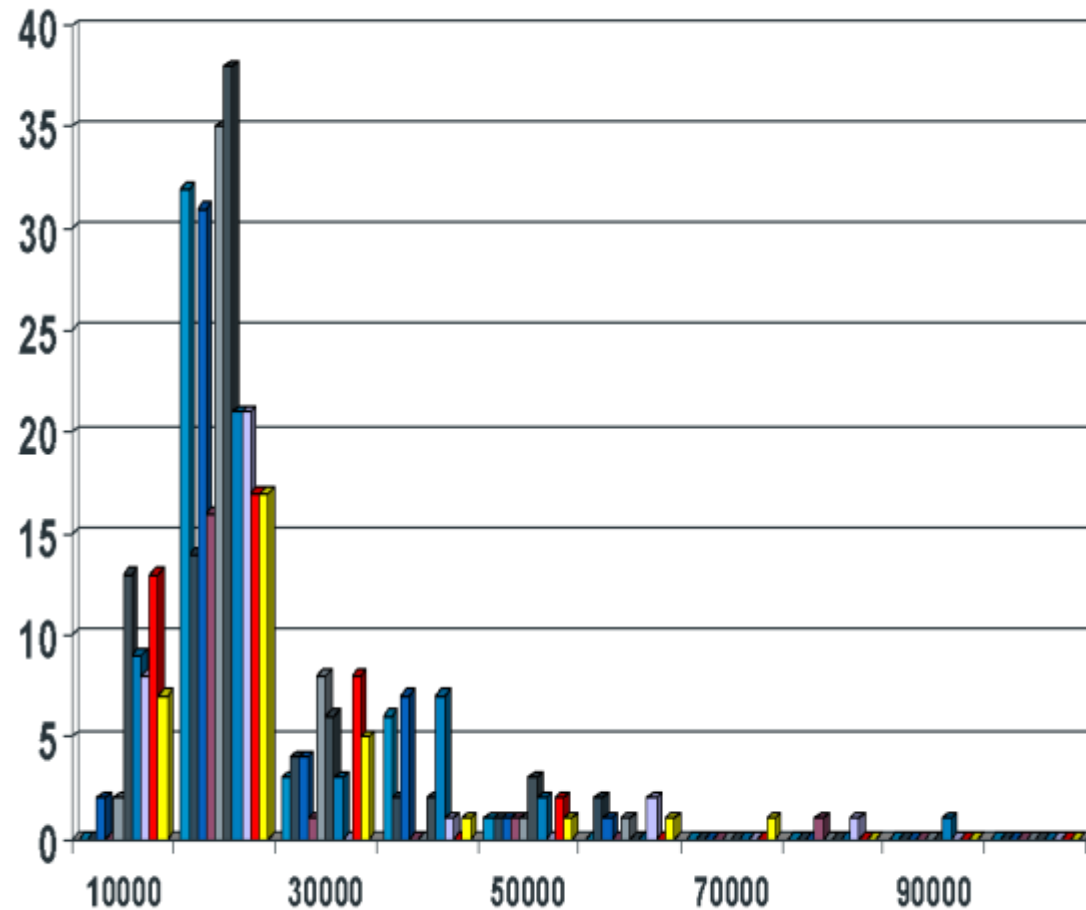
# Histograms

- Approximate data distributions
- Divide data into buckets and store average (sum) for each bucket
- A bucket represents an attribute-value/frequency pair
- Can be constructed optimally in one dimension using dynamic programming
- Related to quantization problems.



# Histogram

- Divide data into buckets and store average (sum) for each bucket
- Partitioning rules:
  - Equal-width: equal bucket range
  - Equal-frequency (or equal-depth)



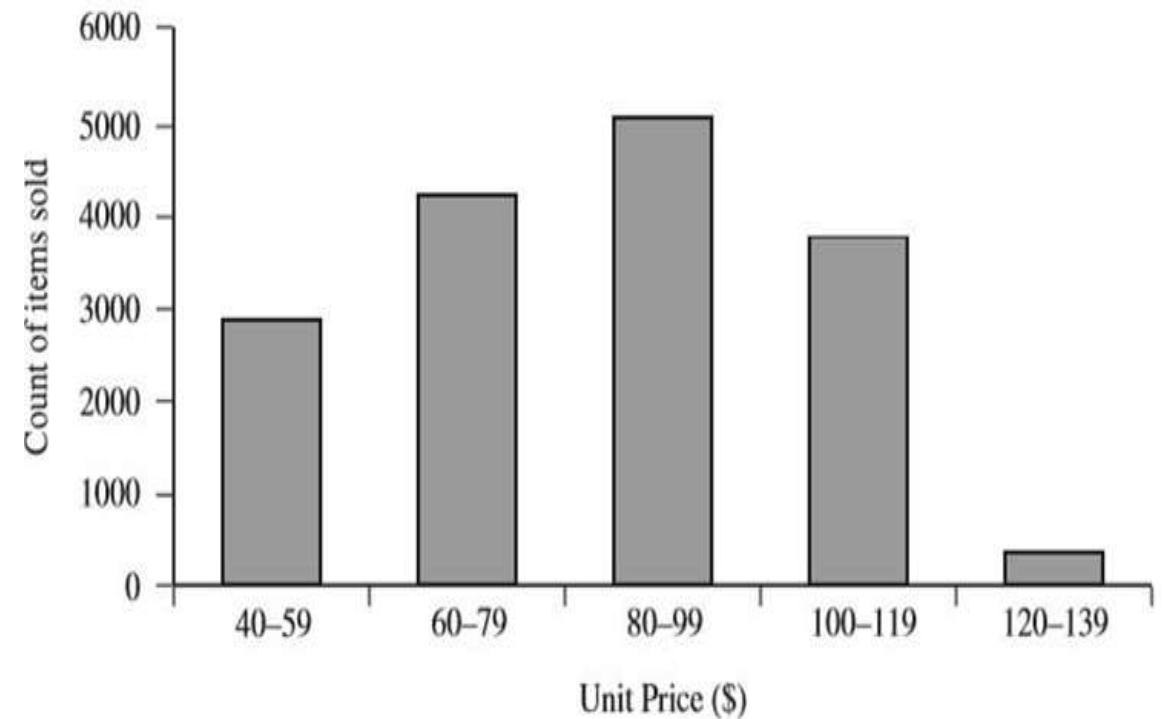
# Histogram Analysis- Explanation & Example

- Plotting histograms, or frequency histograms, is a graphical method for summarizing the distribution of a given attribute.
- A histogram for an attribute  $A$  partitions the data distribution of  $A$  into disjoint subsets, or buckets. Typically, the width of each bucket is uniform.
- Each bucket is represented by a rectangle whose height is equal to the count or relative frequency of the values at the bucket.
- The following data are a list of AllElectronics prices for commonly sold items (rounded to the nearest dollar).

# Histogram Analysis- Explanation & Example

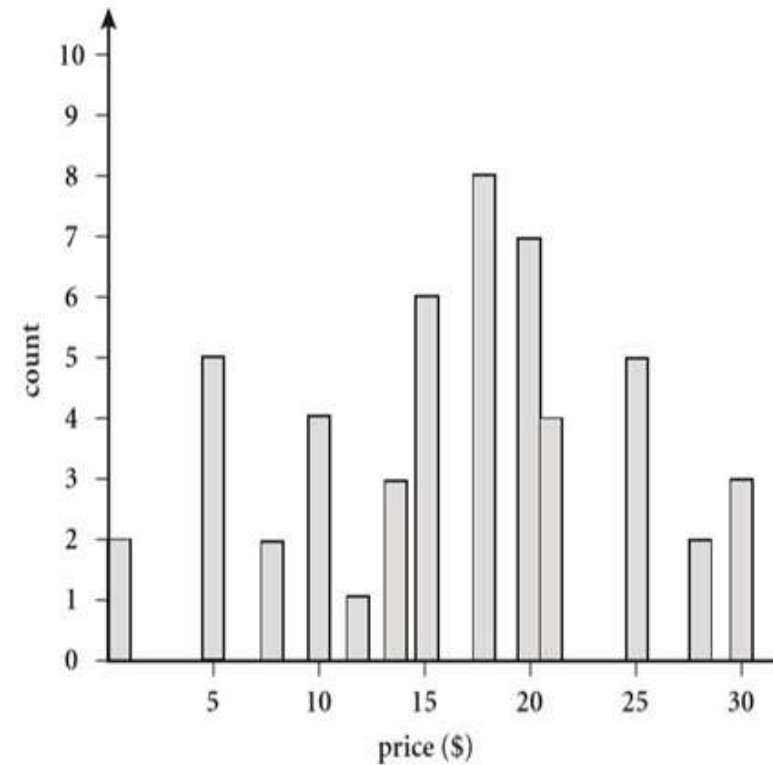
A set of unit price data for items sold at a branch of *AllElectronics*.

Unit price (\$)	Count of items sold
40	275
43	300
47	250
..	..
74	360
75	515
78	540
..	..
115	320
117	270
120	350



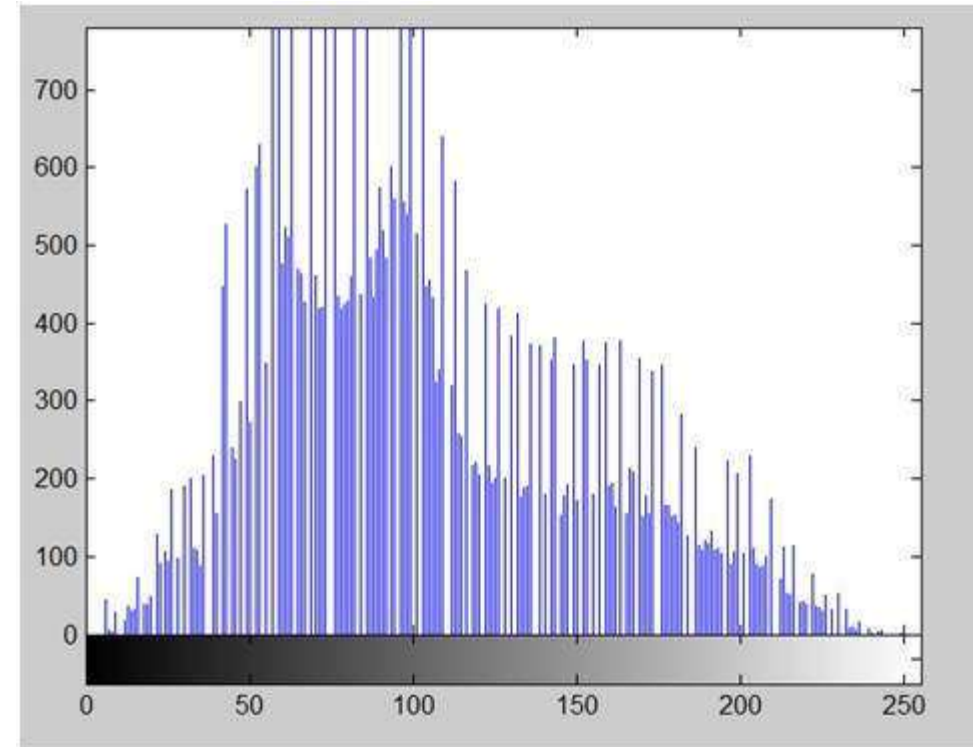
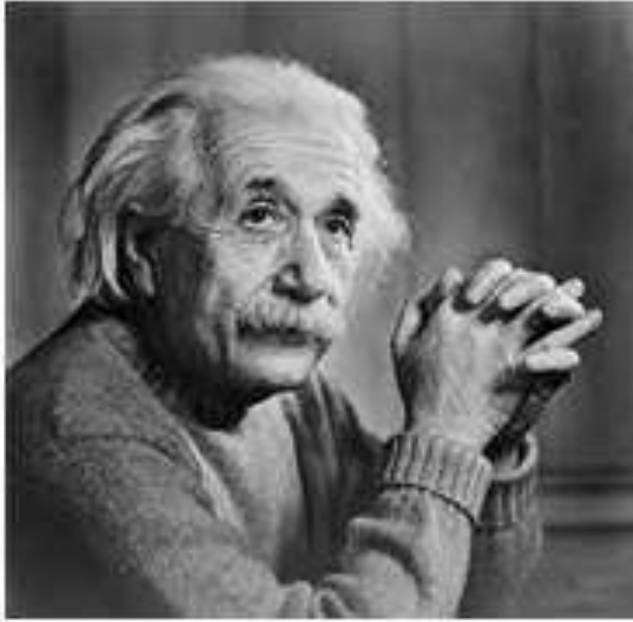
# Histogram Analysis- Explanation & Example

- The following data are a list of prices of commonly sold items at AllElectronics (rounded to the nearest dollar).
- The numbers have been sorted:
- 1, 1, 5, 5, 5, 5, 5, 8, 8, 10, 10, 10, 10, 12, 14, 14, 14, 15, 15, 15, 15, 15, 15, 18, 18, 18, 18, 18, 18, 18, 18, 18, 20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 25, 25, 25, 25, 25, 28, 28, 30, 30, 30.



A histogram for *price* using singleton buckets—each bucket represents one price-value/frequency pair.

# Histogram of an image : Application



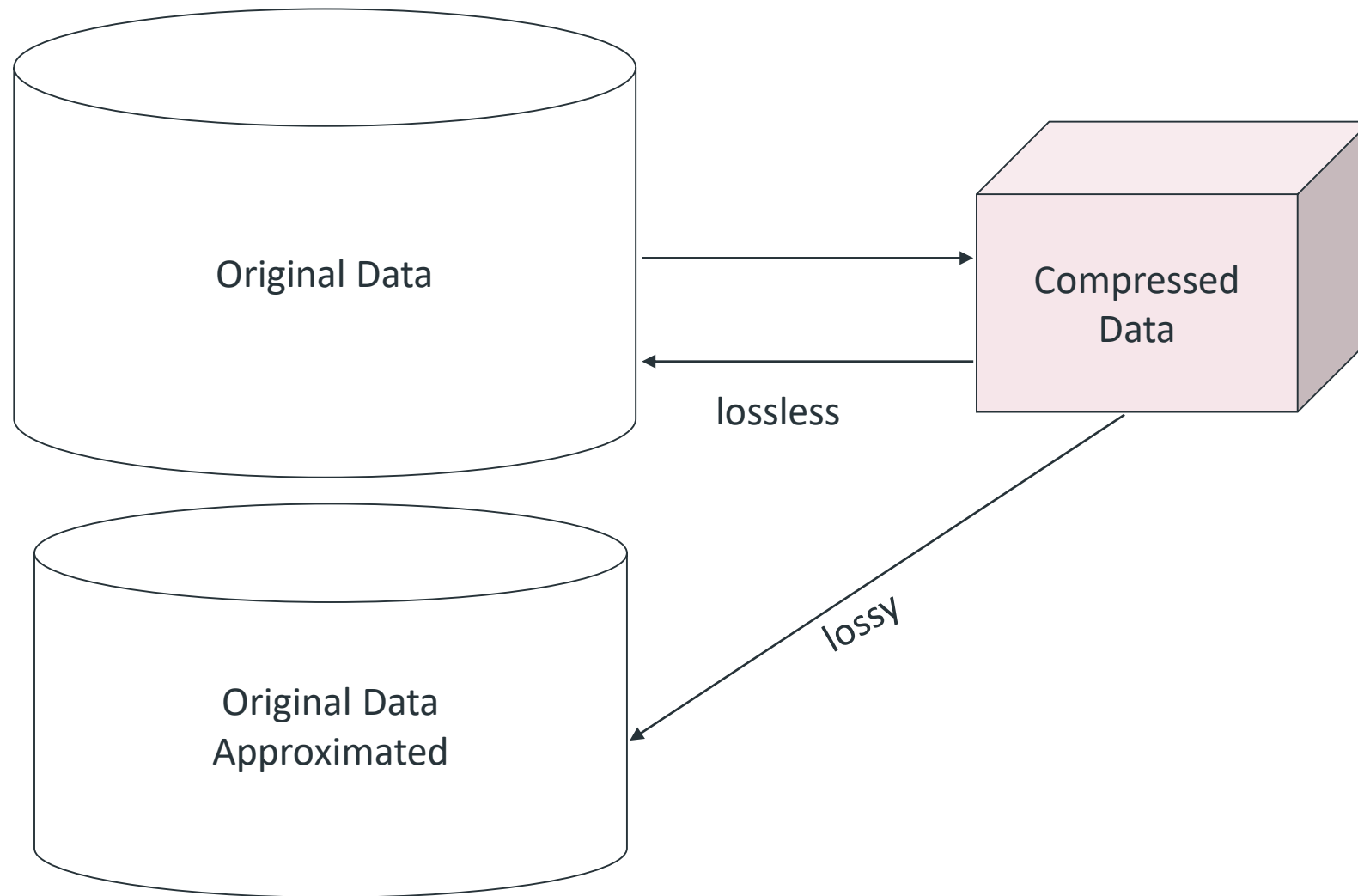
As you can see from the graph, that most of the bars that have high frequency lies in the first half portion which is the darker portion. That means that the image we have got is darker



# Data Compression

- String compression
  - There are extensive theories and well-tuned algorithms
  - Typically **lossless**
  - But only limited manipulation is possible without expansion
- Audio/video, image 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

# Data Compression



# Clustering

- **Partition** data set into clusters, and store cluster representation only
- **Quality** of clusters measured by their **diameter** (max distance between any two objects in the cluster) or centroid distance (avg. distance of each cluster object from its centroid)
- Can be very effective if data is clustered but not if data is “smeared”
- Can have hierarchical clustering (possibly stored in multi-dimensional index tree structures (B+-tree, R-tree, quad-tree, etc))
- There are many choices of clustering definitions and clustering algorithms (further details later)

# 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)

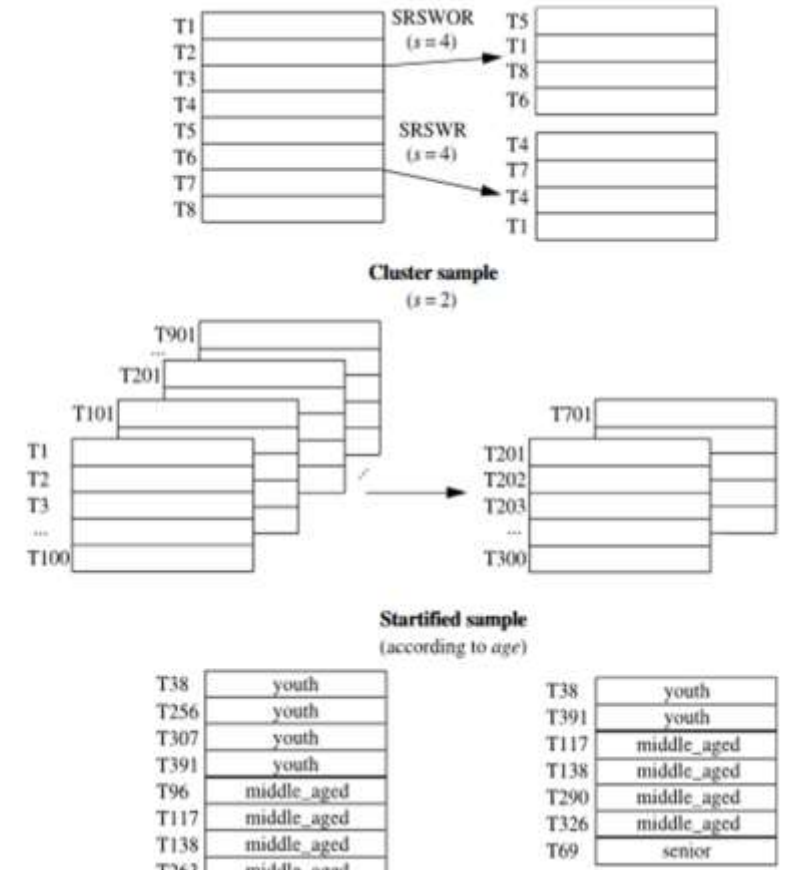
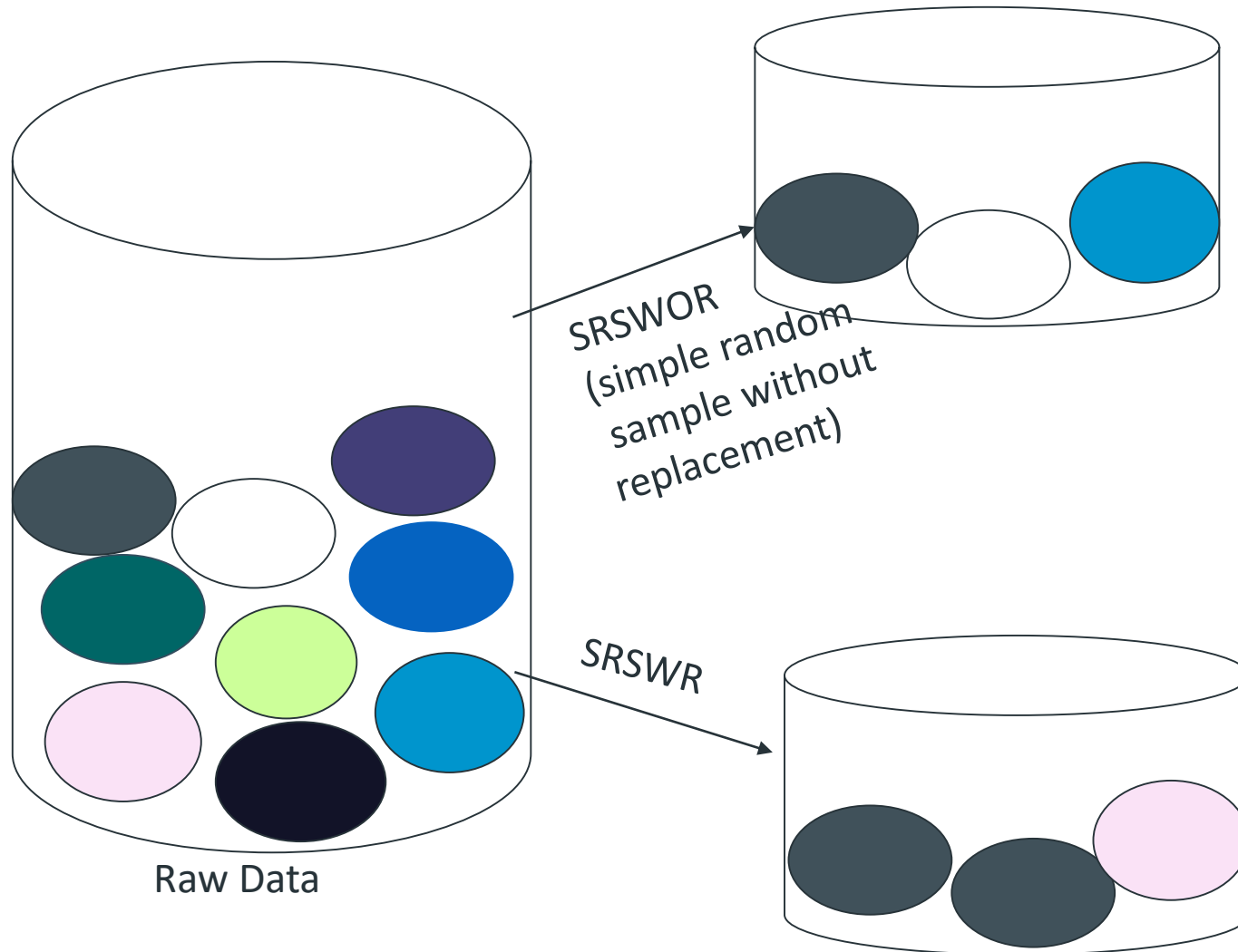
# Sampling

- Allow a mining algorithm to run in complexity that is potentially sub-linear to the size of the data
- Cost of sampling: proportional to the size of the sample, increases linearly with the number of dimensions
- Choose a representative subset of the data
  - Simple random sampling may have very poor performance in the presence of skew
- Develop adaptive sampling methods
  - Stratified sampling:
    - Approximate the percentage of each class (or subpopulation of interest) in the overall database
    - Used in conjunction with skewed data
- Sampling may not reduce database I/Os (page at a time).
- Sampling: natural choice for progressive refinement of a reduced data set.

# Types of Sampling

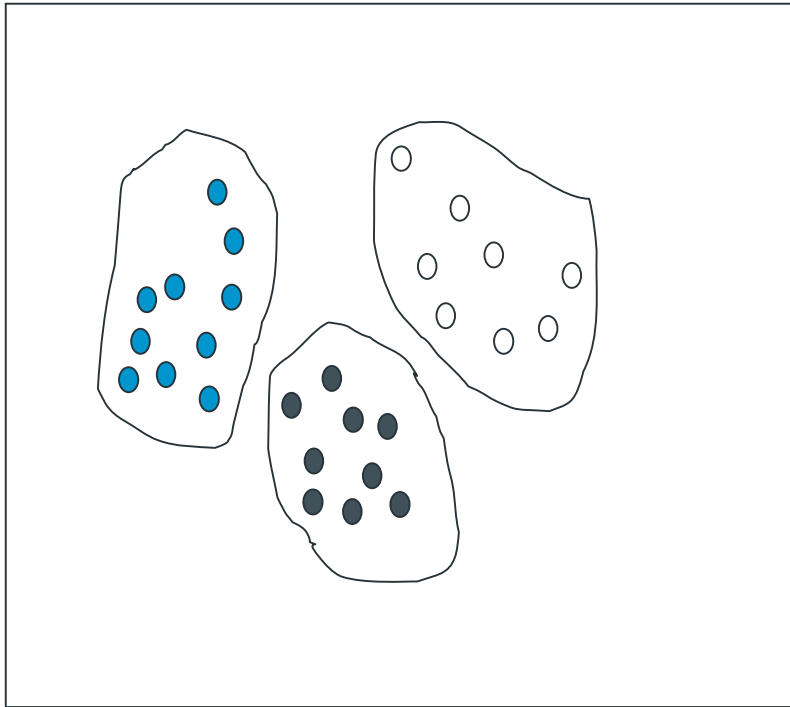
- **Simple random sampling**
  - There is an equal probability of selecting any particular item
- **Sampling without replacement**
  - Once an object is selected, it is removed from the population
- **Sampling with replacement**
  - A selected object is not removed from the population
- **Stratified sampling:**
  - Partition the data set, and draw samples from each partition (proportionally, i.e., approximately the same percentage of the data)
  - Used in conjunction with skewed data

# Types of Sampling

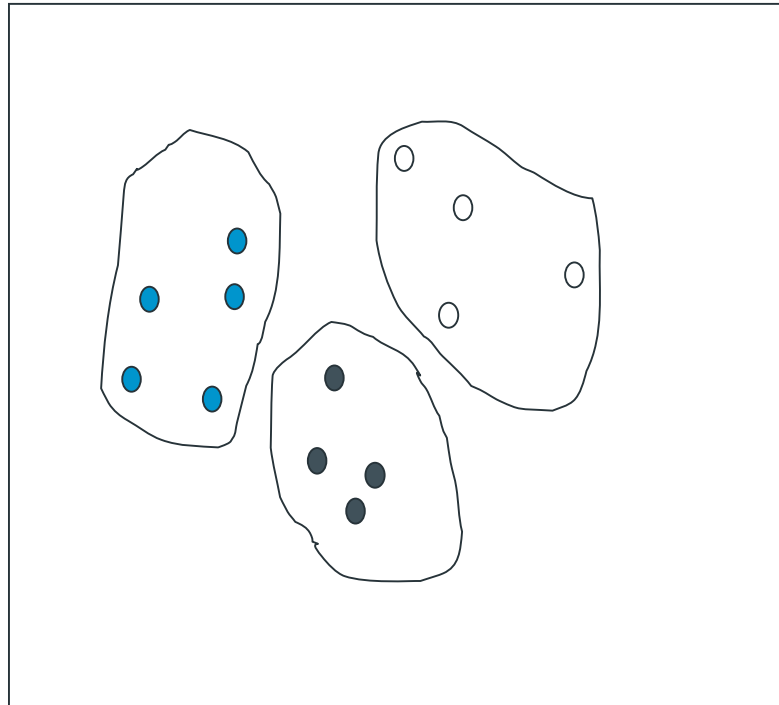


# Sampling

Raw Data



Cluster/Stratified Sample



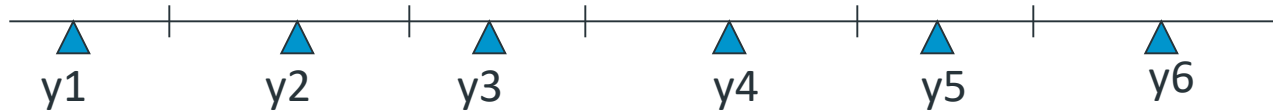


# Discretization

- 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

- **Three types of attributes:**
  - Nominal — values from an unordered set
  - Ordinal — values from an ordered set
  - Continuous — real numbers
- **Discretization/Quantization:**
  - divide the range of a continuous attribute into intervals



- Some classification algorithms only accept categorical attributes.
- Reduce data size by discretization
- Prepare for further analysis

# Discretization and Concept Hierarchies

- Discretization
  - **reduce the number of values** for a given continuous attribute by dividing the range of the attribute into intervals. Interval labels can then be used to replace actual data values.
- Concept Hierarchies
  - reduce the data by collecting and **replacing low level** concepts (such as numeric values for the attribute age) **by higher level concepts** (such as young, middle-aged, or senior).

# Discretization and Concept Hierarchies : Numerical data

- Hierarchical and recursive decomposition using:
  - Binning (data smoothing)
  - Histogram analysis (numerosity reduction)
  - Clustering analysis (numerosity reduction)
- Entropy-based discretization
- Segmentation by natural partitioning

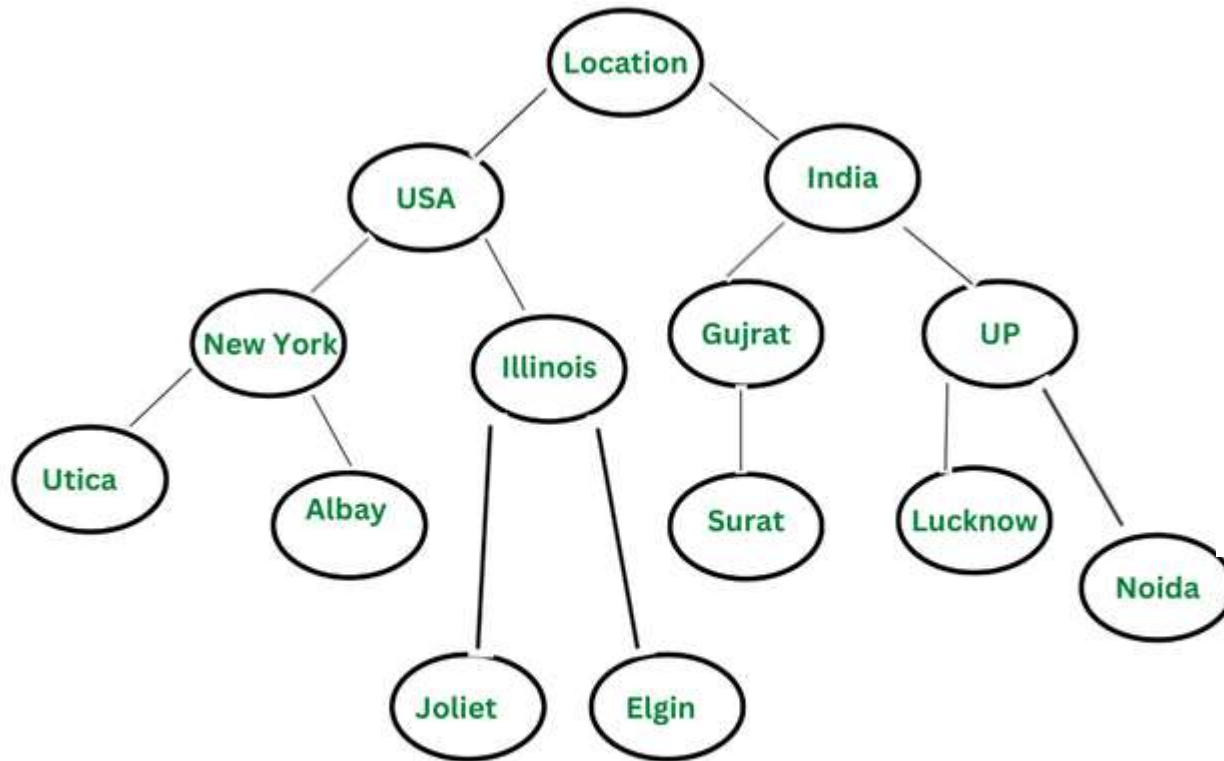
# Data Discretization Methods

- 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.,  $\chi^2$ ) analysis (unsupervised, bottom-up merge)

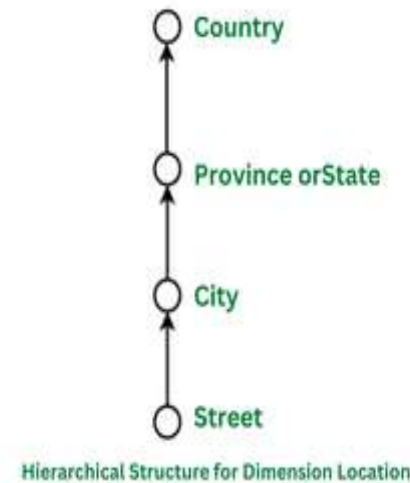
# Concept of Hierarchy in Nominal Data

- In data mining, the concept of a concept hierarchy refers to the organization of data into a tree-like structure, where each level of the hierarchy represents a concept that is more general than the level below it.
- This hierarchical organization of data allows for more efficient and effective data analysis, as well as the ability to drill down to more specific levels of detail when needed.
- The concept of hierarchy is used to organize and classify data in a way that makes it more understandable and easier to analyze.
- The main idea behind the concept of hierarchy is that the same data can have different levels of granularity or levels of detail and that by organizing the data in a hierarchical fashion, it is easier to understand and perform analysis.

# Concept of Hierarchy in Nominal Data



Concept Hierarchy for Dimension Location



The top of the tree consists of the main dimension location and further splits into various sub-nodes. The root node is located, and it further splits into two nodes countries i.e. USA and India. These countries are further then splitted into more sub-nodes, that represent the province states ie. New York, Illinois, Gujarat, UP.

# Need of Concept Hierarchy in Data Mining

- **Improved Data Analysis:** A concept hierarchy can help to organize and simplify data, making it more manageable and easier to analyze. By grouping similar concepts together, a concept hierarchy can help to identify patterns and trends in the data that would otherwise be difficult to spot.
- **Improved Data Visualization and Exploration:** A concept hierarchy can help to improve data visualization and data exploration by organizing data into a tree-like structure, allowing users to easily navigate and understand large and complex data sets.



# Need of Concept Hierarchy in Data Mining

- **Improved Algorithm Performance:** The use of a concept hierarchy can also help to improve the performance of data mining algorithms. By organizing data into a hierarchical structure, algorithms can more easily process and analyze the data, resulting in faster and more accurate results.
- **Data Cleaning and Pre-processing:** A concept hierarchy can also be used in data cleaning and pre-processing, to identify and remove outliers and noise from the data.
- **Domain Knowledge:** A concept hierarchy can also be used to represent the domain knowledge in a more structured way, which can help in a better understanding of the data and the problem domain.

# Applications of Concept Hierarchy

- **Data Warehousing:** Concept hierarchy can be used in data warehousing to organize data from multiple sources into a single, consistent and meaningful structure. This can help to improve the efficiency and effectiveness of data analysis and reporting.
- **Business Intelligence:** Concept hierarchy can be used in business intelligence to organize and analyze data in a way that can inform business decisions. For example, it can be used to analyze customer data to identify patterns and trends that can inform the development of new products or services.
- **Online Retail:** Concept hierarchy can be used in online retail to organize products into categories, subcategories and sub-subcategories, it can help customers to find the products they are looking for more quickly and easily.
- **Healthcare:** Concept hierarchy can be used in healthcare to organize patient data, for example, to group patients by diagnosis or treatment plan, it can help to identify patterns and trends that can inform the development of new treatments or improve the effectiveness of existing treatments.
- **Natural Language Processing:** Concept hierarchy can be used in natural language processing to organize and analyze text data, for example, to identify topics and themes in a text, it can help to extract useful information from unstructured data.
- **Fraud Detection:** Concept hierarchy can be used in fraud detection to organize and analyze financial data, for example, to identify patterns and trends that can indicate fraudulent activity.

# 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

# Reference

- **Data Mining: Concepts and Techniques**, Jiawei Han, Micheline Kamber, and Jian Pei, 3<sup>rd</sup> edition