



# Why Data Preprocessing?

- Data in the real world is dirty
    - **incomplete**: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
      - e.g., occupation=""
    - **noisy**: containing errors or outliers
      - e.g., Salary="-10"
    - **inconsistent**: containing discrepancies in codes or names
      - e.g., Age="42" Birthday="03/07/1997"
      - e.g., Was rating "1,2,3", now rating "A, B, C"
      - e.g., discrepancy between duplicate records
-



# What is Data?

- 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				
Tid	Refund	Marital Status	Taxable Income	Cheat
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



# Types of Attributes

- There are different types of attributes
    - **Nominal**
      - Examples: ID numbers, eye color, zip codes
    - **Ordinal**
      - Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height in {tall, medium, short}
    - **Interval**
      - Examples: calendar dates, temperatures in Celsius or
    - **Ratio**
      - Examples: temperature, length, time, counts
-



# Discrete and Continuous Attributes

- Discrete Attribute
    - Has only a finite or countably 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.
-

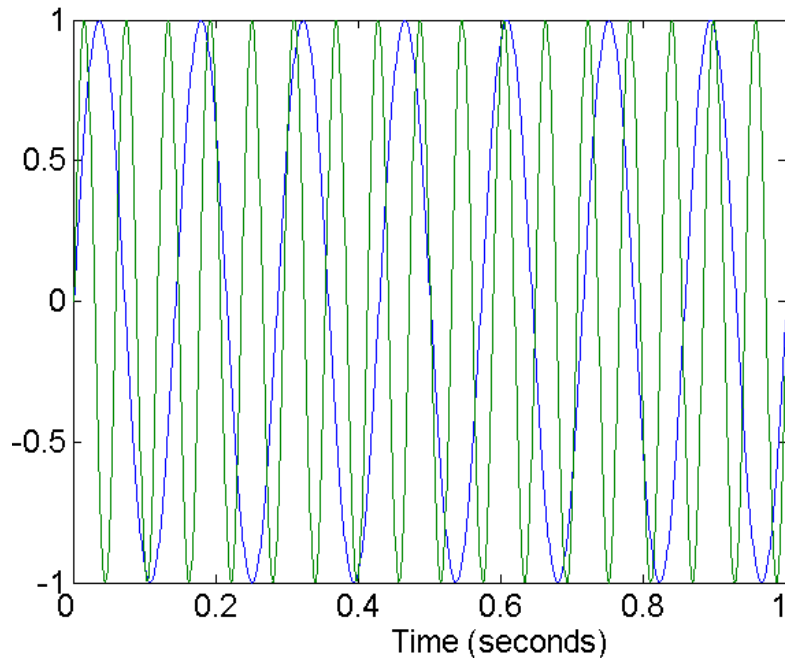


# Data Quality

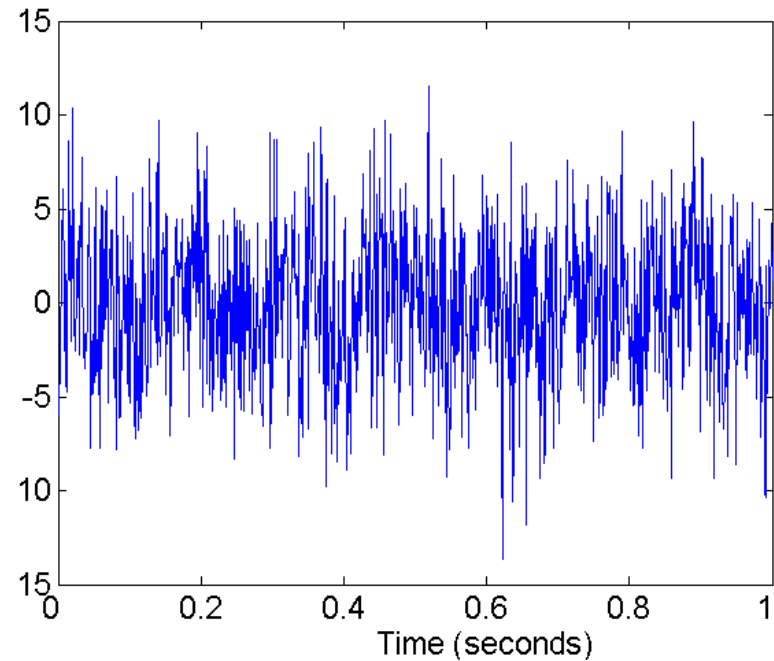
- What kinds of data quality problems?
  - How can we detect problems with the data?
  - What can we do about these problems?
  
  - Examples of data quality problems:
    - Noise and outliers
    - missing values
    - duplicate data
-



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



**Two Sine Waves**

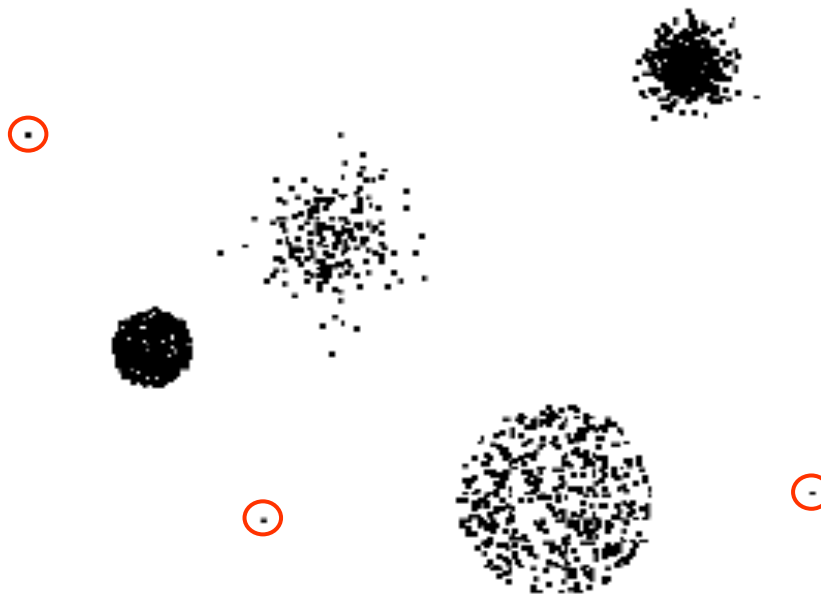


**Two Sine Waves + Noise**



# 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
    - 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)
  - Handling missing values
    - Eliminate Data Objects
    - Estimate Missing Values
    - Ignore the Missing Value During Analysis
    - Replace with all possible values (weighted by their probabilities)
-





# 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
  - Examples:
    - Same person with multiple email addresses
  - Data cleaning
    - Process of dealing with duplicate data issues
-



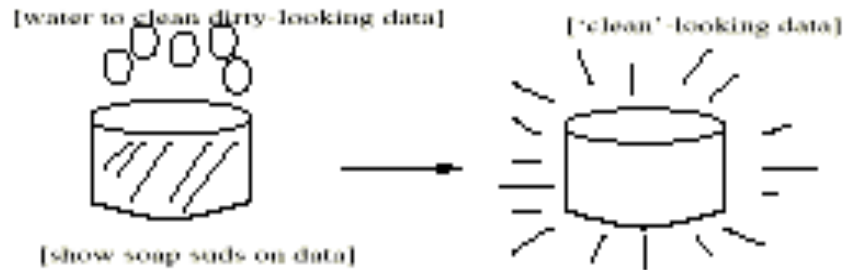
# Major Tasks in Data Preprocessing

- 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
-

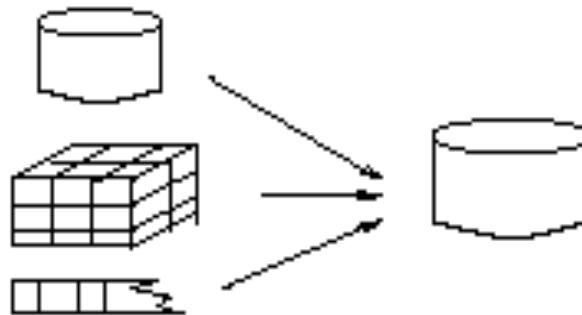


# Forms of Data Preprocessing

## Data Cleaning



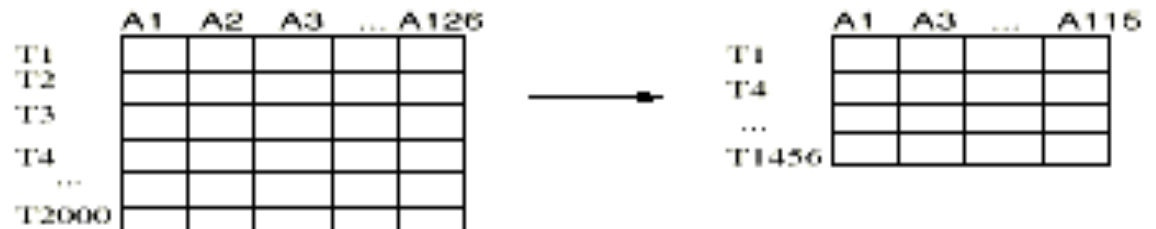
## Data Integration



## Data Transformation

-2, 32, 100, 59, 48 → -0.02, 0.32, 1.00, 0.59, 0.48

## Data Reduction





# Data Cleaning

- Importance

- “Data cleaning is one of the three biggest problems in data warehousing”—Ralph Kimball
- “Data cleaning is the number one problem in data warehousing”—DCI survey

- Data cleaning tasks

- Fill in **missing** values
  - Identify outliers and smooth out **noisy** data
  - Correct inconsistent data
  - Resolve redundancy caused by data integration
-



# Data Cleaning •

## : How to Handle Missing Data? •

- **Ignore the tuple**: 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**
  - 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 regression**
-



# Data Cleaning •

## : 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 •

## : Binning Methods •

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

\* Partition into equal-frequency (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:

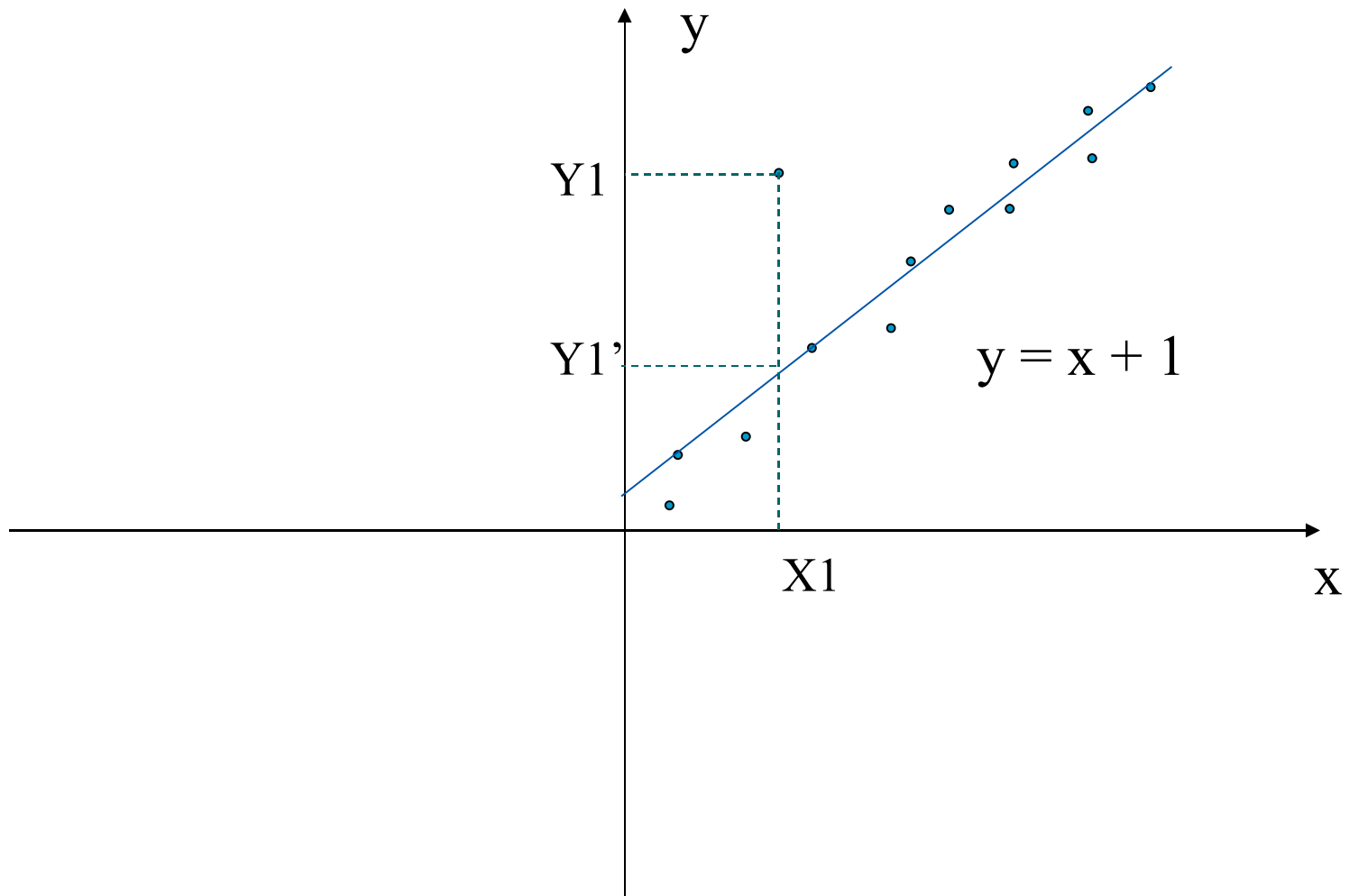
- Bin 1: 9, 9, 9, 9
- Bin 2: 23, 23, 23, 23
- Bin 3: 29, 29, 29, 29

\* Smoothing by bin boundaries:

- Bin 1: 4, 4, 4, 15
  - Bin 2: 21, 21, 25, 25
  - Bin 3: 26, 26, 26, 34
-



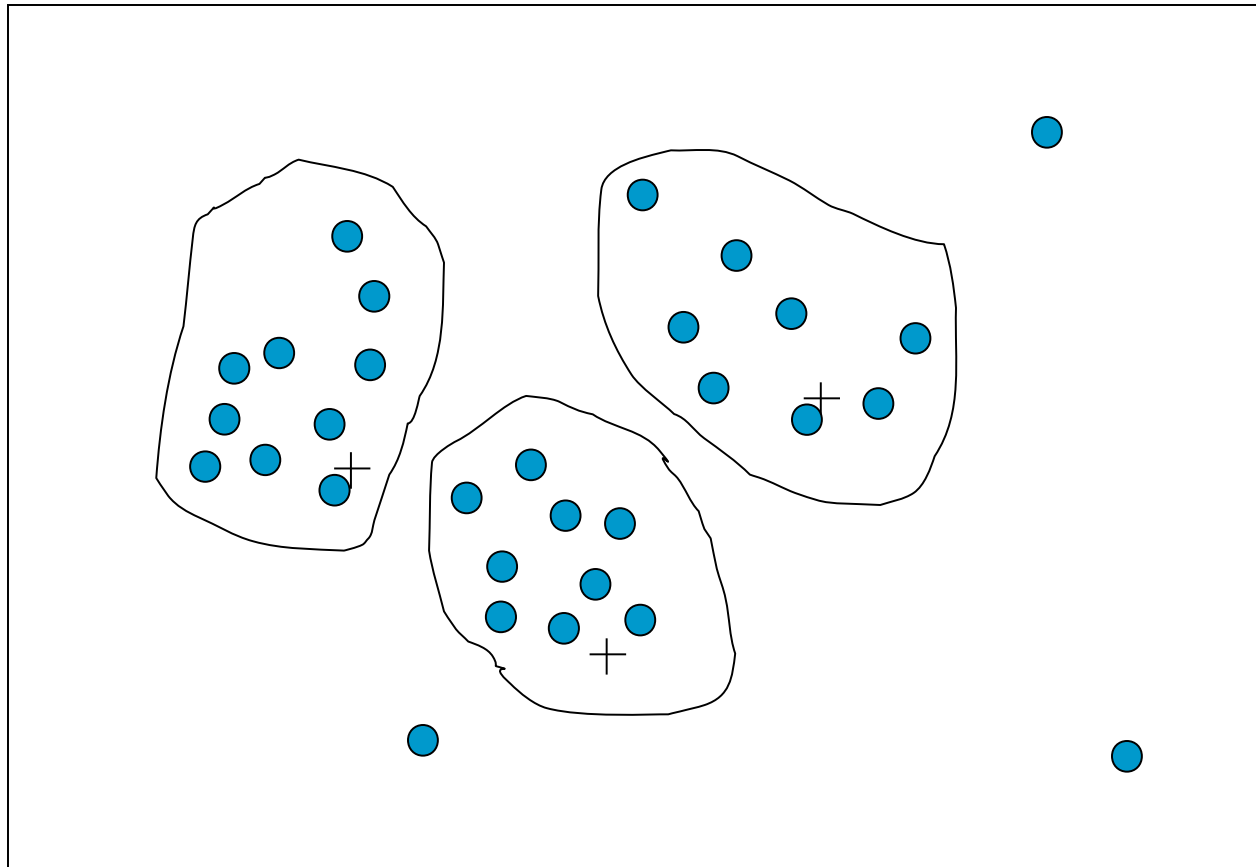
# Data Cleaning : Regression







# Data Cleaning : Cluster Analysis





# Data Integration

- Data integration:
    - Combines data from **multiple sources into a coherent store**
  - Schema integration: e.g.,  $A.\text{cust-id} \equiv B.\text{cust-}\#$ 
    - 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
-



# Data Integration •

## : 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**
  - Careful integration of the data from multiple sources may help reduce/avoid redundancies and inconsistencies and improve mining speed and quality
-



# Data Integration : •

## Correlation Analysis (Numerical Data) •

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

$$r_{A,B} = \frac{\sum (A - \bar{A})(B - \bar{B})}{(n-1)\sigma_A\sigma_B} = \frac{\sum (AB) - n\bar{A}\bar{B}}{(n-1)\sigma_A\sigma_B}$$

where  $n$  is the number of tuples,  $\bar{A}$  and  $\bar{B}$  are the respective means of  $A$  and  $B$ ,  $\sigma_A$  and  $\sigma_B$  are the respective standard deviation of  $A$  and  $B$ , and  $\sum(AB)$  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;  $r_{A,B} < 0$ : negatively correlated
-



# Data Integration •

## : Correlation Analysis (Categorical Data) •

- $\chi^2$  (chi-square) test

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

- The larger the  $\chi^2$  value, the more likely the variables are related
  - The cells that contribute the most to the  $\chi^2$  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
-



# Data Transformation

- Smoothing: remove noise from data
  - 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
-



# Data Transformation •

## : 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 \$73,000 is mapped to  $\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$

- Z-score normalization ( $\mu$ : mean,  $\sigma$ : standard deviation):

$$v' = \frac{v - \mu}{\sigma}$$

- Ex. Let  $\mu = 54,000$ ,  $\sigma = 16,000$ . Then  $\frac{73,600 - 54,000}{16,000} = 1.225$

- 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 Reduction Strategies

- Why data reduction?
    - A database/data warehouse may store terabytes of data
    - Complex data analysis/mining may take a very long time to run on the complete data set
  - 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
  - **Data reduction strategies**
    - Aggregation
    - Sampling
    - Dimensionality Reduction
    - Feature subset selection
    - Feature creation
    - Discretization and Binarization
    - Attribute Transformation
-





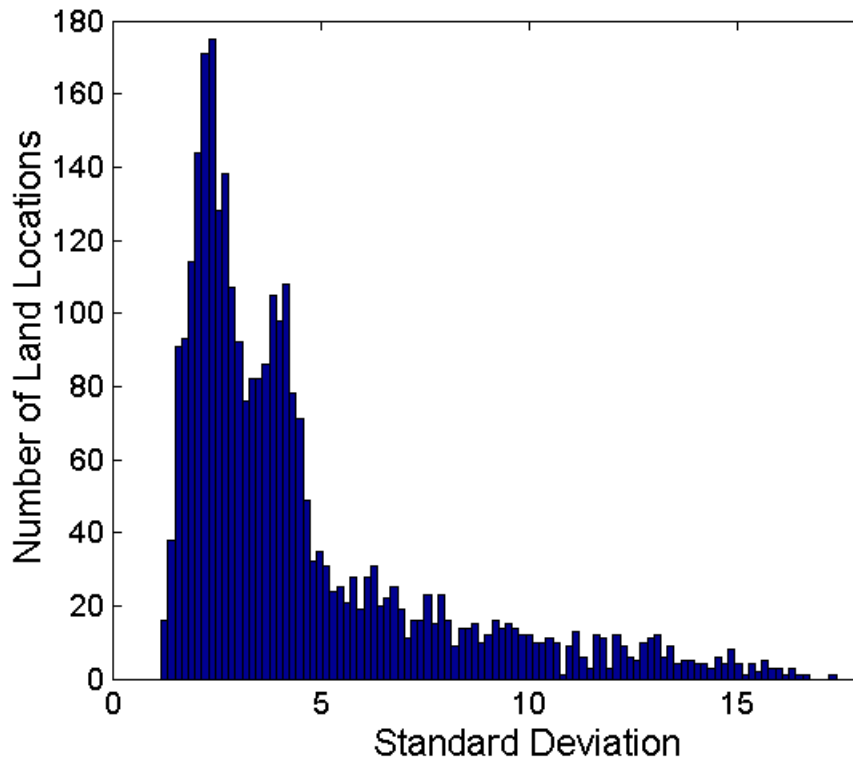
# Data Reduction : Aggregation

- 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
-

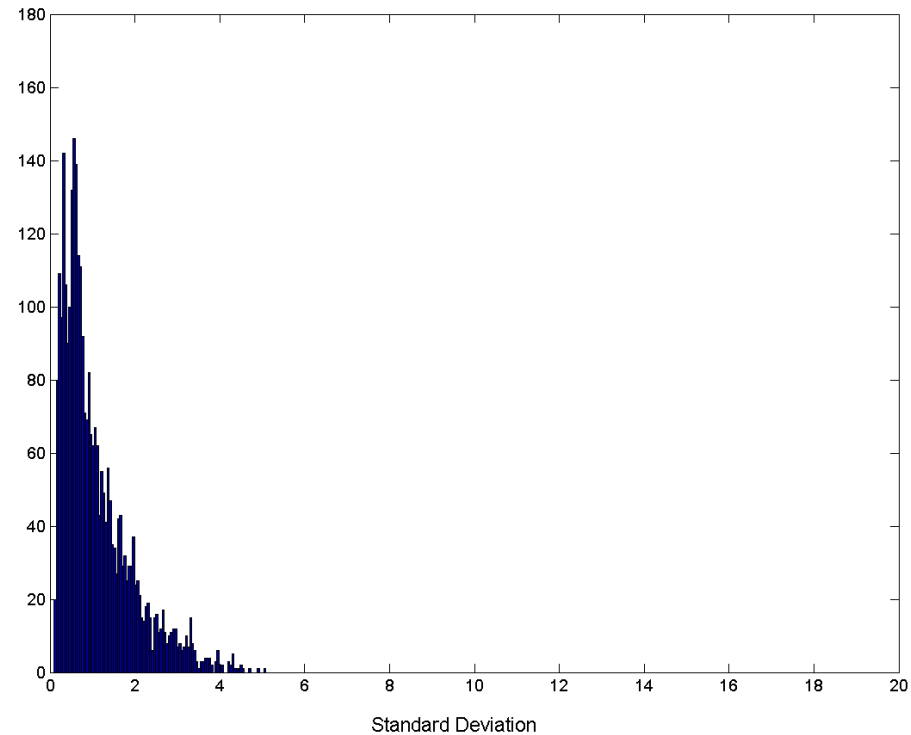


# Data Reduction : Aggregation

## Variation of Precipitation in Australia



**Standard Deviation of Average  
Monthly Precipitation**



**Standard Deviation of Average  
Yearly Precipitation**



# Data Reduction : 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 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.
-



# Data Reduction : 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
-



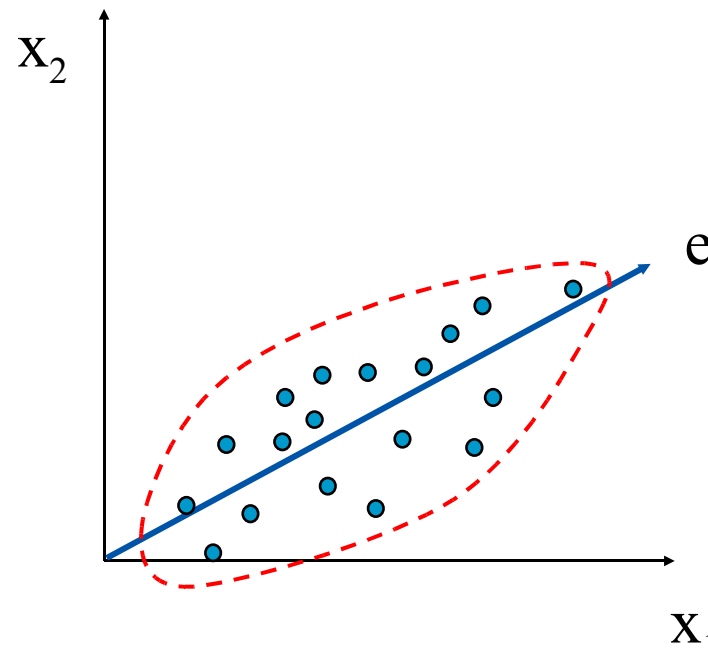
# Data Reduction • : Dimensionality Reduction •

- 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
    - Principle Component Analysis
    - Singular Value Decomposition
    - Others: supervised and non-linear techniques
-



# Dimensionality Reduction : PCA

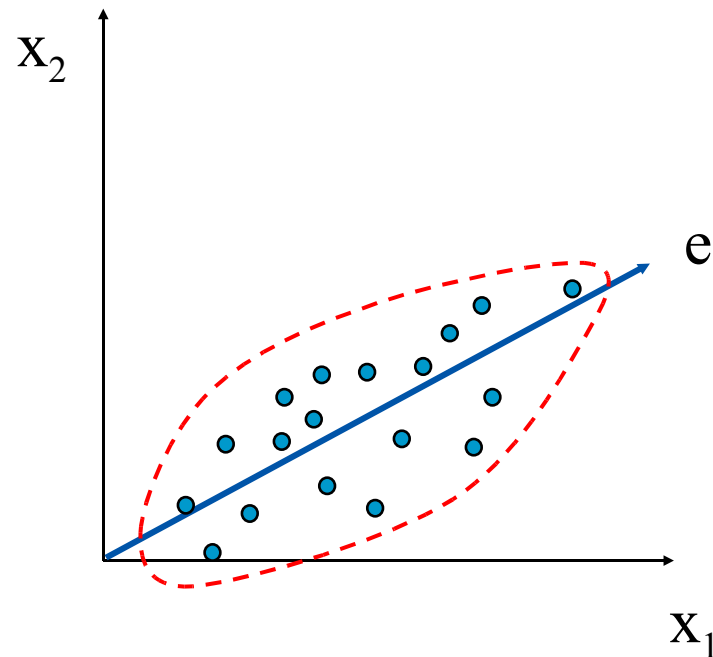
- Goal is to find a projection that captures the largest amount of variation in data





# Dimensionality Reduction : PCA

- Find the eigenvectors of the covariance matrix
- The eigenvectors define the new space





# Data Reduction •

## : Feature Subset Selection •

- Another way to reduce dimensionality of data
  - 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
-





# Data Reduction • : Feature Subset Selection •

- Techniques:
    - Brute-force approach:
      - Try all possible feature subsets as input to data mining algorithm
    - **Filter approaches:**
      - Features are selected before data mining algorithm is run
    - **Wrapper approaches:**
      - Use the data mining algorithm as a black box to find best subset of attributes
-



# Data Reduction • : 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
      - domain-specific
    - Mapping Data to New Space
    - Feature Construction
      - combining features
-



# Question & Answer

---



# 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

- $\chi^2$  (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