# Data Mining: Data Preprocessing

**I211: Information infrastructure II**

# 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**

**Objects**

| 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 |

# Data Preprocessing

- Why preprocess the data?

- Descriptive data summarization (covered!)

- Data cleaning

- Data integration and transformation

- Data reduction

- Discretization and concept hierarchy generation

- Summary

# 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

# Why Is Data Dirty?

- Incomplete data may come from
  - "Not applicable" data value when collected
  - Different considerations between the time when the data was collected and when it is analyzed.
  - Human/hardware/software problems

- Noisy data (incorrect values) may come from
  - Faulty data collection instruments
  - Human or computer error at data entry
  - Errors in data transmission

- Inconsistent data may come from
  - Different data sources
  - Functional dependency violation (e.g., modify some linked data)

- Duplicate records also need data cleaning

# Why Is Data Preprocessing Important?

- No quality data, no quality mining results!
  - Quality decisions must be based on quality data
    - e.g., duplicate or missing data may cause incorrect or even misleading statistics.
  - Data warehouse needs consistent integration of quality data
- Data extraction, cleaning, and transformation comprises the majority of the work of building a data mining system

# Multi-Dimensional Measure of Data Quality

- A well-accepted multidimensional view:
  - Accuracy
  - Completeness
  - Consistency
  - Timeliness
  - Believability
  - Value added
  - Interpretability
  - Accessibility
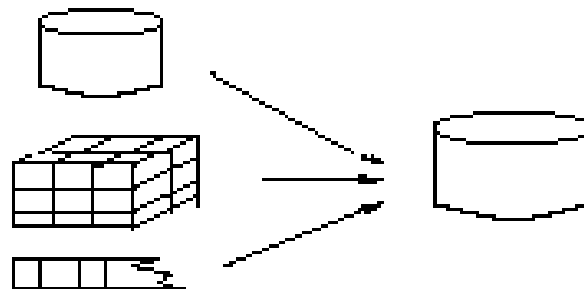
# Major Tasks in Data Preprocessing

- **Data cleaning**
  - Fill in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies

- **Data integration**
  - Integration of multiple databases 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

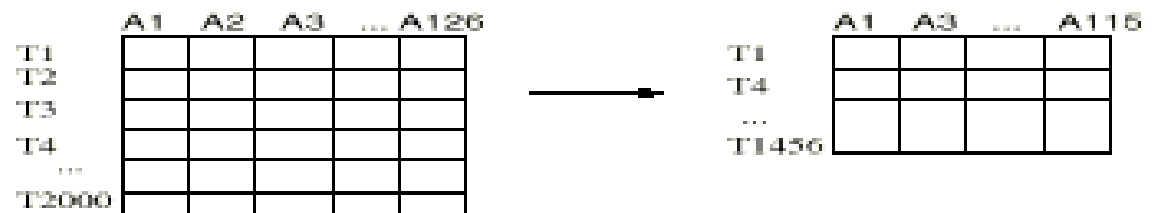# Forms of Data Preprocessing

Data Cleaning

[water to clean dirty-looking data]

[show soap suds on data]

['clean'-looking data]

Data Integration

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

Data Reduction

| | A1 | A2 | A3 | ... A126 |
|---|---|---|---|---|
| T1 | | | | |
| T2 | | | | |
| T3 | | | | |
| T4 | | | | |
| ... | | | | |
| T2000 | | | | |

| | A1 | A3 | ... | A115 |
|---|---|---|---|---|
| T1 | | | | |
| T4 | | | | |
| ... | | | | |
| T1456 | | | | |

# Data Preprocessing

- Why preprocess the data?

- Descriptive data summarization (covered!)

- Data cleaning

- Data integration and transformation

- Data reduction

- Discretization and concept hierarchy generation

- Summary

# Data Cleaning

- Importance
  - garbage in garbage out principle (GIGO)

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

# 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 (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?

- Fill in it automatically with

  - a global constant : e.g., "unknown", a new class?!

  - the attribute mean

  - the attribute mean for all data points 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 due to
  - faulty data collection instruments
  - data entry problems
  - data transmission problems
  - technology limitation
  - inconsistency in naming convention
- Class label noise is hard to deal with
  - sometimes we don't know whether the class label is correct or it is simply unexpected
- Noise demands robustness in training algorithms, that is, training should not be sensitive to noise

# 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)

# Simple Discretization Methods: Binning

- Equal-width (distance) partitioning

    – Divides the range into $N$ intervals of equal size: uniform grid

    – if $A$ and $B$ are the lowest and highest values of the attribute, the width of intervals will be: $W = (B - A)/N.$

    – The most straightforward, but outliers may dominate presentation

    – Skewed data is not handled well

- Equal-depth (frequency) partitioning

    – Divides the range into $N$ intervals, each containing approximately same number of data points

    – Good data scaling

    – Managing categorical attributes can be tricky

# Binning Methods for Data Smoothing

❑ 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

# Cluster Analysis as Binning

# 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)

# Data Preprocessing

- Why preprocess the data?

- Data cleaning

- Data integration and transformation

- Data reduction

- Discretization and concept hierarchy generation

- Summary

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

# Correlation Analysis (Numerical Data)

- Correlation coefficient (Pearson's correlation 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 $\Sigma(AB)$ 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.

- $r_{A,B} = 0$: uncorrelated; $r_{A,B} < 0$: negatively correlated

# Correlation Analysis (Categorical Data)

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

$$\chi^2_{n-1} = \sum_{i=1}^{n} \frac{(Observed_i - Expected_i)^2}{Expected_i}$$

- $n$ is the number of possible values

- The larger the $X^2$ value, the more likely the variables are related

- The cells that contribute the most to the $X^2$ 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 | 200 | 450 |
| Not like science fiction | 50 | 1000 | 1050 |
| Sum (col.) | 300 | 1200 | 1500 |

**Probability to play chess: P(chess) = 300/1500 = 0.2**

**Probability to like science fiction: P(SciFi) = 450/1500 = 0.3**

**If science fiction and chess playing are independent attributes, then the probability to like SciFi AND play chess is**

**P(SciFi, chess) = P(SciFi) · P(chess) = 0.06**

**That means, we expect 0.06 · 1500 = 90 such cases (if they are independent)**

# 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 |

- $X^2$ (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

# Data Transformation

- Smoothing: remove noise from data

- Aggregation: summarization

- 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

# Aggregation

## Variation of Precipitation in Australia



**Standard Deviation of Average Monthly Precipitation**

**Standard Deviation of Average Yearly Precipitation**

# 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

# Attribute 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):

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

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

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

# Data Preprocessing

- Why preprocess the data?

- Data cleaning

- Data integration and transformation

- Data reduction

- Discretization and concept hierarchy generation

- Summary

# Data Reduction Strategies

- Why data reduction?
  - A database 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
  - Data Compression
  - Sampling
  - Discretization and concept hierarchy generation
  - Dimensionality reduction — e.g. remove unimportant attributes

# 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

# Data Compression

# Data Compression (via PCA)

Dimensions = 206

# Data Reduction Method: 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
- 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

# 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

# Sampling: with or without Replacement



SRSWOR
(simple random sample without replacement)

SRSWR

Raw Data

# Sample Size



**8000 points**          **2000 Points**          **500 Points**
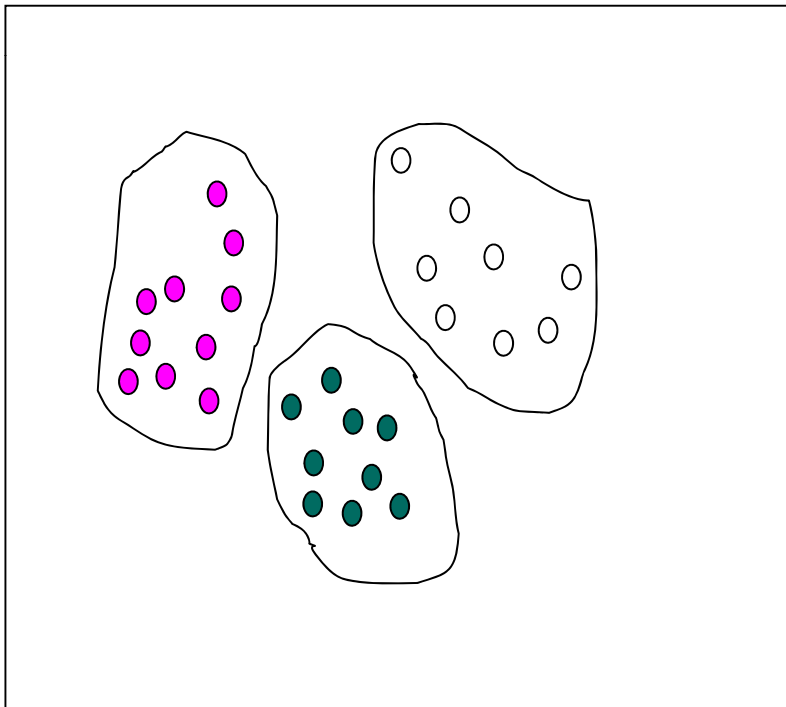
# Sample Size

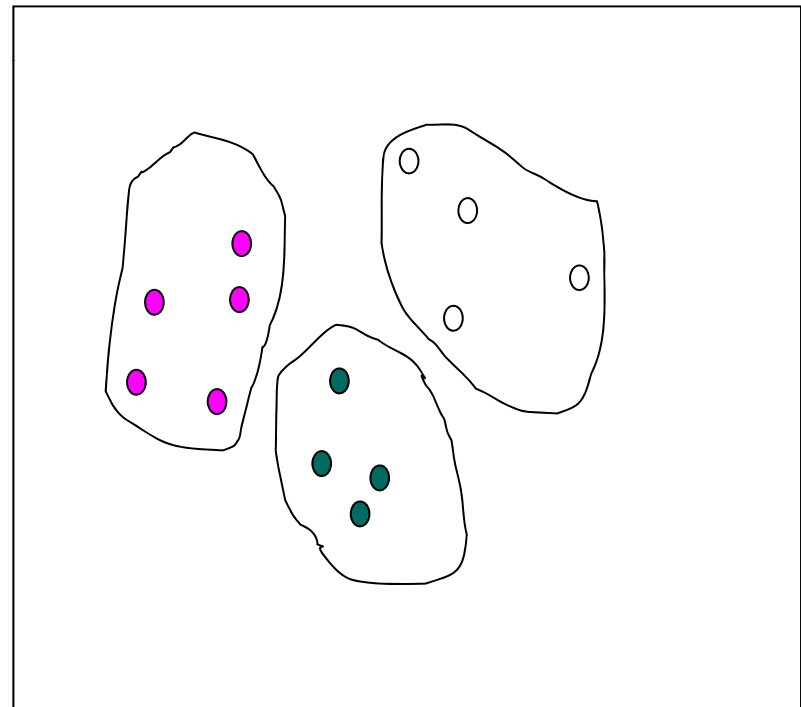- **What sample size is necessary to get at least one object from each of 10 groups.**

# Sampling: Cluster or Stratified Sampling

Raw Data

Cluster/Stratified Sample

# 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

# Feature Subset Selection

- Techniques:
  - Brute-force 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
  - 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

# Feature Creation

- Create new attributes that can capture the important information in a data set much more efficiently than the original attributes

- Methodologies:
  - Mapping Data to New Space
    - Feature construction by combining features
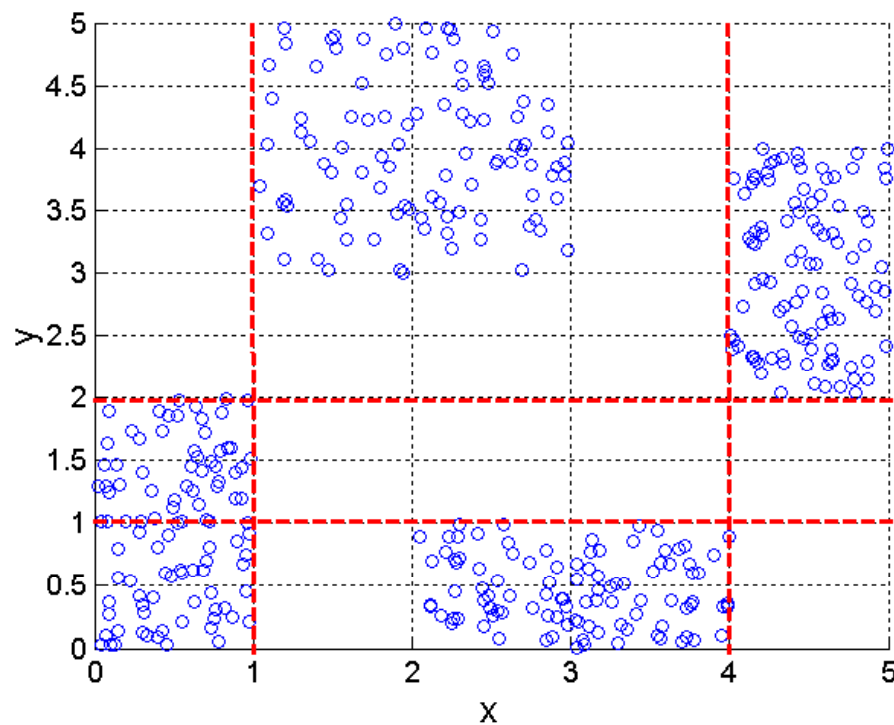
# Data Preprocessing

- Why preprocess the data?

- Data cleaning

- Data integration and transformation

- Data reduction

- Discretization and concept hierarchy generation
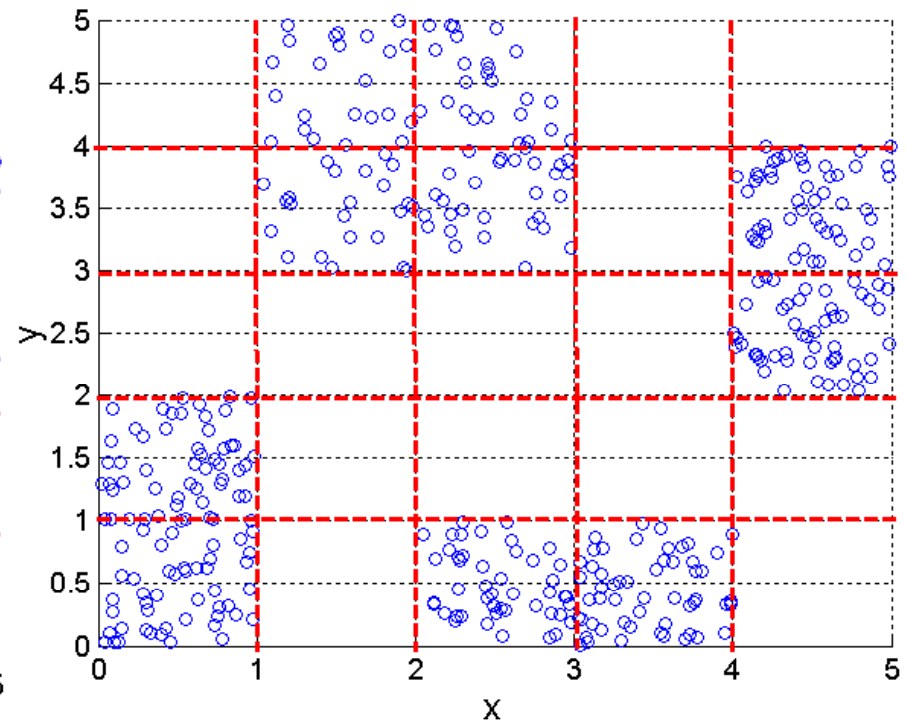
- Summary

# 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

  – Continuous — real numbers, e.g., integer or real numbers (here we aggregated interval and ratio attributes into *continuous*)

- Discretization:

  – 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 Using Class Labels

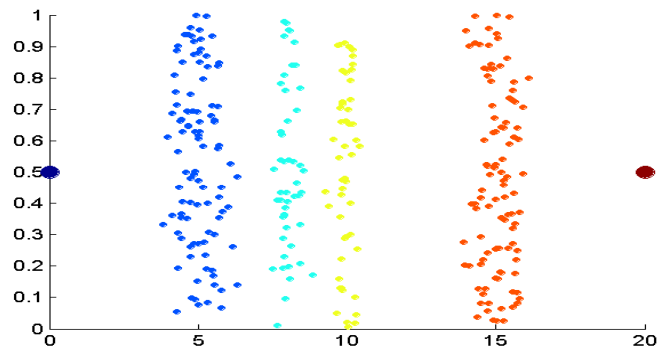● **Entropy based approach**



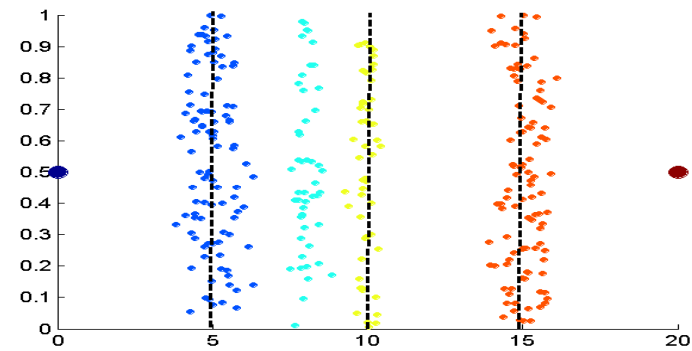3 categories for both x and y
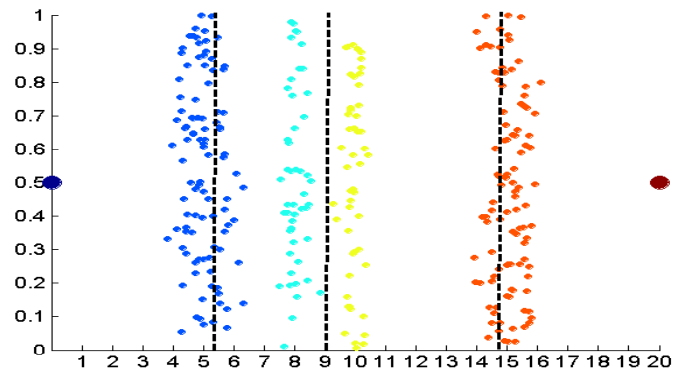
5 categories for both x and y

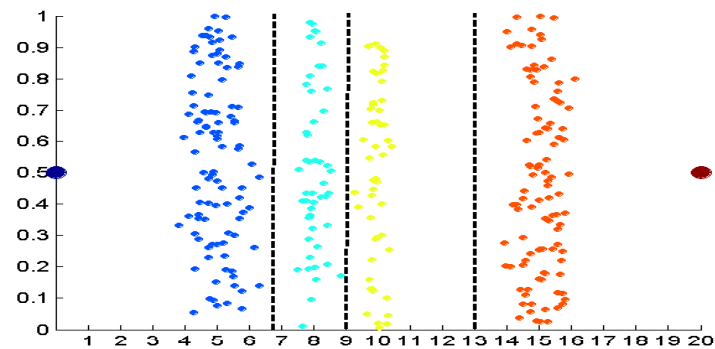# Discretization Without Using Class Labels



Data

Equal interval width

Equal frequency

K-means

# Discretization and Concept Hierarchy

- 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

  - Supervised vs. unsupervised (use class or don't use class variable)

  - Split (top-down) vs. merge (bottom-up)

- Concept hierarchy formation

  - Recursively reduce the data by collecting and replacing low level concepts (such as numeric values for age) by higher level concepts (such as young, middle-aged, or senior)

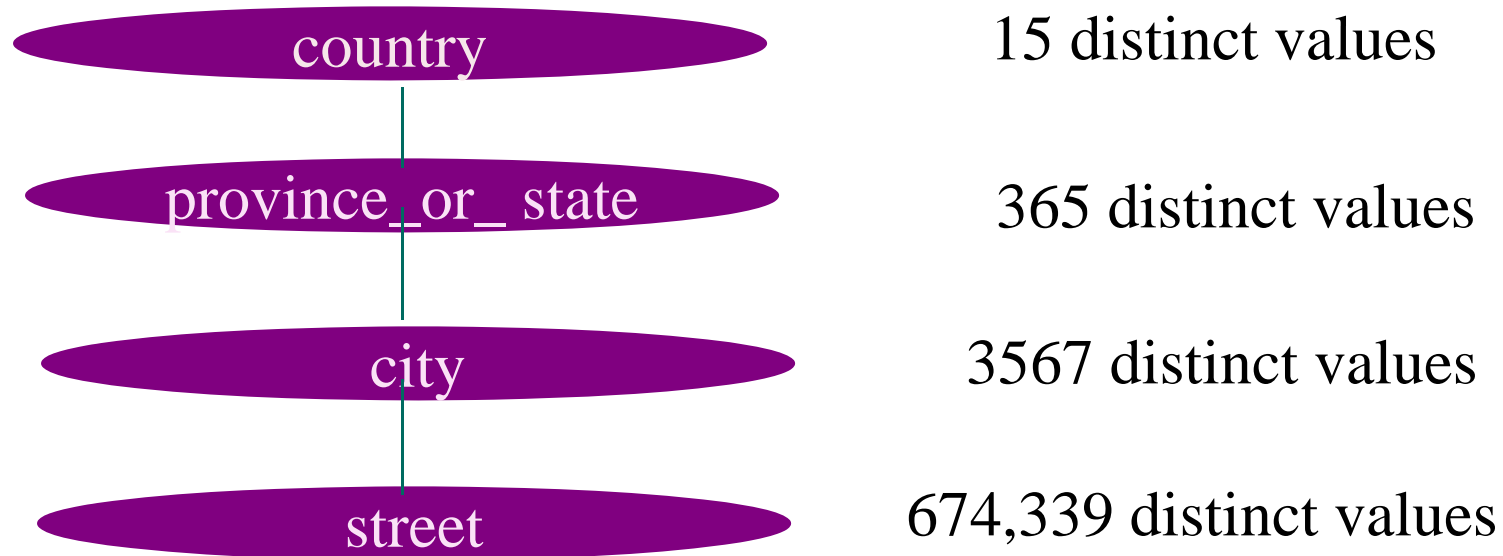# Discretization and Concept Hierarchy Generation for Numeric Data

- Typical methods: All the methods can be applied recursively

  – Binning (covered earlier)

    ◆ Top-down split, unsupervised,

  – Histogram analysis (covered earlier)

    ◆ Top-down split, unsupervised

  – Clustering analysis (covered earlier and in more detail later)

    ◆ Either top-down split or bottom-up merge, unsupervised

  – Entropy-based discretization: supervised, top-down split

  – Interval merging by $\chi^2$ Analysis: unsupervised, bottom-up merge

  – Segmentation by natural partitioning: top-down split, unsupervised

# Concept Hierarchy Generation for Categorical Data

- Specification of a partial/total ordering of attributes explicitly at the schema level by users or experts
  - street < city < state < country
- Specification of a hierarchy for a set of values by explicit data grouping
  - {Urbana, Champaign, Chicago} < Illinois
- Specification of only a partial set of attributes
  - E.g., only street < city, not others
- Automatic generation of hierarchies (or attribute levels) by the analysis of the number of distinct values
  - E.g., for a set of attributes: {street, city, state, country}

# Automatic Concept Hierarchy Generation

- Some hierarchies can be automatically generated based on the analysis of the number of distinct values per attribute in the data set
  - The attribute with the most distinct values is placed at the lowest level of the hierarchy
  - Exceptions, e.g., weekday, month, quarter, year

| | |
|---|---|
| country | 15 distinct values |
| province_or_ state | 365 distinct values |
| city | 3567 distinct values |
| street | 674,339 distinct values |

# Data Preprocessing

- Why preprocess the data?

- Data cleaning

- Data integration and transformation

- Data reduction

- Discretization and concept hierarchy generation

- Summary

# Summary

- Data preparation or preprocessing is a big issue for data mining

- Descriptive data summarization is need for quality data preprocessing

- Data preparation includes

  – Data cleaning and data integration

  – Data reduction and feature selection

  – Discretization

- A lot a methods have been developed but data preprocessing still an active area of research