

Data Preprocessing

Slides by: Shree Jaswal

Topics to be covered

- Why Preprocessing? Data Cleaning; Data Integration;
- Data Reduction: Attribute subset selection,
- Histograms, Clustering and Sampling;
- Data Transformation & Data Discretization: Normalization, Binning, Histogram Analysis and Concept hierarchy generation.

Why Data Preprocessing?

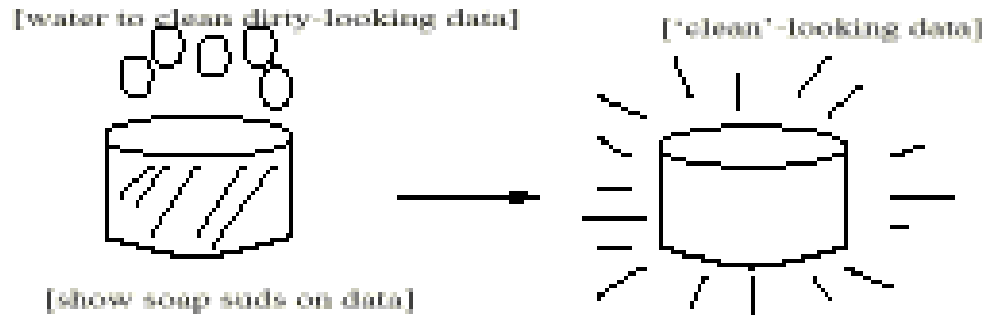
- Data in the real world is dirty
 - **incomplete**: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
 - **noisy**: containing errors or outliers
 - **inconsistent**: containing discrepancies in codes or names
- No quality data, no quality mining results!
 - Quality decisions must be based on quality data
 - Data warehouse needs consistent integration of quality data
- A multi-dimensional measure of data quality:
 - A well-accepted multi-dimensional view:
 - accuracy, completeness, consistency, timeliness, believability, value added, interpretability, accessibility
 - Broad categories:
 - intrinsic, contextual, representational, and accessibility.

Major Tasks in Data Preprocessing

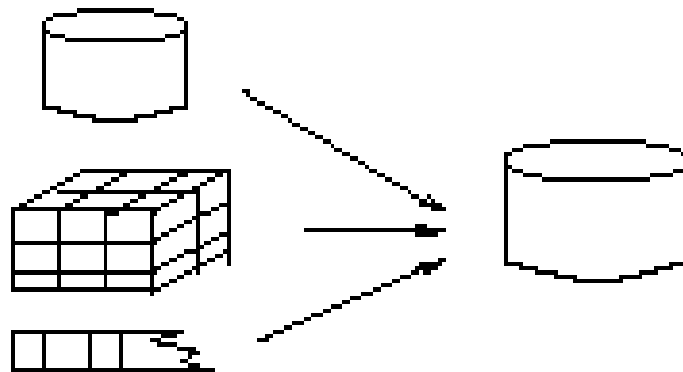
- **Data cleaning**
 - Fill in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies
- **Data integration**
 - Integration of multiple databases, data cubes, files, or notes
- **Data transformation**
 - Normalization (scaling to a specific range)
 - Aggregation
- **Data reduction**
 - Obtains reduced representation in volume but produces the same or similar analytical results
 - Data discretization: with particular importance, especially for numerical data
 - Data aggregation, dimensionality reduction, data compression, generalization

Forms of data preprocessing

Data Cleaning



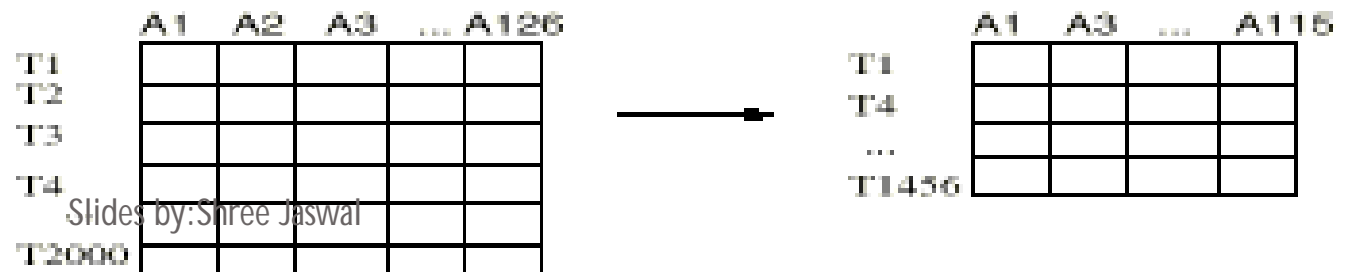
Data Integration



Data Transformation

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

Data Reduction



Data Cleaning

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

Missing Data

- Data is not always available
 - E.g., many 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 task is classification—not effective in certain cases)
- Fill in the missing value **manually**: tedious + infeasible?
- Use a **global constant** to fill in the missing value: e.g., “unknown”, a new class?!
- Use the **attribute mean** to fill in the missing value
- Use the **attribute mean for all samples of the same class** to fill in the missing value: smarter
- Use the **most probable value** to fill in the missing value: inference-based such as regression, Bayesian formula, decision tree

Noisy Data

- Q: What is noise?
- A: Random error 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 requires data cleaning
 - duplicate records
 - incomplete data
 - inconsistent data

How to Handle Noisy Data?

- Binning method:
 - first sort data and partition into (equi-depth) bins
 - then one can **smooth by bin means**, **smooth by bin median**, **smooth by bin boundaries**, etc.
 - used also for discretization (discussed later)
- Clustering
 - detect and remove outliers
- Semi-automated method: combined computer and human inspection
 - detect suspicious values and check manually
- Regression
 - smooth by fitting the data into regression functions

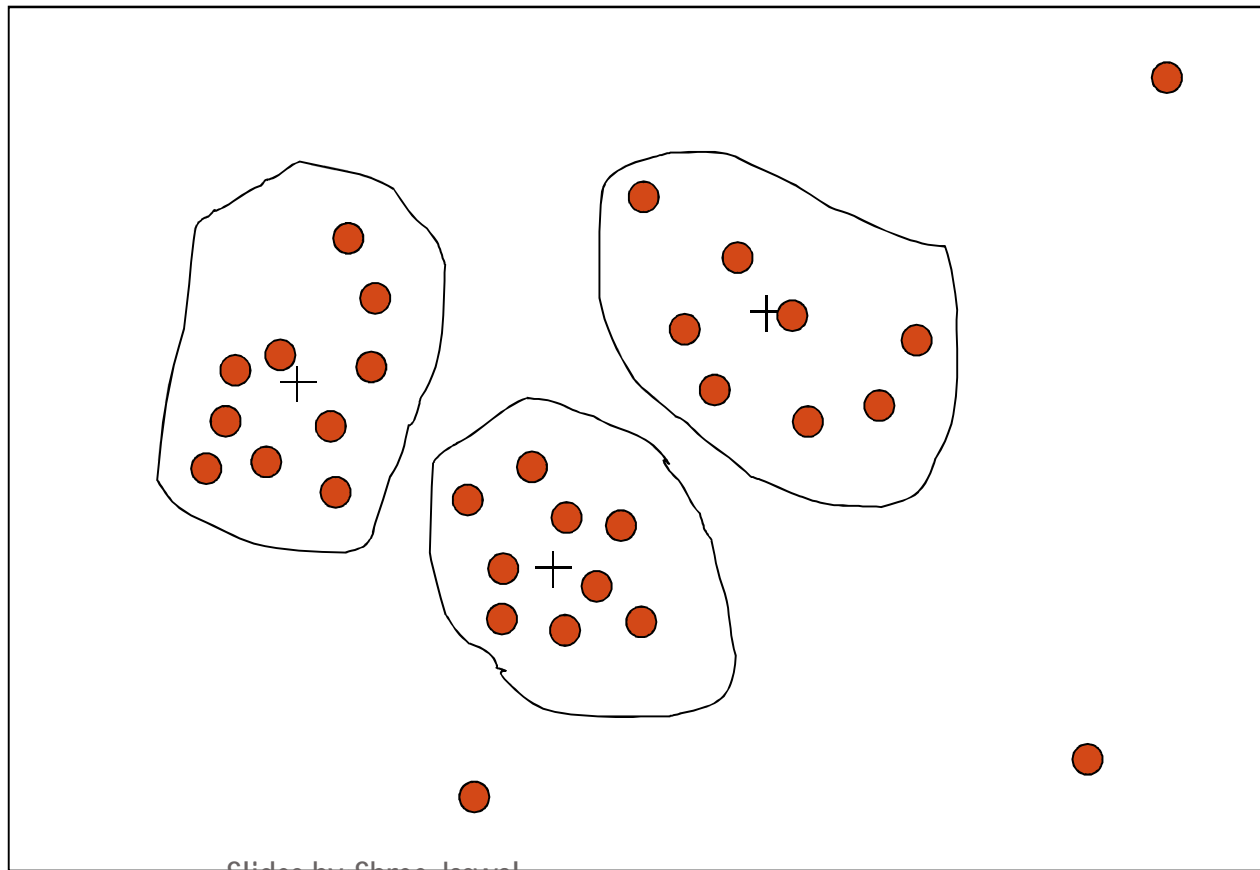
Simple Discretization Methods: Binning

- **Equal-width** (distance) partitioning:
 - It 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:
 - It divides the range into N intervals, each containing approximately same number of samples
 - Good data scaling
 - Managing categorical attributes can be tricky.

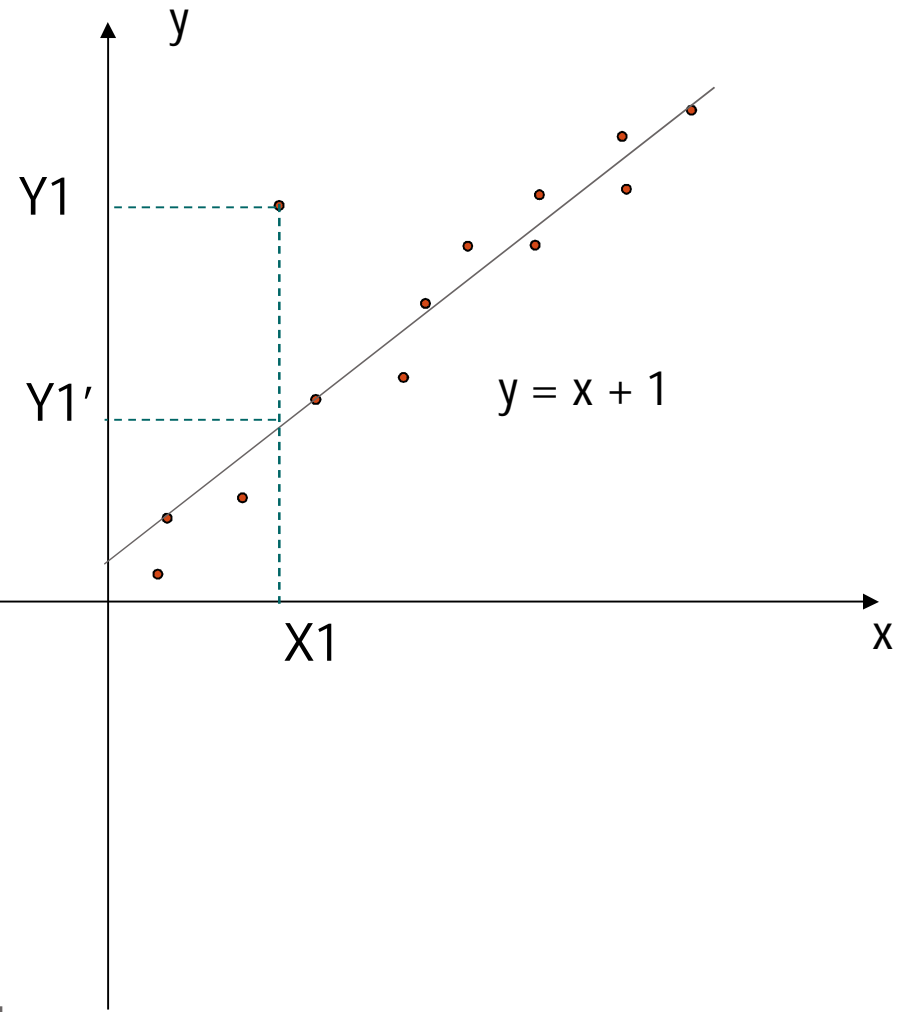
Binning Methods for Data Smoothing

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



Regression



- Linear regression (best line to fit two variables)
- Multiple linear regression (more than two variables, fit to a multidimensional surface)

How to Handle Inconsistent Data?

- Manual correction using external references
- Semi-automatic using various tools
 - To detect violation of known functional dependencies and data constraints
 - To correct redundant data

Data Integration

- Data integration:
 - combines data from multiple sources into a coherent store
- Schema integration
 - integrate metadata from different sources
 - Entity identification problem: identify real world entities from multiple data sources, e.g., A.cust-id \equiv B.cust-#
- 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, different currency

Handling Redundant Data in Data Integration

- Redundant data occur often when integrating multiple DBs
 - The same attribute may have different names in different databases
 - One attribute may be a “derived” attribute in another table, e.g., annual revenue
- Redundant data may be able to be detected by correlational analysis
- **Careful integration** can help reduce/avoid redundancies and inconsistencies and improve mining speed and quality

Topics to be covered

- Why Preprocessing? Data Cleaning; Data Integration;
- **Data Reduction: Attribute subset selection,**
- Histograms, Clustering and Sampling;
- Data Transformation & Data Discretization: Normalization, Binning, Histogram Analysis and Concept hierarchy generation.

Data Reduction

- **Problem:**

Data Warehouse may store terabytes of data: Complex data analysis/mining may take a very long time to run on the complete data set

- **Solution?**

- Data reduction...

Data Reduction

- 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

Dimensionality Reduction

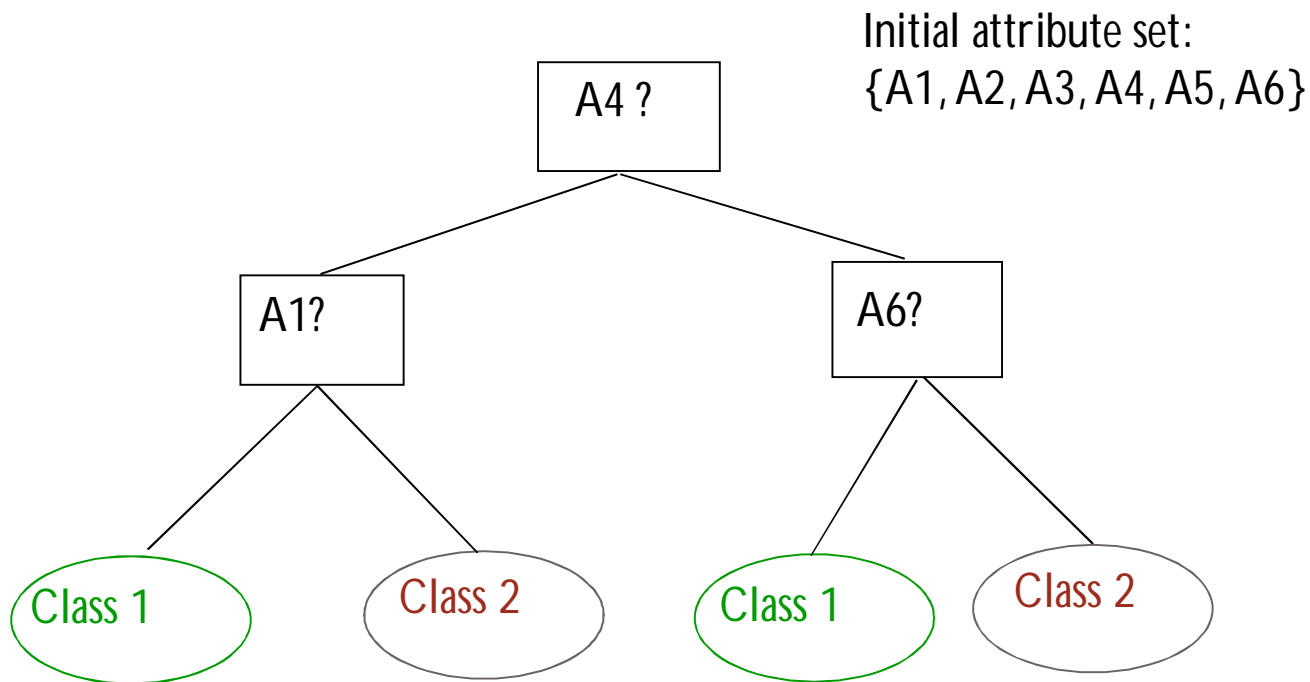
- **Problem:** Feature selection (i.e., **attribute subset selection**):
 - Select a **minimum set of features** such that the probability distribution of different classes given the values for those features is as close as possible to the original distribution given the values of all features
 - **Nice 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

Example of Decision Tree Induction

nonleaf nodes: tests

branches: outcomes of tests

leaf nodes: class prediction



-----> Reduced attribute set: $\{A1, A4, A6\}$

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**

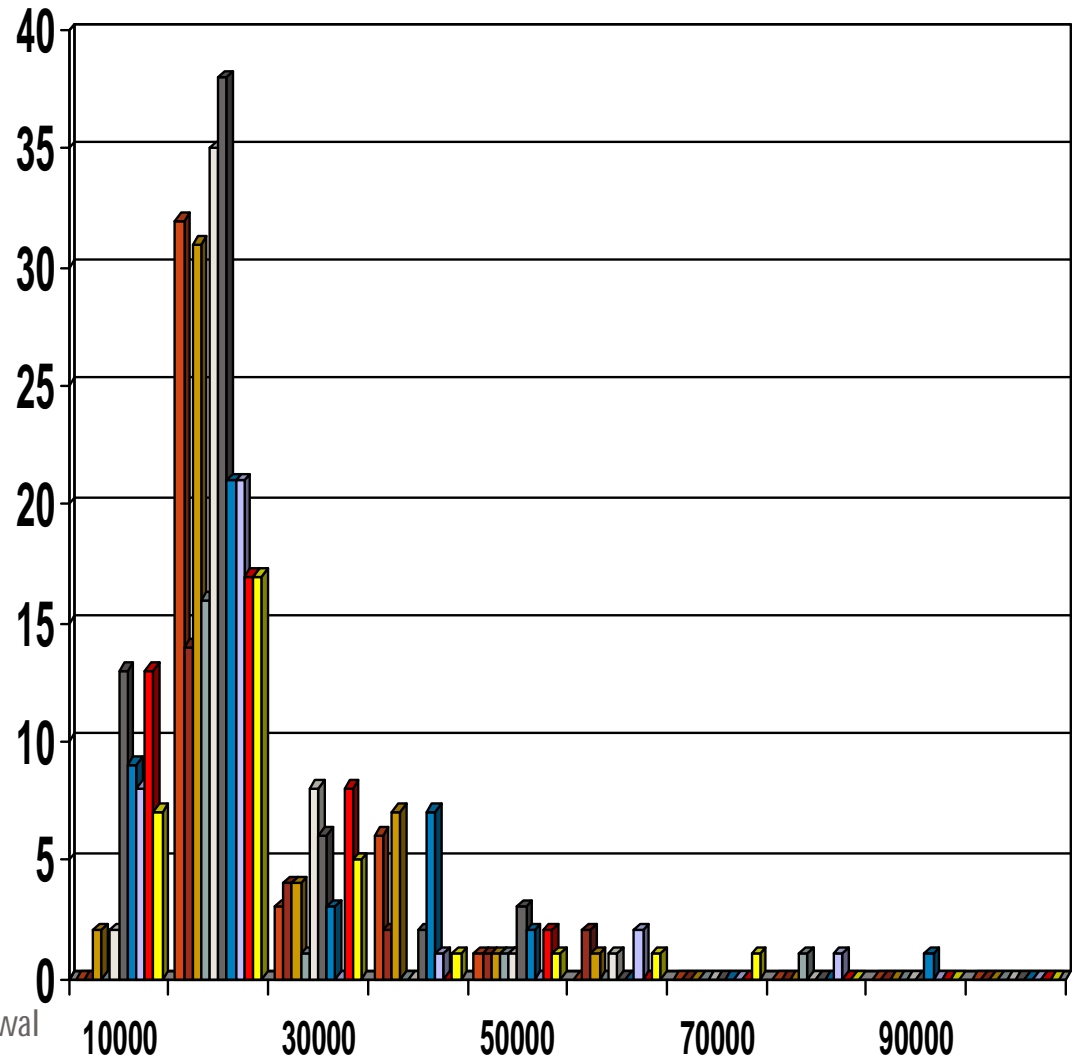
- Do not assume models
- Major families: histograms, clustering, sampling

Topics to be covered

- Why Preprocessing? Data Cleaning; Data Integration;
- Data Reduction: Attribute subset selection,
- **Histograms, Clustering and Sampling;**
- Data Transformation & Data Discretization: Normalization, Binning, Histogram Analysis and Concept hierarchy generation.

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.



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

- 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
- Common ways of sampling are:
 - Simple random sample without replacement (SRSWOR)
 - Simple random sampling may have very poor performance in the presence of skew
 - Simple random sample with replacement (SRSWR)
 - Cluster sample: Eg. page
 - Stratified sample: Approximate the percentage of each class (or subpopulation of interest) in the overall database
 - Used in conjunction with skewed data

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.**

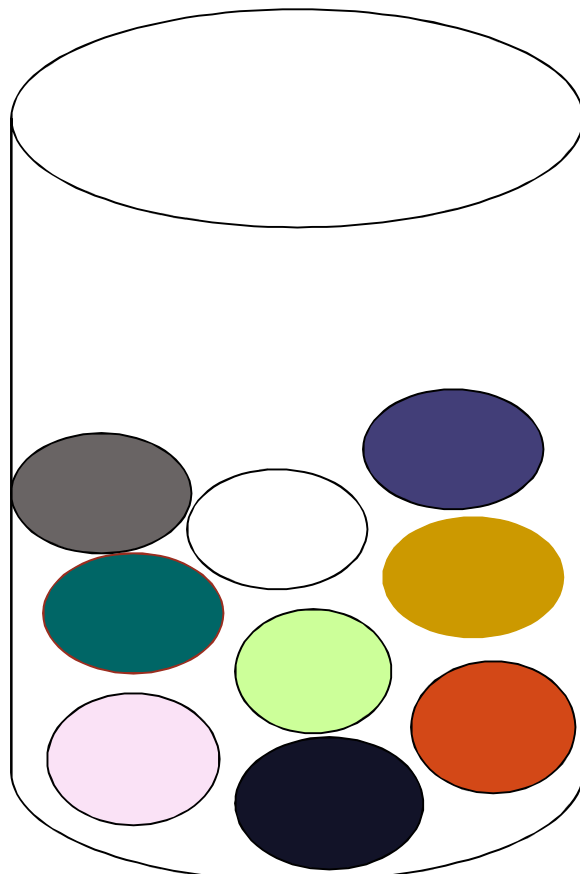
Sampling ...

- The key principle for effective sampling is the following:
 - using a sample will work almost as well as using the entire data sets, if the sample is representative
 - A sample is representative if it has approximately the same property (of interest) as the original set of data

Types of Sampling

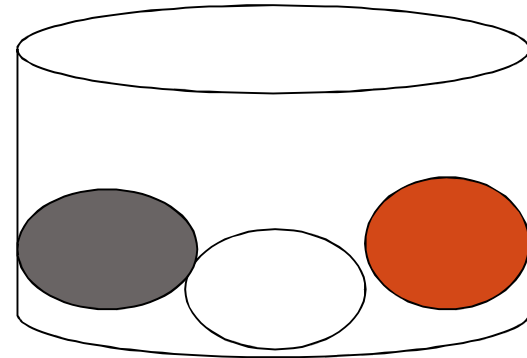
- Simple Random Sampling
 - There is an equal probability of selecting any particular item
- Sampling without replacement (SRSWOR)
 - As each item is selected, it is removed from the population
- Sampling with replacement (SRSWR)
 - Objects are not removed from the population as they are selected for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

Sampling

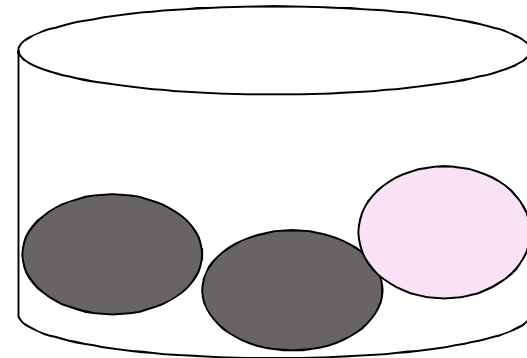


Raw Data

SRSWOR
(simple random
sample without
replacement)

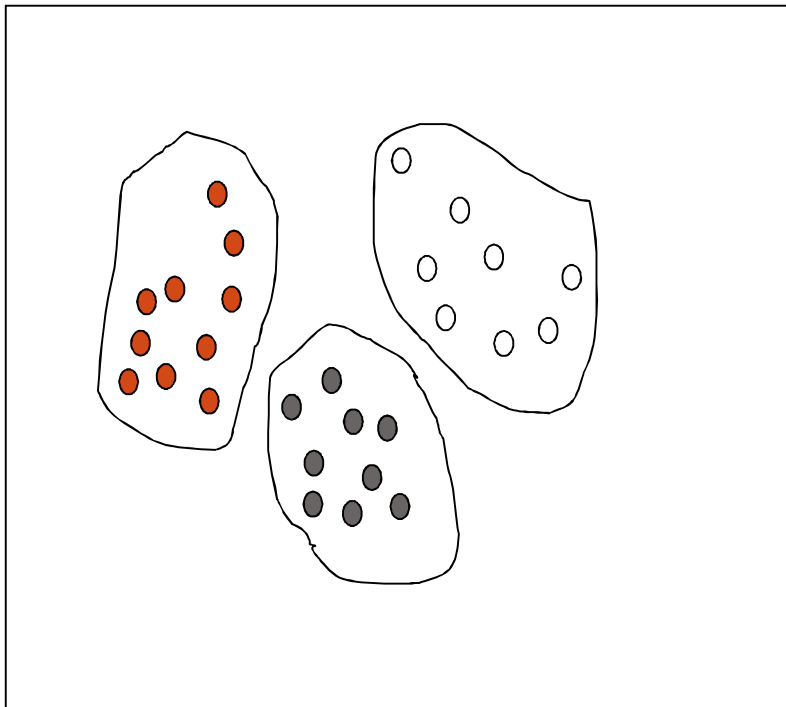


SRSWR

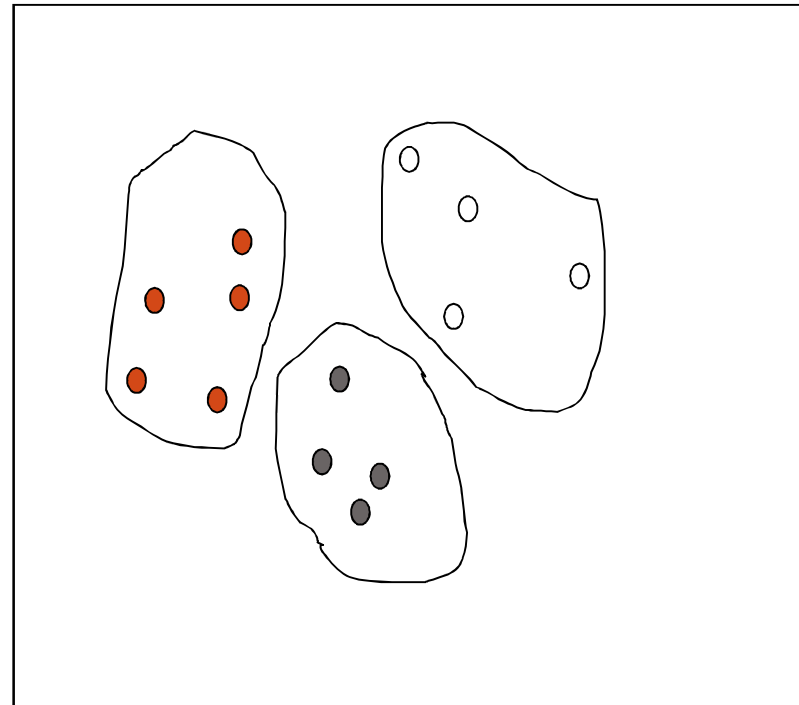


Sampling

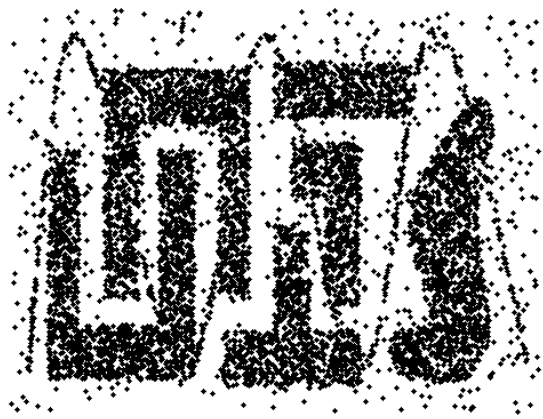
Raw Data



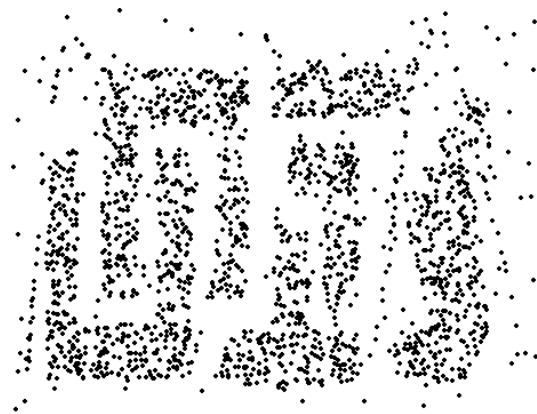
Cluster/Stratified Sample



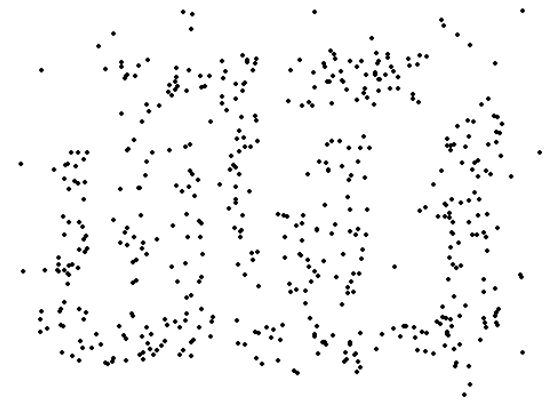
Sample Size



8000 points



2000 Points



500 Points

Topics to be covered

- Why Preprocessing? Data Cleaning; Data Integration;
- Data Reduction: Attribute subset selection, Histograms,
- Clustering and Sampling;
- Data Transformation & Data Discretization: Normalization, Binning, Histogram Analysis and Concept hierarchy generation.

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

Data Transformation: Normalization

Particularly useful for classification (NNs, distance measurements, nn classification, etc)

- min-max normalization

$$v' = \frac{v - \text{min}_A}{\text{max}_A - \text{min}_A} (\text{new_max}_A - \text{new_min}_A) + \text{new_min}_A$$

- z-score normalization

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

- normalization by decimal scaling

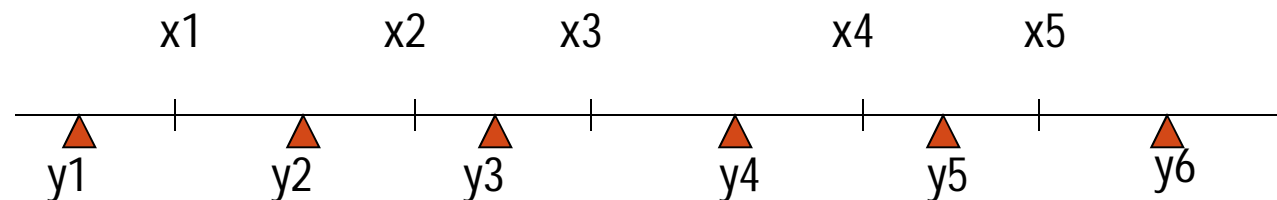
$$v' = \frac{v}{10^j} \quad \text{Where } j \text{ is the smallest integer such that } \text{Max}(|v'|) < 1$$

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.
- 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/Quantization

- 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 hierarchy generation for numeric data

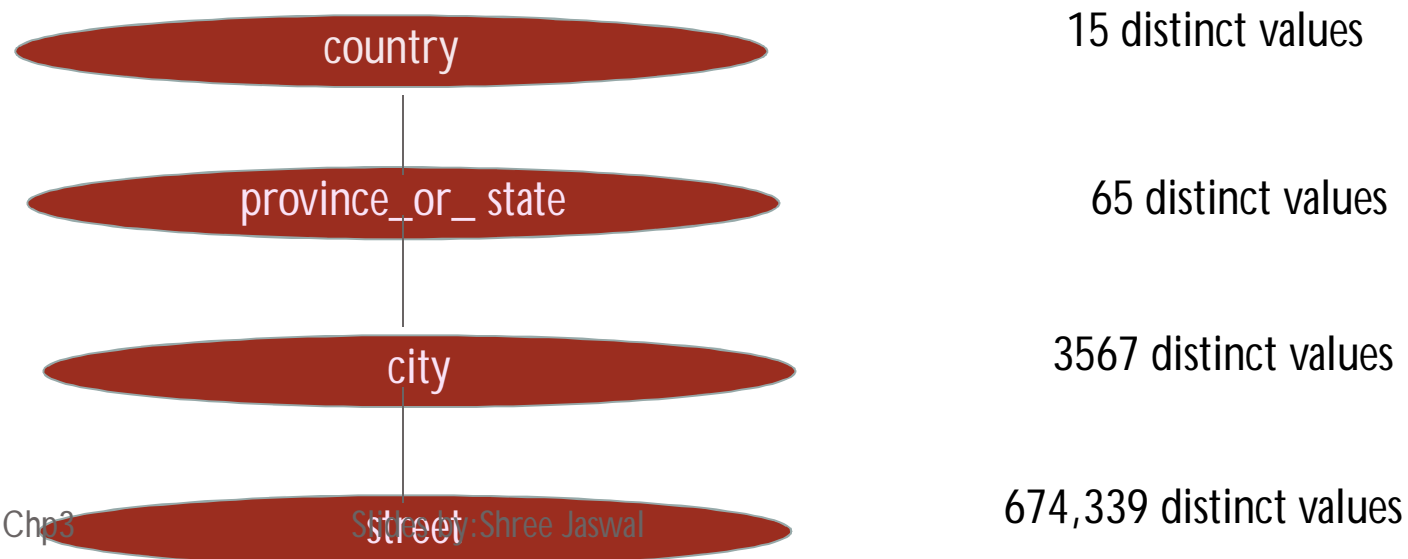
- Hierarchical and recursive decomposition using:
 - Binning (data smoothing)
 - Histogram analysis (numerosity reduction)
 - Clustering analysis (numerosity reduction)
- Entropy-based discretization
- Segmentation by natural partitioning

Concept hierarchy generation for categorical data

- Categorical data: no ordering among values
- Specification of a partial ordering of attributes explicitly at the schema level by users or experts
- Specification of a portion of a hierarchy by explicit data grouping
- Specification of a set of attributes, but not of their partial ordering
- Specification of only a partial set of attributes

Concept hierarchy generation w/o data semantics - Specification of a set of attributes

Concept hierarchy can be automatically generated based on the number of distinct values per attribute in the given attribute set. The attribute with the most distinct values is placed at the lowest level of the hierarchy (**limitations?**)



Summary

- Data preparation is a big issue for both warehousing and mining
- Data preparation includes
 - Data cleaning and data integration
 - Data reduction and feature selection
 - Discretization
- A lot a methods have been developed but still an active area of research