

Data Cleaning

Real world data → tend to be incomplete, noisy and inconsistent.

→ Data cleaning ~~and~~ routines attempt to fill in missing values, smooth out noise while identifying outliers and correct inconsistencies in data, incomplete • noisy • inconsistent • intentional

◆ missing values:

Data is not always available,

e.g.: customer income = "",

missing values due to

- equipment malfunction
- inconsistent thus deleted
- not entered due to misunderstanding
- not imp while entry

Methods:

1) Ignore the tuple: usually done when class label is missing.

→ not effective

→ poor when % of missing values per attribute varies considerably.

→ By ignoring we do not make use of remaining attribute's values in the tuple.

2) Fill in the missing value manually: Time-consuming
It may not be feasible given a large data set in many missing values.

3) Use a global constant to fill in missing values: "Unknown" or $-\infty$
This method is simple not foolproof

- 4) Use a measure of central tendency for the attribute to fill in the missing value:
 Normal data distribution \rightarrow mean can be used
 skewed data distribution \rightarrow median.
- 5) Use the attribute mean or median for all samples belonging to same class as the given tuple.
 ex: credit card risk classification.
- 6) Use the most probable value to fill in the missing value.
 \rightarrow determined by with regression, decision tree, Bayesian formulism.



Noisy data

Noise - random error or variance in a ~~measure~~ variable

causes

- faulty data collection instruments,

- data entry problems,

- data transmission problems,

- technology limitation,

- inconsistency in naming convention.

methods to handle noisy data:

Binning

1) Sort data

2) partition into equal frequency bins.

3) a) smooth by bin means, each ~~out~~ value in bin is replaced by mean value of bin.

b) smooth by bin medians; replaced by median.

c) Smooth by bin boundaries; replaced by closest boundary value.

- Regression, ~~multiple linear regression~~
Smooth by fitting data into regression function.
 - linear regression
 - multiple linear regression
- Outlier analysis: Detected by "Clustering"
- Combined computer & human inspection.
Detect suspicious values & check by human.

Data Integration

- merging of data from multiple data stores.
- Careful integration can help reduce & avoid redundancies & inconsistencies.
- improve speed & accuracy of Dm. process.

(1) Entity identification problem:

- data analysis involves data integration which combines data from multiple sources into data warehouse.

Issues for data integration → schema integration
→ object matching

These are called as Entity Identification problem

When matching attributes from one db to another db etc., special attention must be paid to "structure" of data.

→ This ensures any attributes func. dependences & constraints match to target system.

ex: discount in each order → discount in each item

SOURCE → TARGET sys

② Redundancy and Correlation Analysis.

Redundant data can occur when integration

- Object identification : Same attribute, different names.

- Derivable data : One attribute may be "derived" attribute in another table.

Solution:

Some redundancies can be detected by
Correlation Analysis OR AND
Covariance Analysis.

Such analysis can measure how strongly
 one attribute implies the other.

For nominal data $\rightarrow \chi^2$ test

numeric attributes \rightarrow Correlation coefficient
 and covariance.

③ Tuple duplication.

\rightarrow Duplication should also be detected at tuple level.

\rightarrow Use of denormalized tables is another source of data redundancy.

\rightarrow Inconsistencies arise between duplicates

ex: purchase order db contains attributes
 for the purchaser's name and address
 instead of key to this information,

\rightarrow discrepancies can occur.

④ Data value Conflict Detection and Resolution:

→ ~~* Data Attribute values from different sources may differ.~~

This maybe due to differences in representation, scaling or encoding.

ex: One university may adopt a quarter system, offer 3 courses on DBMS and assign grades A to F while other

University may opt semester system, 2 courses DBMS, assign grades 1 to 10.

It is difficult to work out precise course to grade transformation rules between universities.

→ Attributes can also differ at abstraction level.

ex: totalsales in one db → refers to one branch of company

while other db

total_sales refers to all branches

Data Reduction

→ Techniques can be applied to obtain a reduced representation of the data set that is much smaller yet closely maintains the integrity of original data.

∴ mining on reduced data should be more efficient

yet produce the same results.

Data Reduction Strategies :

1
2
3

- Dimensionality reduction
- Numerosity reduction
- Data compression

1

Process of reducing the number of random variables or attributes under consideration.

- wavelet transforms
- principal component analysis, Transforming data onto a smaller space.

- Attribute subset selection -

↳ irrelevant, weakly relevant, or redundant attributes or dimensions are detected and removed.

2

Replaces the original data volume by alternative, smaller forms of data representation.

ex: regression → ↳ parametric - only data params stored

ex: histograms, → ↳ non parametric -

↳ clustering, sampling, data cube aggregation.

3

Transformations are applied so as to obtain a reduced, or "compressed" representation of the data.

→ if original data can be reconstructed from the compressed data w/o any info loss



lossless.

→ we can reconstruct only approx. of original data ⇒ lossy.

Attribute Subset Selection (Refer ref book)

→ Data sets may contain hundreds of attributes, many are irrelevant to mining task or redundant. Attribute Subset Selection removes it.

goal → Find min set of attributes such that resulting probability distribution of the data classes is as close as possible to orig distribution.

⇒ greedy heuristic methods that explore a reduced search space are commonly used for Attr. Subsel. selec..

methods:

1. Stepwise forward selection:

- ① empty set
- ② Best of orig attributes is determined and added to reduced set.

2. Stepwise backward elimination

- ① Full set of attributes
- ② removes worst attribute.

3. Combination of forward & backward selection elimination.

- ① at each step procedure selects the best attribute & removes worst.

4. Decision tree induction.

ID 3, C4.5, CART originally intended for classification,

→ flowchart like

• internal node - denotes test on an

attribute, branch denotes outcome.

• leaf node - class prediction.

Histograms

→ Use binning

If each bucket represents only a single attribute-value / frequency pair, buckets are called singleton buckets.

Partition Rules:

Equal width → equal bucket range

Equal freq → (or equal depth)

Adv: effective at approximating both dense & sparse data, highly skewed & uniformed data.

Clustering (Refer PPT)

→ consider data tuples as Obj.

① partition objects into groups / clusters based on

② (similarity) ~~is defined~~

③ (quality of) a cluster by diameter.

~~Centroid~~

→ can be effective if data is clustered but not if data is ~~smeared~~

→ Can have hierarchical clustering & be stored in multidimensional index tree structures.

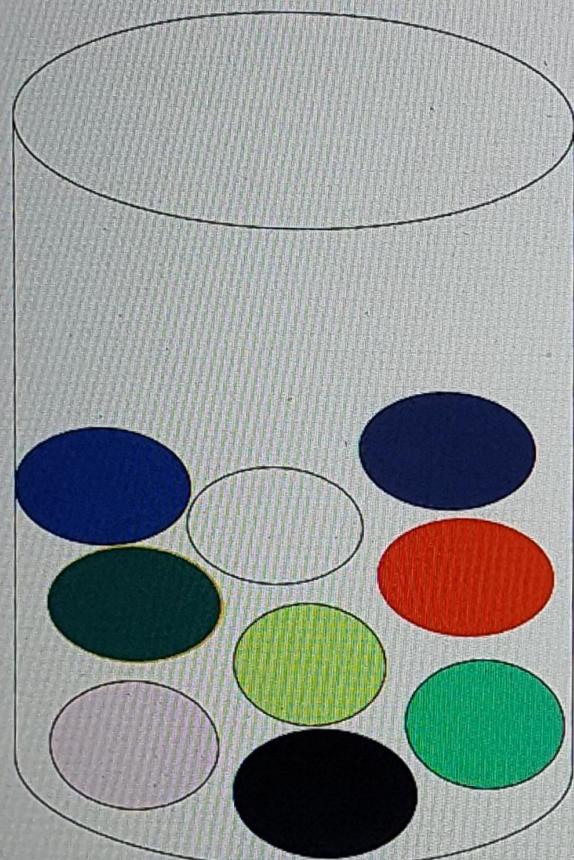
Sampling

- 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)

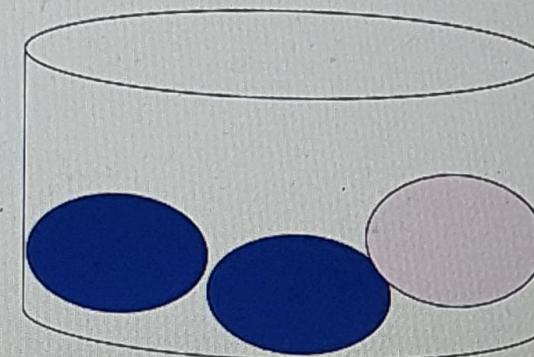
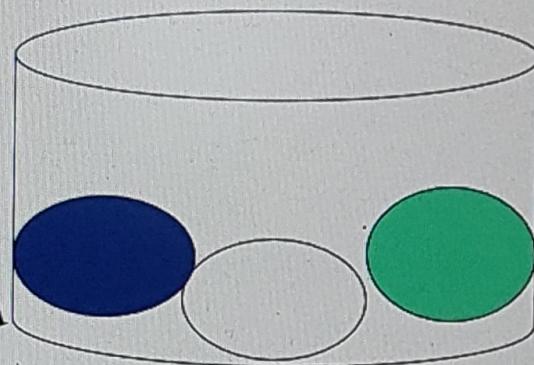
Types of Sampling

- **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

Sampling: With or without Replacement



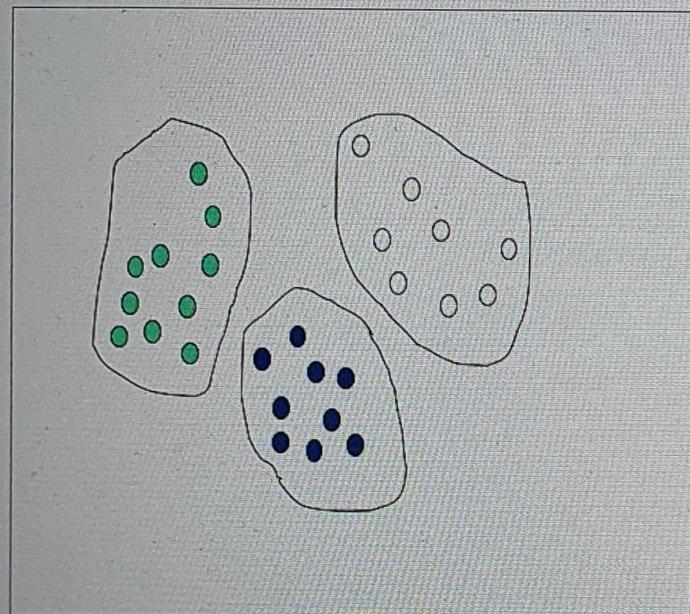
SRSWOR
(simple random sample without replacement)



Raw Data

Sampling: Cluster or Stratified Sampling

Raw Data



Cluster/Stratified Sample

