

Sampling

→ Obtaining a small sample s to represent the whole data set N

□ Simple random sampling

□ Sampling w/o replacement

□ n with n

□ ~~Stratified sampling~~

↳ Partition data

↳ draw samples from each partition.

Data Transformation

→ a func. 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.

methods.

- Smoothing
- Attribute / feature construction
- Aggregation, - summary
- Normalization ; data is scaled
- Discretization,
- concept hierarchy generation

Data Transformation by Normalization

→ measurement unit can affect data analysis.

→ To avoid dependence on the choice of measurement units, the data should be normalized or standarized.

(1)

min max normalization.

perform linear transformation on og data.

13, 15, 16, 16, 19, 20, 23, 29, 35, 41, 44, 53, 62,
69, 72.

Value

for age ~~range~~ to [0:0, 1:0]

$$\text{Transform point } v' = \frac{(v - \text{min})}{(\text{max} - \text{min})} \times (\text{newmax} - \text{newmin}) + \text{newmin}$$

Given range \rightarrow 0:0 to 1:0,Data-transform $v = 45$ $\text{min} = 13$ ~~new~~ new min = 0 $\text{max} = 72$

$$\text{newmax} = 1 \quad v' = \frac{45 - 13}{72 - 13} \times (1 - 0) + 0$$

Transform point. $\rightarrow v' = 0.5423$, [0,1] range.

(2)

Z-score normalization (zero score)

\hookrightarrow transformation is done by values conversion
to a common scale where an avg
number = 0 & standard deviation = 1

$$z = \frac{x - \mu}{\sigma}$$

Z score tells you how many SD from
the mean your score is

ex: mean = 54000 (μ)
 $SD = 14$

$$SD = 16000$$

$$x = \text{Salary} \quad D = 73600$$

$$Z = \frac{73600 - 54000}{16000}$$

$$Z = 1.225$$

③ Decimal scaling

D A = -986 to 917

মাইক্রোপ্লেট টেস্টিং? $\max \text{ abs}(A) = 986$.

→ divide by 1000

$$= -0.986 \pm 0.917$$

Discretization by Binning.

→ top-down splitting technique (bins).

Dissemination by Histogram Analysis!

Same as prev.

Concept Hierarchy Generation for Nominal Data

PPT

Concept Hierarchy Generation

- **Concept hierarchy** organizes concepts (i.e., attribute values) hierarchically and is usually associated with each dimension in a data warehouse
- Concept hierarchies facilitate drilling and rolling in data warehouses to view data in multiple granularity
- 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 *youth*, *adult*, or *senior*)
- Concept hierarchies can be explicitly specified by domain experts and/or data warehouse designers
- Concept hierarchy can be automatically formed for both numeric and nominal data. For numeric data, use discretization methods shown.

Concept Hierarchy Generation for Nominal 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
 - $\{\text{Urbana, Champaign, Chicago}\} < \text{Illinois}$
- 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: $\{\text{street, city, state, country}\}$

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

