

efficient implementation of (numerical data mining, predictive modeling, decision support, pattern analysis etc.)

- frequently used information mining & DMT (parallel computing, distributed) in DBMS

11/10/19 Right coupling

- can be (mostly integrated) into DBMS
- can implement business as organizational component of an information system
- can handle a large data parallel

11/10/19 Concept Descriptions

Data Mining :  Descriptive DM
Predictive DM

DDM - Summary of original dataset
(More & concise view is obtained)

PDM - trying to predict the behavior of your data (\cong ML)

concept - collection of dataset

\Rightarrow Since, datasets are huge and it is difficult to derive any information from them, summarization of data is important

concept description

\hookrightarrow characterization

\hookrightarrow comparison (dissemination)

representing lower level concepts
into higher level concept

Generalization of data using 'age' attribute is much easier than using 'DOB' attribute.

Attributes such as all phone no. can be completely discarded as it serves no visible useful purpose.

- ⇒ OLAP deals with numeric data only (generally), whereas concept description (CD) deals with complex datatype including text & image.

Since, OLAP works with numeric data, aggregation functions such as sum(), avg(), count() etc. are used.

OLAP is user-defined (controlled) where the dimensions & measures of a schema are given by the user. CD is automated & it discovers the useful attributes by itself.

Functions such as roll-up, & drill down are specified by the user in OLAP.

→ based on some threshold (on degree of generalization)

- ⇒ attribute such as gender can never be generalized and attribute such as major (CS, IT, EC etc.) can be generalized to various degrees such as engineering / science based on the

- entropy

Page No:	/ /
Date:	/ /

threshold and requirement

Data generalization and summarization based characterization
Consider a dataset of sales

Sales dataset



item information

(item_id, name, brand, category,
supplier, place, price)

If these lower level details can be generalized into some higher level detail like seasonal sale (say, in Christmas), then attributes such as category, supplier, place & price must be used and analyzed.

AOI (Attribute oriented induction)

Approach proposed for data generalization & summarization in 1989

- based on characterization

AOI

- Online (during runtime)
- Before query processing (thus, slow)

DW

- Offline ↗
 - after query processing (thus, fast)
- precomputed data is stored (data cube)

Steps for ADI

- i) Collect fact relevant data (TRD)
- ii) Perform generalization on TRD
- iii) Apply aggregation
- iv) Interactive presentation to user

if some attribute can be derived from some other attribute, it is removed as there is no point in storing both of them.

attribute removal (irrelevant attributes) which aren't related to TRD

8/8/19 The first step is also called data focusing.

On performing generalization, over-generalizing and under-generalizing can pose previous problems.

<city, state, country>
On dropping state & country attributes leads to under-generalization

DOB, Age
This leads to over-generalization

If large difference between training data and predicted data shows high variance, whereas if it matches to a great extent, it may lead to high biasing.

Generalization \rightarrow Attribute Removal
 \rightarrow Attribute generalization

Attribute can be removed if it has a large set of distinct values (for attribute there is no concept hierarchy defined). The ~~of~~ uniqueness leads to no interesting patterns and nothing much can be predicted. Ex- Names (can differ even in spelling).

Attribute generalization depends on the level of generalization required and the information which needs to be derived.

<city, state, country>

5 entries of various cities of the same state can straight away be removed/reduced by dropping 'city' attribute & keeping a count of these entries (here, 5).

How large should be the no. of distinct values?

Attribute generalization control

\uparrow generalization - over generalization
 \downarrow - under "

May lead to
misclassification

\uparrow
- preferring specialization
- interestingness lost

Process Name	
Date	/ /

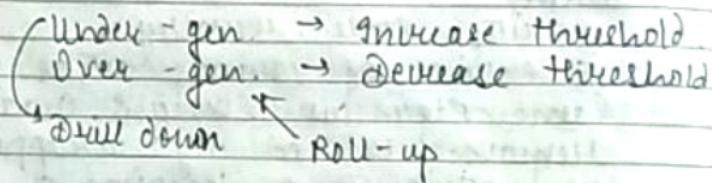
Attribute generalization threshold (contd)

- user decides on the threshold on any attribute you control.

Ex - If a column has (2-8) distinct values then generalization may not be required.

The threshold can be same for all attributes or different for every attribute.

If the values in a column exceed the threshold, attribute removal and generalization should be performed.



Generalized selection threshold (contd)

No. of distinct values are counted per tuple (distinct tuples).

After generalization, how many rows should be present; is answered here.

In older times, 10 to 30 distinct tuples was set as a threshold.

Ex. Initial less relevant data

= Name | gender | major | birth-place | residence |

Phone no. | GPA (attribute) | birth-date |

name Gender Major Birthplace Birth-date Residence phone# GPA
 jim m cs Vancouver 9-18-1997 331,XYZ,PQR 8.4
 (Canada)

On application of ADI:

- 'name' is removed as it does not yield any pattern.
- 'gender' has only two values and thus cannot be generalized further.
- 'major' can be generalized under arts, science, commerce (as per requirement). Done by climbing the concept hierarchy.
- 'Birth-place' can be generalized into state or country based on req. & the threshold value isn't satisfied.
- 'Birth-date' can be converted into 'Age' and then stratified into senior, middle-aged, young etc.
- Street no. & Street generate more specialization & thus, can be discarded.
- Phone no. can be discarded as it does not generate any interestingness.
- 'GPA' needs to be normalized based on the scale & the grading system it uses (out of 4/8/10) and then stratified into { poor, average, good, extraordinary }

After applying all these steps, a 'count' attribute is maintained with every tuple.

9/18/19

After applying ADI

Gender	Major	Birth-Country	Age Range	Residence
			20-30	
			40-50	

GPA	COUNT
good	
excelling	

Algorithm

1. $w \leftarrow \text{get_taskRelevantData}(\text{DMQuery}, \text{DB})$
2. Prepare-for-generalization (w)
 - a. Scan $w \rightarrow$ get count of distinct value a_i - distinct attribute value
 - b. for every a_i : decide whether to generalize or not based on threshold
3. $P \leftarrow \text{generalization}(w)$

→ constructing ^{Datcube} DB on the fly for given DM query
v/s
Predefined Datcube

ON THE FLY

collect TRD

If it is very small and specialized, it is better to go for precomputed datcube.

Roll up & drill down operations are well because we don't drop the original

data everytime

More response time

Predefined

- Granularity ~~data cube~~ matches query
- Suitable when TRD is huge/large
- The computational overhead should not be overlooked

Presentation of derived generalization

Deal with how the data should be displayed after the steps of generalization have been applied.

They can be represented as:

- Cross tab
- Bar chart
- Pie chart
- 3D - cube

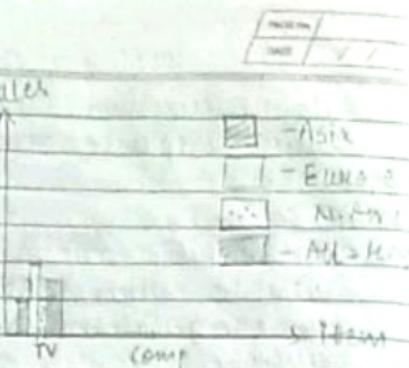
Ex:	location	item	sales	count (in K)
	Asia	TV	15	300
	Europe	TV	12	250
	N America	TV	28	450
	Asia	Cmp.	120	1000
	Europe	Cmp.	150	1200
	N. America	Cmp.	280	1800

optional

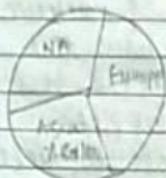
Cross Tab:

location/item	TV	Cmp		Both items		
		Sales	Count	Sales	Count	
Asia	15	300	120	1000	135	1300
Europe	12	250	150	1200		
N America	28	450	200	1800		
All Regions	55	1000				

Bar chart: Sales



Pie-Chart



19/8/19

t-weight (Typicality weight) Almost similar
to Support

$$t\text{-weight} = \text{count}(q_n)$$

$$\sum_{i=1}^n (\text{count}(q_i))$$

n: total
no. of
topics

Range: [0-1]

$\forall x, \text{target-class}(x) \Rightarrow \text{condition}(x) [t; w_1]$

v ...

v $\text{condition}(x) [t; w_n]$

Ex:

$\forall x, \text{item}(x) = \text{"computer"} \Rightarrow$

$(\text{location}(x) = \text{"Asia"}) [t: 25\%]$

v $(\text{location}(x) = \text{"Europe"}) [t: 80\%]$

Page No.	
Date	12/1

$\forall (\text{location}(x) = \text{"North America"}) [t: 45\%]$

↑
1806
4000

After evaluating tree, we can set a threshold, say 26%. Thus, the tuple with t-weight less than the threshold can be discarded.

Analytical characterization : Attribute Relevance Analysis

The attribute should follow covariance ≈ 0 or lesser

& variance should be high.

Inherently, attribute relevance analysis tries to achieve this goal

Ex- color of a car does not matter when predicting price of the car.

By applying some statistical measure, the method aims to discard the weakly relevant / related or irrelevant attribute

For predicting salary, the date & month from DOB can be dropped, an approximate age can be calculated from the year and then divided into groups of decades like 10..20, 20..30 etc.

Depends on which

Ex: If target class is underaged, group becomes
the contrasting class. (when dealing with
these two only)

Target class
contrasting class - tuple which aren't
in target class

Method for Attribute Relevance Analysis (ARA)

- Information gain & entropy.

Expected information.

$$I(S_1, S_2, \dots, S_m) = - \sum_{i=1}^m \frac{S_i}{S} \log \left(\frac{S_i}{S} \right)$$

Entropy:

$$E(A) = \sum_{j=1}^v \frac{S_{ij}}{S} + \frac{S_{2j}}{S} + \frac{S_{3j}}{S} + \dots \frac{S_{mj}}{S}$$

Information Gain:

$$\text{Gain}(A) = I(S_1, S_2, \dots, S_m) - E(A)$$

→ S = samples in training data, i.e. total tuples.

$$\rightarrow \{S_1, S_2, \dots, S_m\}$$

↓
tuples belonging to target class C_1

Probability of
an arbitrary sample from S belonging
to class $i = \frac{S_i}{S}$

Entropy analyzes every attribute available against the target class.

Expected info identifies the no. of tuples which can uniquely describe the target class.

$$0 \leq s_i \leq 1$$

$$\sum s_i = 1$$

$$\log\left(\frac{s_i}{s_j}\right) < 0$$

$$\text{To compensate } I(s_1, s_2, \dots, s_m) = -\sum_{i=1}^m \left(\frac{s_i}{S} \right) \log\left(\frac{s_i}{S}\right)$$

Undergraduate (contrasting class)

gender	major	birth-country	age-range	gpa	count
M	Science	Foreign	≤ 20	very-good	18
F	Business	Canada	≤ 20	fair	20
M	Business	Canada	≤ 20	fair	22
F	Science	Canada	$21-25$	fair	24
M	Engg.	Foreign	$21-25$	very-good	22
F	Engg.	Canada	≤ 20	excellent	24
<u>130</u>					

Graduate (target class)

gender	major	birth-country	age-range	gpa	count
M	Science	Canada	$21-25$	very-good	16
F	Science	Foreign	$26-30$	excellent	22
M	Engg.	Foreign	$26-30$	excellent	18
F	Science	Foreign	$26-30$	excellent	25
M	Science	Canada	$21-25$	excellent	21
F	Engg.	Canada	$21-25$	excellent	18
<u>120</u>					

Expected Info:

$$\begin{aligned}
 I(S_{11}, S_{21}) &= -\sum_{i=1}^2 \frac{S_i}{S} \log\left(\frac{S_i}{S}\right) \\
 &= -\frac{S_1}{S} \log\left(\frac{S_1}{S}\right) - \frac{S_2}{S} \log\left(\frac{S_2}{S}\right) \\
 &= -\frac{120}{250} \log\left(\frac{120}{250}\right) - \frac{130}{250} \log\left(\frac{130}{250}\right) \\
 &= 0.9988
 \end{aligned}$$

22/8/14

To check whether 'major' is a relevant attribute or not.

Here, $v = 3$ [Science, Engg., Business]

Expected info. for Science:

$$S_{11} = 84 \quad (\text{for grad.})$$

$$S_{21} = 42 \quad (\text{for undergrad.})$$

$$\begin{aligned}
 I(S_{11}, S_{21}) &= -\frac{84}{126} \log\left(\frac{84}{126}\right) - \frac{42}{126} \log\left(\frac{42}{126}\right) \\
 &= 0.9183
 \end{aligned}$$

For Engg.

$$S_{12} = 36 \quad (\text{grad.})$$

$$S_{22} = 46 \quad (\text{undergrad.})$$

$$\begin{aligned}
 I(S_{12}, S_{22}) &= -\frac{36}{82} \log\left(\frac{36}{82}\right) - \frac{46}{82} \log\left(\frac{46}{82}\right) \\
 &= 0.9892
 \end{aligned}$$

For Business:

$$S_{13} = 0 \quad (\text{grad.})$$

$$S_{23} = 42 \quad (\text{undergrad.})$$

Page No.	
Date	/ /

$$I(S_{13}, S_{23}) = - \frac{42}{42} \log(1) = 0$$

$$\begin{aligned} E(\text{major}) &= \frac{126}{250} I(S_{11}, S_{21}) \\ &\quad + \frac{82}{250} I(S_{12}, S_{22}) \\ &\quad + \frac{42}{250} I(S_{13}, S_{23}) \\ &= 0.7873 \end{aligned}$$

$$\text{Gain}(\text{major}) = 0.9988 - 0.7873 = 0.2115$$

If threshold had been set to, let say, 0.1 then, this attribute is considered and not discarded.

Attribute relevant analysis steps

- 1) Data collection.
- Determine target and contrasting class
- 2) Preliminary Relevance Analysis using Univariate A.O.I
- 3) Remove irrelevant or weakly relevant attributes using Selected relevance analysis
- 4) Generalize concept description using A.O.I

Ex. 4

(x=2005)

Page No.	/ / /
Date	/ / /

Year	t	y (Sales)	$t \cdot y$	t^2
2005	0	12	0	0
2006	1	19	19	1
2007	2	24	48	4
2008	3	37	111	9
2009	4	45	180	16

$$\sum t = 10 \quad \sum y = 142 \quad \sum xy = 368 \quad \sum x^2 = 30$$

using $y = at + b$

$$\text{Now, } a = \frac{\sum ty - \sum t \sum y}{\sum t^2 - (\sum t)^2}$$

$$= \frac{5(368) - (10)(142)}{5(30) - (10)^2}$$

$$= 8.4$$

$$b = \left(\frac{1}{n}\right) \left(\sum y - a \sum t\right)$$

$$= 11.6$$

$y = 8.4t + 11.6$ can be used to predict any future value.

In the year 2012, Sales = ?

$$\text{Rele, } t = 2012 - 2005 = 7$$

$$\therefore y = 8.4(7) + 11.6$$

$$= 70.4$$

26/5/19

Mining class comparison

Sometimes the target class is of not much importance to the user rather the attributes and their relation may be more relevant.

Computer Science class & Physics class are not related but still are comparable (how students are performing in their respective domains)

Here, the target & contrasting class share some similar dimension or attribute.

Now class comparison is performed?

- i) Data collection
- ii) Dimension relevance analysis
 - The dimensions which are highly relevant are kept & others discarded (done using information gain)
- iii) Synchronous generalization
 - generalize the target class by using the user specified threshold.

If target class is generalized to country level then, contrasting class should also be generalized to the same level. This is called synchronous generalization

- iv) Presentation of derived compilation
 - Excel tables, pie-charts and other representations can be used.

→ d-weight

$$d\text{-weight} = \frac{\text{count}(q_a \in C_i)}{\sum_{j=1}^m \text{count}(q_a \in C_j)}$$

m = total target + contrasting classes
 $d\text{-weight}$ for $q_a \Rightarrow$ ratio of the no.

tuple in target class
 to total no. of tuples in target as well as contrasting classes

Range: [0% - 100%] OR [0, 1]

If the value of $d\text{-weight}$ is ↑ it can be said that the tuple is derived from the target class

If $d\text{-weight}$ is ↓, it can be said that the particular tuple is derived from the contrasting class.

$\forall x, \text{target-class}(x) \Leftarrow \text{condition}(x)$
 $[d: d\text{-weight}]$

status	major	age-range	GPA	count
graduate	Science	21 - 25	good	90
undergrad	Science	21 - 25	good	210

$\forall x, \text{status}(x) = \text{"graduate"} \Leftarrow \text{major}(x) = \text{"Science"} \wedge \text{age-range}(x) = \text{"21..25"} \wedge$
 $\text{gpa}(x) = \text{"good"} [d: 30\%]$

110	240	670	670
100	200	340	340
1	2	17	17

$$\delta = \frac{90}{90+210} = \frac{90}{300} = 0.3 \quad (30\%)$$

Ex: location/item

	TV	computer	both-item
Europe	80	240	320
North America	120	560	680
Both Region	200	800	1000

	TV	computer	Both
location/item	count	count	count
Europe	80	240	320
NA	120	560	680
Both Region	200	800	1000
	$\frac{80}{320} \times 100$	$\frac{80}{200} \times 100$	

t-weight:

$$\forall x, \text{target}(x) \Rightarrow \text{condition}_1(x) [t:w_1] \vee \dots \vee \text{condition}_m(x) [t:w_m]$$

d-weight:

$$\forall x, \text{target}(x) \Leftarrow \text{condition}_1(x) [d:w_1] \vee \dots \vee \text{condition}_m(x) [d:w_m]$$

combining:

$$\forall x, \text{target_all}(x) \Leftarrow (\text{condition}_1(x) [t:w_1, d:w_1] \vee \dots \vee \text{condition}_m(x) [t:w_m, d:w_m])$$

1999

$$\rightarrow \forall x, \text{location}(x) = "Europe" \Leftarrow \\ (\text{item}(x) = "TV") [t:25\%, d:40\%] \vee \\ (\text{item}(x) = "computer") [t:75\%, d:30\%]$$

Probability of the item being TV & coming from Europe is 25%.

If the item is TV, the probability that it was sold in Europe is 40%. (And since here we have only 1 discriminating class (North America) 60% of the TV were sold in NA.)

Measuring the dispersion of data

Method:

- | | |
|--------------------------|-------------------------------------|
| i) Standard deviation | Degree of spread of particular data |
| ii) Inter-quartile range | (How it varies along the axis) |
- Quartiles, outliers & box-plot

$Q_1 = 25$ percentile

$Q_2 = 50$ percentile (median)

$Q_3 = 75$ percentile

$Q_4 = 100$ percentile

IQR - Inter quartile range

$$= Q_3 - Q_1$$

To observe the behavior of the data in the middle of the distribution

→ Particularly useful for skewed data

Using IQR, we cannot say anything

mean	
mode	/ /

about the data residing at extremes.
Thus, 5 point ^{number} summary is used.

Steps:

- | | | |
|------|------------------------|--------------|
| I) | Observe Q ₁ | <u>order</u> |
| II) | " Median | (2) |
| III) | " Q ₃ | (3) |
| IV) | " minimum | (1) |
| V) | " maximum | (5) |

Box-plot (Representation of 5 number summary)

