

Data Mining and warehousing

Converting DBMS to data warehousing - data pre-processing., ETL tool is used (extract, transform, load)

* Why data mining - All DM will be used in ML
There is explosive growth of data from terabyte to petabyte.

→ activity of user - help in recommendation system.

Ex) case study of 2003 is considered, when internet not used so much

* Pattern retrieval chart & data set for - analytical.

* Evolution of Sciences

① Data collection & database creation (before DBMS whatever you learnt) → Primitive file processing.

1960s & earlier. No GUIs, windows. So, config file के अंदर, data store को ही सो, data retrieval को को application programming lang like

C++, Java के program में SQL problem if,

C++ के user के & DBMS के user का access form सो, intermediate mode (user के से, problem, (1970-1980 early))

② DBMS comes - hierarchical & network DS.

Relational DS most common

ER, indexing, SQL, report,

first subquery and
rest products stored (number of tuples)

Page No. _____
Date: _____

query processing, transaction, online transaction process

(3) Advanced Database system (mid ADs - present)
spatial - geoinformative, satellite, sequence, structural
semantic web meaning of attributes you type in google
Page is ranked on the basis of hit ratio.

(4) Advanced data analysis - OLAP → for retrieving
data from data warehouse. OLAP tool uses
classification, clustering, outlier, association
correlation.

multidimensional model of data warehouse -

* In cube, sales are shown $x = \text{item}, y = \text{time}$,
 $z = \text{city}$.

* data is classified in various categories.

Ex) 2003 case study, no net, phone, no OS RAM
with less memory, suppose person start e-commerce
he started 4 major shop in India. He used
barcode scanner. He took frontend .NET & MySQL
in backend, developed interface & database
dataset (item ID, product name, salesperson). 4
CITE software started,

In 2005, frontend nightmare, backend SQL
server. Product ID, attribute = barcode. tabularized
after scope of analytics com. summarized
report dataset, 2 databases with 2 databases
in ETL link with it, convert it to flat
data at

data mining concept & techniques

Page No. _____
Date: _____

multiple datasets with their help find in which
region sale is less

1000x1000x10 products per day so create on server side
summarized info or result is stored in
separate store which is called data warehouse
(no write can be done, only read). Only these
attributes stored which are important for decision
making.

* temporal - time factor like weather
forecasting data
* spatial → space related, sequence - weather forecasting,
all frames in sequence.

* Data mining → (knowledge discover from data)
Processed data → information,
combining multiple info together & finding pattern
→ knowledge.
→ extraction of interesting

* Data mining as step in process of knowledge discovery

* Basic flow of ML in data analytic part -

- (1) databases (multiple datasource ex. SQL, SAP, ERP) -
delta fetch from sources
- (2) Data warehouse (ETL data 3D or 4D) -
only sale is stored in cube after labelling is done

Page No. _____
Date: _____

Stone sale in data warehouse.

Relational database

ETL is used to transform database to data warehouse.

→ Data warehouse → 3 types of schema is stored,

→ we need sale, item, time, location.

Select count of item(id) from table where city = ' ' and time = ' '.

✓ Sale is stored in this way (multi-dimensional).

Quarter = 3 month.

3 array is divided by 1 away for Jan other care of Feb. March

100 is sale of 40, 100, 150
Jan Feb March

Jan also divided into 4 week, E-T-I-O
Week also divided into daily basis

* Data sale daily base it segregate among all sale at different abstraction level & store one by one storing sale in different level

* Data warehouse always single value stored in 1 cube.

Page No. _____
Date: _____

Database of info metadata in size 3GB with help of metadata with error

We clean data first, last 5 years after 5000 years & remove data now we keep only city, time, item, sale.

+ Now, in datawarehouse → count dimension all values

② Datawarehouse is OLAP to implement this so that patterns find & using data mining algo

③ On basis of patterns, knowledge base is created

* Roll up from city → countries if convert drill down - quarters month → month → 3(3)

* Decision making historical data is used to make product selling

* Normal database is knowledge is you & thought

① Database (different sources are there) IDW

② Data warehouse data from different source databases.

ETL used to transform database to data warehouse

Schemas in DW → Star schema & fact table
→ Snowflake fact table
→ Repetitive Dimension table
Act. constant

- depend on us
- | | |
|----------|--|
| Page No. | |
| Date: | |
- ~~Cube it city of concept hierarchy. (n1311)~~
~~city store base it is scale etc~~
~~city break into street.~~
- * DW is subject oriented.
 - (1) Data mining algo \rightarrow patterns found in DW operate on DW (OLAP tool get needed) on basis of pattern we develop pattern, knowledge.
 - * multidimensional link \rightarrow find pattern (n1311)
- lecture 3
- * Determining algo deployed in machine learning lang using R, Python.
 - + mining algo operate on -
 - (1) Data stream (youtube) & sensor data (weather forecasting)
 - (2) Temporal
 - (3) Spatial data
 - (4) multimedia database
 - (5) Text database, newly Generalization 1UNI3
 - * Data Mining Function - (1) Class/concept description
- Data mining categorize table into -
- operating query on dw
- | | |
|----------|--|
| Page No. | |
| Date: | |
- (1) descriptive \rightarrow characterize properties of data in target dataset
 - (2) Predictive \rightarrow predict data. It's induction in order to make predictions.
 - * whatever algo we create are object oriented conaff, real life problem we're solving
 \rightarrow works on basis of attendance
 \rightarrow If we summarizing dataset of class \rightarrow data characterized
 \rightarrow 1 class at same time compare with \rightarrow discrimination.
 - * Suppose Coke, Pepsi are selling. If sale of coke then Pepsi. this association can be find out on basis of discrimination.
 - (2) mining frequent Patterns, \rightarrow on e-commerce you do shopping, suggestions are made
- market Product 1 - most sale 2 - less sale so, they keep combo
so, they kept items ① & ② together
so that ③ can get easily sold.
- * Transactional data if item occurrence. count of item, count of item 1 & count of item 2.
 - * Subsequence \rightarrow item type occurred 2nd item is
 \rightarrow 1st user laptop then fine webcam for memory card.
 - * frequent pattern \rightarrow 2nd 21/8 sell milk & bread.
 - * substitution \rightarrow milk \leftarrow full cream less cream. In railway station

in all cities use full cream milk, all juice shops near hospital for all city → Substructure.
include frequent set as subsequent

→ support & confidence both in frequent set.
out of 100, 50 milk & bread buy it (then
confidence)

- 3.21 association rule \Rightarrow del \Rightarrow \Rightarrow
for min item observe min - min threshold

③ classification - (supervised) class level $\bar{y}_N \pm$
using ML train tree we predict all patterns used $\bar{y}_N \pm$
cural
non classif clustering - unsupervised.

After 5 attributes (e.g. 10, 12, 13 etc.), communication skills, technical, placement → we train model on basis of these attributes then we predict at what company you will be placed.

class level - output
final & final cc at upper & i. loop ch(120)
such always at first placed of

↗ **yes** (107.1.7 CS) internal node - condition with
branches - outcome ↗
yes (127.7 CS) class A / class B - spot
yes (376.4.7 CS) ↗
② Neural (numerical data)
hidden layer there, weight ϵ

DT Cluster analysis - when no class label, we need to find pattern, we use that
→ we apply colour to all on basis of similarity,
dissimilarity we make clusters

* Kmean algo \rightarrow with clusters तो (we select 3 dots) representative. then we calculate euclidean distance of this point from all the points ($n-1$).
at close will be 1 boundary हो जाएगा calculated mean value then make cluster वह again calculated distance) new cluster will change.

- 1 class to 3162 minimize inter class
maximize similarity in class
 - outside cluster elements are outliers // fraud detection 1. unauthorized account

in all cities use full cream milk, all juice stops near hospital for all day → substructure.

include frequent set as subsequent

→ support & confidence factor in frequent set.
out of 100, 50 milk & bread buy it (then
confidence → supports 50% (actual if we buy milk buy the bread)
confidence set with 50% next day)

- Bayes classification rule $\hat{y} = \arg \max_{\hat{y}} P(\hat{y}|x)$ →
for min error observe max - min threshold

③ classification - (supervised)
in train we predict class level $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n$
in test we predict $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n$
clustering - unsupervised

After 5 attributes (e.g., Block %, communication skill, technical, placement → we train model on basis of these attributes then we predict at what company you will be placed)

class level - output
final & final as if upper & lower limit
and always if first placed of

if we don't have class level then it comes in unsupervised learning

Categorical Numeric

Page No. _____
Date: _____

internal node - condition with
branching outcome \hat{y}_1
class A/B/C - root
yes (100, 1, 7CS)
no (120, 7CS)
yes (100, 1, 7CS)
no (120, 7CS)
yes (100, 1, 7CS)
no (120, 7CS)

② Neural Numeric data, \hat{y}_1
Hidden layer three, weight \hat{w}_1

④ Cluster analysis - when no class label, we
need to find pattern, we use this
→ we apply column no col on basis of similarity,
dissimilarity we make clusters

• Kmean age → with clusters and (we select 3
dots) representative. then we calculate Euclidean
distance of this point from all the points ($n-1$).
at class still with 1 boundary if \hat{y}_1 is first
calculated mean value then make cluster here.
again calculate distance) now, cluster will
change.

- 1 class \hat{y}_1 minimize inter class
maximize similarity in 1 class
- outside cluster elements are outlier / fraud
detection!, unauthorized access

- * If numerical value - production (manufacturing)
- categorical (alphabet) - decision tree
- * outlier → object that does not comply with general behaviour of data
- * sequential pattern mining -
- * Data warehouse usage

- Two data warehouse applications
- ① Information processing
 - ② Analytical processing
 - ③ Data mining
- From OLAP to OLAM
- ① High quality of data in data warehouses
 - ② Available info processing
- Whenever we have to fetch data from warehouse, we use OLEDB.

OLAP system architecture

- ① Relational OLAP (ROLAP) - close to relational DBs
 - ② Multidimensional OLAP (MOLAP) - storing multidimensional array. If we aggregate data then many models over time. Multiple abstraction levels are there, all data stored at different abstraction level.
 - ③ Hybrid → combo of ROLAP & MOLAP
 - ④ Specialized SQL servers - tool depend upon company, company specific domain at first at tool provide analysis
- We were using Data warehouse system to fetch data from DB & ETL → ROLAP
 - Specialisation - At first at different level it store and ETL ex UG → BTech → CSE, Civil.
 - Relational database → multidimension

p-90

Unit 2

Data Preprocessing

- ① Data cleaning - 1 Numerical
- ② Data integration - 2 Numerical
- ③ Data reduction (data file is very large, hold that data which is helpful)
- ④ Data transformation - 2,3 numerical
 - ↳ normalize, convert hierarchy.

we use ETL tool to implement all (4) part.
transform data

Information, multiple companies provide ETL

ETL → relational data at factual & descriptive
5th convert.

- ① fill in missing value, smooth noisy data (binning)
remove outliers & resolve inconsistencies.

- ① whenever dataset is combined from sources, data is faulty.
 - incomplete
 - noise → errors at data file
 - inconsistent → rating 1,2,3 now A₁B₁C₁
 - intentional → as enemy's data.
- Data is not available always.

* How to handle missing data

- ② ignore tuple where data missing.

- ③ fill missing value manually.

- ④ fill it automatically with
 - global constant (not proper)
 - attribute mean (total mean / n → overall % $\frac{\text{count}}{\text{total}}$)
 - attribute mean for all samples belonging to same class. e.g. for class 10% of grade A = 8th
 - overall 10% of avg student
 - class #1 fill for first

- ④ most probable value → fits in classification model

- Noisy data

Noise - random errors or variances in measured variable.

at table at B attribute different at data CTE at

→ To handle noisy data

- ① Binning
- ② Regression
- ③ Clustering
- ④ combined computer & human inspection

- * Data Transformation Strategies →
 - we transform data from 1 form to other.
 - transform transactional data → factorize & dimension data
 - (value to coefficient)
 - ① Smoothing → Techniques binning, regression, clustering (value to cluster)
 - remove noise from data
 - clustering → datapoint that are not in particular cluster are noise & are removed. 1 median point.
 - Ex → boys hostel, girls, localities, pg. are of different clusters. Those who are out of clusters are outliers.
 - 1 cluster - whole feature of that datapoint
 - pattern of value does not reflect individual.

- ② attribute construction → new attributes are constructed, added from given set of attributes
- Ex) sale attribute
- ③ aggregators - In database, sale is derived according to location, annually..
sale is stored on basis of time also.
means storing on different abstraction level.

numerical

④ Normalization - attribute data are scaled so as to fall within smaller range.
→ Large range will transform small & smaller range \Rightarrow
Ex) graph \Rightarrow 1 pm drawn 1 cm in nb.

(i) min-max normalization -

Ex) attribute A value is 100, 200, 300, 400.

min_A = 100

max_A = 400

Transform these values in range of -2 to +2
using $v_i' = \frac{v_i - \text{min}_A}{\text{max}_A - \text{min}_A} \times 2 + 2$

new min_A = -2, new max_A = +2

$v_1' = -2 - 2$

$v_1 = 100 - 400$

$v_1 = 100, v_2 = 200, v_3 = 300, v_4 = 400$

Put these values in formulae

$$v_i' = \frac{v_i - \text{min}_A}{\text{max}_A - \text{min}_A} (\text{new max} - \text{new min}) + \text{new min}$$

$$v_1' = \frac{100 - 100}{400 - 100} (0.8 - -2) + -2$$

$$v_1' = 0 + -2 = -0.2$$

$$v_2' = \frac{200 - 100}{400 - 100} (0.8 - -2) + -2$$

$$= \frac{100}{300} \times 0.6 + -2$$

$$= 0.4$$

$$V_3' \quad V_4' \quad \text{calculate}$$

$$V_3' = \frac{300-100}{300} (-0.6) + 0.2$$

$$= \frac{200 \times -0.2}{300} + 0.2$$

$$V_4' = \frac{400-100}{300} (-0.6) + 0.2$$

$$\boxed{V_3' = 0.8}$$

Mean

Geant

(i) Z-score normalisation \rightarrow when large dataset is
large from min, max value may differ much more than
we use this.

mean
normalisation

use also mean value ; standard deviation

$$\text{mean of } \bar{x} = 100 + 200 + 300 + 400 = 250$$

$$\sigma_A = \sqrt{\frac{1}{n} ((100-250)^2 + (200-250)^2 + (300-250)^2 + (400-250)^2)}$$

$$\boxed{\frac{V_i' - \bar{V}_i - \bar{A}}{\sigma_A}}$$

$$V_1' = \frac{100-250}{\sigma_A}, \quad V_2' = \frac{200-250}{\sigma_A}$$

$$V_3' = \frac{300-250}{\sigma_A}, \quad V_4' = \frac{400-250}{\sigma_A}$$

Book 8 Data Mining concepts & Techniques

Page No. _____
Date: _____

- * Data reduction → data size and reduce transactional database → warehouse data reduce
Only important data is true
- (i) dimensionality reduction -
 - 1 value (cell) listed with 3 dimension
If you need only 1 dimension rest have to be discarded i.e. dimensionality reduction.

- obtain reduced data smaller in volume
here maintain integrity of dataset
- process of removing attributes
whenever you plot graph, you reduce from 3D to 2D
major pattern, sufficient is same but closer & scaled
- * Some all numbers in dataset are not responsible for pattern generation, we took those only which can reflect

- (ii) Numerosity reduction → replace original data with small forms of data
 - ↳ parametric
 - ↳ nonparametric → histogram, cluster

- format it → data that represent are set
- parameters store data & not data
 - ex → linear regression in store for x & y not x, y .
- storing data but in reduced form.

- (iii) Data compression → reduced or compressed representation of data
 - ↳ lossless → original data from compressed data
 - ↳ lossy

* Dimensionality Reduction → 3rd component used
↓ 3rd reduce it.

(1) Wavelet Transform -

discrete wavelet transform (DWWT) -

We have data vector X , we converted it into X' , of wavelet coefficient. For this multiply data vector with wavelet coefficient i.e. $\mathbf{x} = (x_1, x_2, \dots)$

→ higher degree term n^2 smaller coefficients remove. e.g. $n^2 + 2n + 1 \Rightarrow n^2$ pattern can be accessed.

* In ML, packages like Wavelet, Daubechies, Daubechies 6 for reduction of dimension.

algo

step 1) length L , of input data vector must be integer power of 2, this is met by padding data vector with zeros as $(L \geq n)$

2) Each transform in vector 2 function,

1 → data smoothing (mean or avg)
2 → weighted difference which pick detailed features of data

as we do in optimise series ques (find pattern)

→ length get reduced by 2.

(3) functions are applied to pairs of data point in X
1 side → low-frequency version of X
2 side → high frequency.

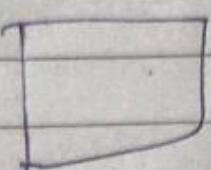
* If data in matrix form, then all data must be in orthonormal (90°)

→ highest coefficient $\approx 90^\circ$ matrix

Page No. _____
Date: _____

→ used for higher dimension
Date: _____

Ex)



matrix find sum 182 value, 284 value from another matrix.

find difference 182, 384 another matrix

In this way 1 single value from pattern

(K-L method).

(2) PCA (Principal Component Analysis) - converting data into orthogonal form

x_1 x_2 pattern cutting with y_1 & y_2 components spread \rightarrow reduce x_1 & x_2 as 2nd basis spread \rightarrow 1st

(1) select k vectors that can best use to represent data.
suppose n vectors then reduce to best needed i.e. k vectors.

step 2) If data are normalized, so that each attribute falls within same range.

can compute

→ used for lower dimension.