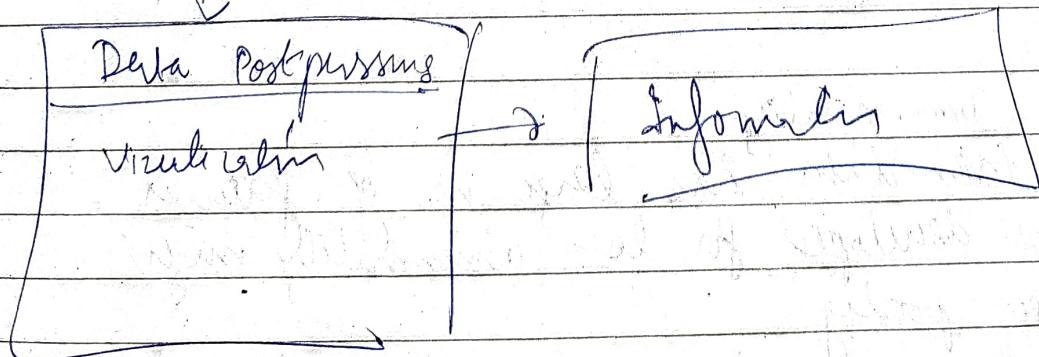
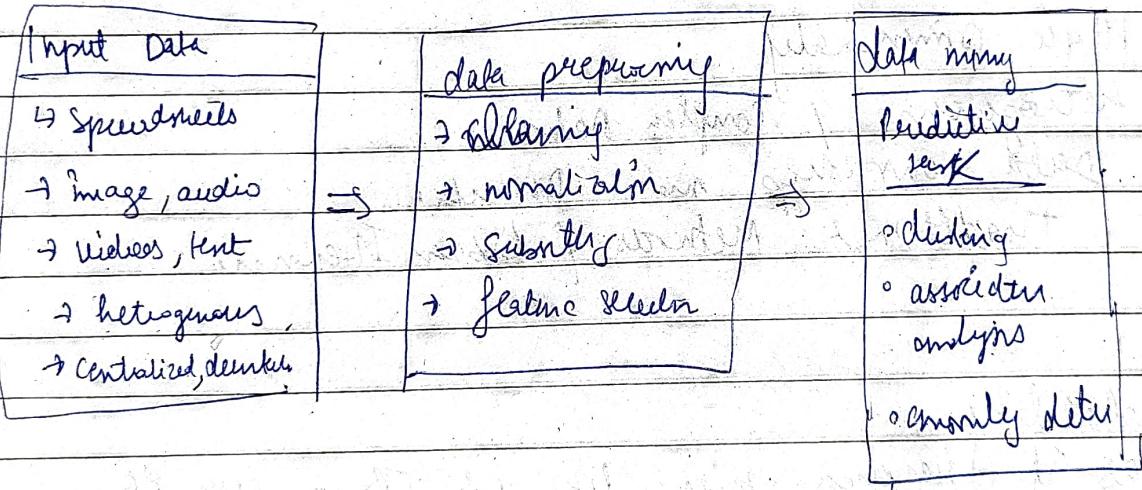
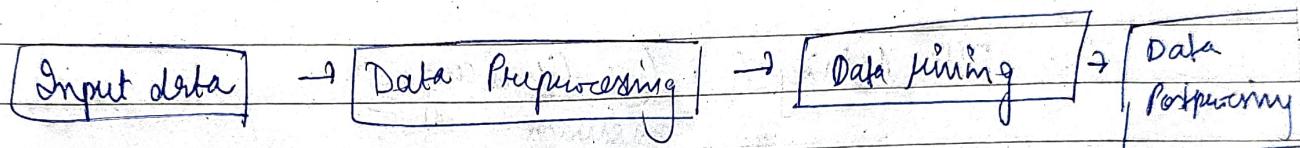
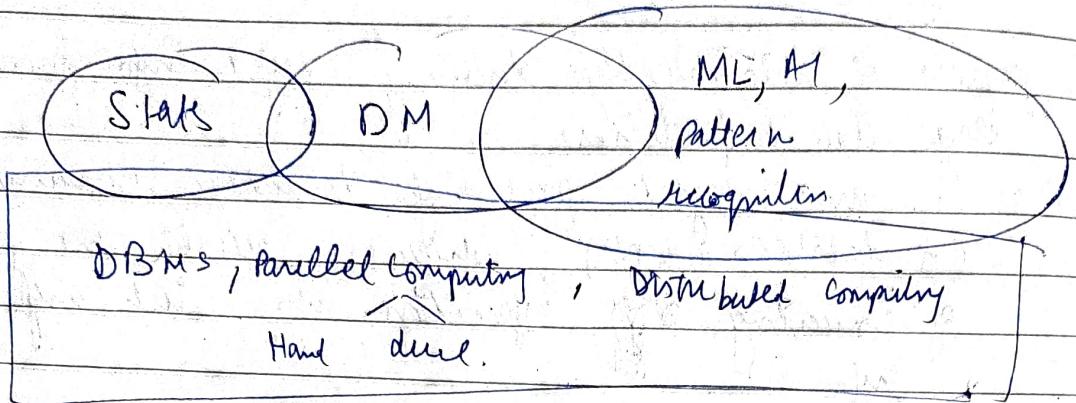


1.1 What is Data Mining?

- process of automatically discovering useful info from large data repositories
- it is an integral part of Knowledge Discovery in databases (KDD) which is overall the process of converting ~~useless~~ information into useful info.



h3 Origin



h2 Why Data Mining? (Challenges)

- Scalability
- High Dimensionality
- heterogeneity / complex Data
- Data ownership and Distribution
- Traditional methods f Labor Intensive

① Scalability

as it reaches sizes like GB, TB or even PB.
so algo's must be scalable enough to handle
these volumes.

② High Dimensionality

as the data has large no. of features.
Algo's developed for low dimensionality might
perform poorly.

③ Heterogeneity and complex Data

Traditional data deals with large no. of
features. ~~Algo's developed for low dimensionality~~

1 attributes of same types (continuous/kategorical)
 so techniques must be able to handle heterogeneous
 attributes and complex objects.

(4) Data ownership and distribution

necessary data may not be stored centrally /
 owned by one organization, requiring development
 of distributed DM techniques.

(5) Traditional Methods

Instead of starting with one hypothesis we let
 data suggest thousands of pattern (rules, clusters,
 correlations)

(6) Data Mining Tasks

Prediction

Descriptive tasks

To predict the value of a particular attribute (target / dependent variable) based on values of other attributes

attribute to be predicted

other

target / dependent variable

attributes used for making prediction

explanatory / independent variable

1 core data mining tasks.

Predictive

Modeling

Classification

Regression

anomaly

detection

association

analysis

clustering

Predictive modeling

building a mathematical / statistical model that can predict an unknown (target) variable based on known values (explanatory variables)

target variable → outcome we want to predict
(Yes/No)

explanatory variable → the factors / inputs that influence (predictor, features) the target

It is of 2 types.

Classification : target variable is discrete
for eg → spam/not spam
species of a flower.

Regression : target variable is continuous (numerical value)
for ex → predicting house price

DT.

Association analysis

discovers co-occurrence patterns among items
in features

expressed A \Rightarrow B

Cluster analysis

groups observations into clusters so that points within a cluster are more similar to each other than to points in other clusters

Anomaly (Outlier) detection

identifies observations that are significantly different from the majority (crime / abnormal patterns)
e.g.: credit card fraud detection

anomaly detection

Supervised

unsupervised

when we do have labels when we don't have reliable labels (fraud vs legit, attack vs normal)

goal	Supervised learning
data required	predict a known label labeled data
task	classification, regression

unsupervised learning
discover hidden structures
unlabeled data

clustering, association analysis,
dimensionality red., density
spiral

1) diff b/w DM & KDD.

2) what is unsupervised learning, explain it with the help of an example

② ~~Diff~~

2

①

Data

Quantitative
(numerical)

Qualitative
(categorical)

spatial data time series text
sequence image

② Data Quality (how accurate, complete, consistent, representative the data is)

Common issues

Noise (random error)

Outliers (rare/extreme occurrences)

Missing data (blank/placeholders like 0 / "missing")

Inconsistency/duplicates (conflicting/repeated rows)

Bias/unrepresentative (sample doesn't reflect the population)

③ Preprocessing (to make data suitable for mining)

Transforming (reshaping data to fit a technique and improve results)

Scaling (z-score)

$$x' = \frac{(x - \bar{x})}{s_x}$$

Min max normalization

$$[0, 1]$$

Discretization

turn numeric into bins (like low / med / high)

Binarization / encoding (categorical into binary)

Aggregation / sampling (fewer rows for speed, higher level summaries)

Dimensionality reduction

4) Analyze via relationships

replace raw objects with their pairwise proximity and analyse those.

2.1 Types of Data

data object is a single record / can be row / instance / vector / sample.

attribute = column / feature / dimension.

• Attribute.

measurable characteristic of an object that can vary across objects over time.

Ex: eye color, temperature.

• Measurement Scale

a rule that maps the true property to a numeric / symbolic value.

attribute type

Nominal (categorical)	Ordinal (categorical)	Interval (numerical)	Ratio (numerical)
Values are names / labels	Values have order / rank	differences play an imp role.	diff) and ratios are imp.

Discrete vs Continuous vs Binary.

Discrete \rightarrow finite / countable values

Binary \rightarrow special case of discrete with two values

Continuous \rightarrow any real value

Asymmetric Attributes and Sparsity

Asymmetric attributes
only presence (non-zero) is informative; absence (0)
is less meaningful

Sparse data.

most values are 0.
attributed are of same type and asymmetric.
non zero are considered important

Resolution

Data can be collected at multiple levels
(fine vs coarse)

more detail

but more noise

Smooth but

may hide patterns

If resolution is too fine then pattern may not be visible or may be confused into noise whereas if pattern is too coarse then it may disappear.

~~Properties of Dataset~~

General characteristics of Dataset that has significant impact on DM.

1) Dimensionality

Curse of Dimensionality

It happens when we have too many features (dimensions) in our data.

for ex 1D graph: data is pretty close

2D graph: Same no. of points feel more spread out

3D: they are farther apart.

100 dimensions ~~are~~ covering feels equally far from everything else.

2) Sparsity It saves computational time & storage requirements

3) Resolution.

Types of Dataset

Record Data

- Market Basket
- Data matrix
- Sparse Data Matrix

Graph based

- Relationship among objects
- Objects form a graph

Order Data

- Sequential / Temporal Data
- Sequence
- Time series
- Spatial

Record Data

- fixed schema, each record has same attributes
- stored in flat files

Market Based

each record is a set of items bought together

Date Matrix / Pattern Matrix

Rows = objects

Columns = numeric attributes

have same fixed number of numeric attribute.

Sparse Data Matrix

Attributes are of same type

(1)

(2)

data objects are mapped to the nodes of the graph

Relationship among the objects are captured by the links b/w objects and their properties such as directions and weights

Graph Based :

relationship among objects
nodes = objects, edges = links.

Objects may be graphs

each object has internal structure (atoms/bonds)

Ordered Data :

Sequential / Temporal data

- there is no timestamp but there are positions in an ordered sequence
 - ordered symbols without sequence (DNA sequences)
- (en)
- sequence of words or letters
 - sequence of nucleotides in a genetic seq data could inform if a person is healthy or not.

Time series :

- measurements over time (stock prices, monthly temp)
- here a series of measurement are taken over time.

(en)

temp during diff years.

Temporal autocorrelation

If two measurements are close in time then their values are often similar to each other.

Spatial

- measurement by location (and time)
- objects have spatial attributes such as position/areas

(e.g.) weather data collected for diff geographical locations

Spatial autocorrelation

objects that are physically closer tend to be similar in other ways as well.

handling non-random data

non random data is where order or relationship is important

General characteristics affect analysis

- dimensionality • no. of attributes
 - curse of dimensionality \rightarrow motivation
 - dimensionality reduction
- distribution frequency / concentration of values.
- resolution level of detail, affects visibility of pattern and noise.

2.2 Data Quality

• DM uses data collected for other purposes
we must detect / correct issues

measurement and data collection issues

↳ us values might differ from true values

• Measurement error

Mixed value + true value

• Data collection error

missing / extra objects / inappropriately including
some data objects

Spatial

Occurs in spatial (by location) or
temporal (timestamp) date

Artifact deterministic alteration of the data

(e.g.) 8 trees in the same place in a set of photos

Noise

random error (sensor fluctuation)

Hard to fully remove.

Precision

closeness of repeated measurements (to each other) of same quantity.

it is measured by standard deviation

$$s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

Bias

it is diff b/w the mean of the set of values and known value of q/y being measured.

systematic deviation from the true value
Goal is to use as many digits to represent the result as one justified by precision of data.

Mass of 1 gm weighing 5 times.

$$\{ 1.015, 0.9901, 0.991, 0.996, 0.931 \}$$

Bias 0.399

Mean = 0.991.

Outliers / Anomalous objects

- Objects / values that are very different from the bulk.
- Values that are unusual / considerably different from the typical values of that attribute.
- They are not the same as noise; they can be legitimate and interesting.

Missing Values

- Info is not collective because people decide to share
- that attribute may not be applicable to all his censys.

Handling Missing Values. (Viva)

- Eliminate the data object / attribute, ~~extending~~ ~~existing~~ missing values
- ignore the missing values ~~as before~~

Inconsistent values

~~zipcode~~

conflicting entries

- for ex :- zip code of a city is correct but it is not appropriate for a particular location
- height of a person cannot be negative

↳ correction of inconsistency requires additional info from a master database.

Duplicate Data

↳ same object appearing multiple times

- Care must be taken while merging data from heterogeneous sources
- there could be two objects that represent a single object

- avoid accidentally combining data objects that are similar but not duplicates

④ deduplication

↓

it requires solving attribute discrepancies while avoiding merging distinct entries with similar info.

Issues related to the application

- timeliness
data becomes stale, refresh required.
- relevance
data must contain relevant variables needed for the goal.
- sampling bias
sample doesn't represent the whole population
- documentation / knowledge about data
a good data dictionary prevents serious misuses

↓

(types, units, encoding, missing-value codes, known relationships, etc.)
the quality of data

2.3 Date Preprocessing

- > Aggregation
- > Sampling
- > Dimensionality reduction
- > Feature subset selection
- > Feature creation
- > Discretization and binarization
- > Variable transformation

Aggregation

Combining two or more objects into single attribute (object)

Purpose

↳ data reduction & reduction in no. of attributes
↳ objects

↳ Change of scale : monthly view of data instead of a daily / per item view.

↳ more stable data : aggregated data tends to have less variability

Advantages

- Smaller datasets after data reduction requires less memory & processing time and hence aggregation may permit the use of more expensive data mining algorithms
- aggregation can act as change of scope / scale by providing high level view instead of low level view

disadvantage

- potential loss of interesting details

Sampling → analyze a representative subset instead of whole data.

- Using sample can work almost as well as using the entire data if the sample is representative.
- It is used for DMR because processing the entire data is too expensive and time consuming.

Sampling

simple random stratified program.

with without

replacement replacement

Simple random sampling (Equal probability of selecting any particular item)

Sampling without replacement Sampling with replacement

as each item is selected it

is removed from the set

of all objects that constitute population.

↳ objects aren't removed from population as they are selected from sample

↳ same object can be picked more than once

↳ simpler to analyze

Spiral

disadvantages of sampling (simple Random)

when populations consist of different types of objects with widely different no. of objects, simple random sampling can fail to adequately represent those types of objects that are less frequent.

Solution

stratified sampling

- Starts with pre-specified groups of objects.
- equal no of objects are drawn from each group even though groups are of diff sizes
- the number of objects drawn from each group is proportional to the size of that group.

Progressive / adaptive sampling

- sample size can be difficult to determine
- these approaches start with a small sample and increase the sample size until a sample of sufficient size has been obtained
- while this advantage eliminates the need to determine the correct sample size initially, it requires that there be a way to evaluate the sample to judge if it is large enough.

Dimensionality Reduction

↓
make data fewer features by projecting data into a lower - dimension - space

advantages

- avoids curse of dimensionality
- eliminate irrelevant features and reduces noise
- reduce amount of time and memory required by data mining
- allows data to get visualized easily
- allows data to be more understandable

Techniques

- ↓
- by creating new attributes that are a combination of old attributes
 - by selecting new attributes that are a subset of old is known as feature subset selection / feature selection

feature subset selection

Keep only useful original features
(drop irrelevant ones)

techniques

- Some irrelevant & redundant attributes can be removed by using domain knowledge
- Try all possible features as input to data mining algorithm.

Standard approaches

- embedded approaches
 - feature selection occurs naturally as a part of the data mining algorithm
 - algo itself decides which attribute to use and which to ignore
- filter approaches
 - features are selected before data mining algo runs.
 - pre select general criteria
- wrapper
 - search subset using final model

Feature creation

↳ Creating a new attribute from existing ones

why? raw inputs are mostly poor shaped

Methods

- 1) Feature extraction (domain transforms)
- 2) map data to a new space
- 3) feature construction

feature extraction : summarizing messy raw data into useful signals

for ex] → we turn pixels into images
words into sentences → text
Raw sound to audio

map data to a new space

↳ frequency data to reveal hidden pattern

fourier transform (time \rightarrow frequency) wavelet transform (patterns over time)

feature construction

use simple formulas on existing features to make better ones

Discretization and Binarization

Discretization

transforming continuous attributes to categorical attributes.

Steps

- 1) Choose split points
- 2) map any value to its interval label.

unsupervised

(no labels)

Supervised

(use label)

- class info is not used
- equal width $\rightarrow w = \frac{\text{max} - \text{min}}{k}$
- equal frequency $\rightarrow n/k$
- clustering
- class info is used
- pure intervals, entropy based

Data

$0, 4, 12, 16, 14, 19, 24, 26, 28$

Width $= \frac{\text{max} - \text{min}}{K}$

$$\frac{28 - 0}{3}$$

$$= 9.333 \approx 10$$

here no of bins is 3 $(0-10) [10, 20] [20, 30]$

$[10, 20]$

$[20, 30]$

eg. say $s \in [0, 10] \rightarrow B_1$

$\frac{n}{k}$

Data $\{ 2, 4, 6, 7, 10, 12, 15, 18, 20, 25 \}$

$$\frac{25 - 2}{3} = \frac{23}{3} = 7.67$$

~~Standardization~~

Binaryization

- Conversion into binary
- binarize categorical attributes

Angry	0	0	0	0
Poor	1	0	0	1
OK	2	0	1	0
Good	3	0	1	1
Great	4	1	0	0

Standardization

$$\text{mean} = 0 \quad \text{SD} = 1$$

($6+2$)

If diff variables are to be combined in some way then it's necessary to avoid having a variable with large no. values.

Min-Max Normalization

$$V' = \frac{V - \text{min}}{\text{max} - \text{min}} \quad (\text{new max} - \text{new min}) + \text{new min}$$

marks

8

0

10

$$\frac{10-8}{20-8} \rightarrow \frac{2}{12} = 0.16$$

15

$$\frac{15-8}{20-8} = 0.38$$

1.2

20

$$\frac{20-8}{20-8}$$

$$\text{Max} = 20 \quad \text{Min} = 8$$

2. Scale normalization

$$\frac{\sum v_i}{n}$$

$$v' = \left(\frac{v - \mu}{\sigma} \right) = \frac{\sum v_i}{n}$$

$$\sqrt{\sum (v_i - \mu)^2}$$

n

Mesures of similarity and dissimilarity

→ Similarity b/w two objects is the numerical measure and is usually non negative.

It is often b/w 0 and 1.

where

0 → no similarity

1 → high similarity

→ Dissimilarity is the numerical measure of how the objects are diff from each other

- ° It is lower for more similar objects.
- ° Distance is used to find it.

Minkowski Distance



$$d(x, y) = \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{\frac{1}{r}}$$

n = no. of dimensions

r = no. of parameters

y $r = 1$,

L_1 (Manhattan) = sum $|x_k - y_k|$

y $r = 2$

L_2 Euclidean = $\sqrt{\sum (x_k - y_k)^2}$

y $r = \infty$,

L_∞ supremum = max $|x_k - y_k|$

n	y
P_1	0
P_2	2
P_3	0
P_4	5

$$(s-2) + (1-0)$$

(L2)

$$\sqrt{(s-2)^2 + (1-0)^2} \approx \sqrt{5}$$