DATA MINING AND WAREHOUSING

CS 402.

28/06/2021
Monday.

christy Varghese

② ⓑ  $A = \{116, 234, 486, 544\}$.

min - max normalization

new_$min_A = 0$

new_$max_A = 1$.

$min_A = 116$

$max_A = 544$.

$$V_i' = \frac{V_i - min_A}{max_A - min_A} (new\_max_A - new\_min_A) + new\_min_A.$$

$$V_{116} = \frac{116 - 116}{544 - 116} (1 - 0) + 0 \quad = \frac{0}{\underline{\underline{\phantom{0}}}}$$

$$V_{234} = \frac{234 - 116}{544 - 116} (1 - 0) + 0 = \frac{118}{428} = \underline{\underline{0.276}}$$

$$V_{486} = \frac{486 - 116}{544 - 116} \times 1 = \frac{370}{428} = \underline{\underline{0.8644}}$$
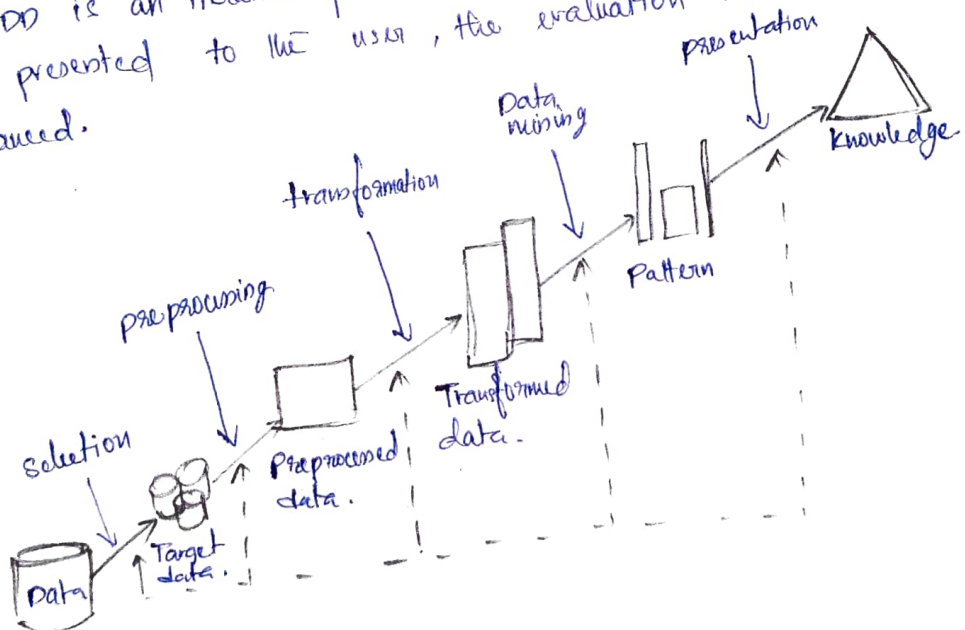
$$V_{544} = \frac{544 - 116}{544 + 116} \times 1 = \underline{\underline{1}}$$

① ⓐ Knowledge Discovery in Databases (KDD).

✦ Data cleaning also known as data cleansing, it is a
phase in which noise data and irrelevant data are removed
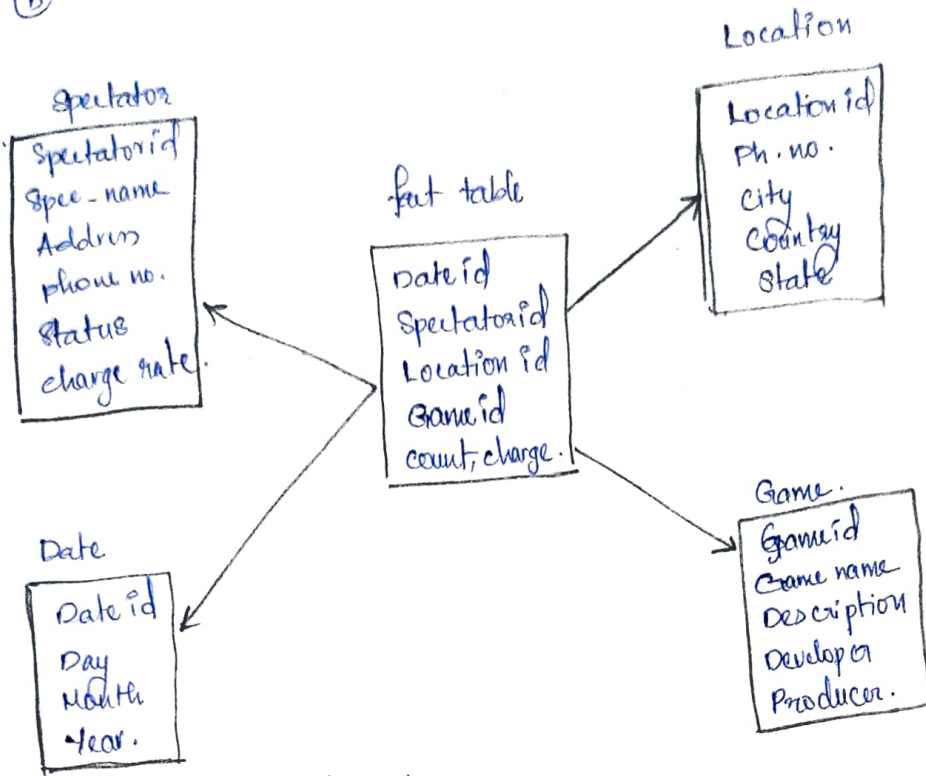from the collection.

christy

※ Data integration at this stage multiple data sources, often heterogeneous, may be combined in a common source.

※ Data mining it is the crucial step in which clever techniques are applied to extract patterns potentially useful.

※ Data selection at this step, the data relevant to the analysis is decided on and retrieved from the data collection.

※ Data transformation it is a phase in which the selected data is transformed into forms appropriate for the mining procedure.

※ Pattern Evaluation strictly interesting patterns representing knowledge are identified based on given measures.

※ Knowledge representation is the final phase in which the discovered knowledge is visually represented to the user.

※ Data selection & data transformation can also be combined where the consolidation of data is result to the selection.

※ KDD is an iterative process. Once the discovered knowledge is presented to the user, the evaluation measures can be enhanced.

① . ⓑ

**Spectator**

| Spectatorid |
|---|
| Spec-name |
| Address |
| phone no. |
| Status |
| charge rate |

**fact table**

| Date id |
|---|
| Spectatorid |
| Location id |
| Game id |
| count, charge |

**Location**

| Location id |
|---|
| Ph. no. |
| city |
| country |
| state |

**Game**

| Game id |
|---|
| Game name |
| Description |
| Developer |
| Producer |

**Date**

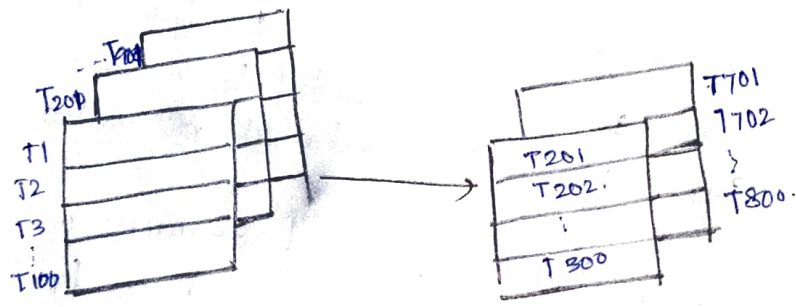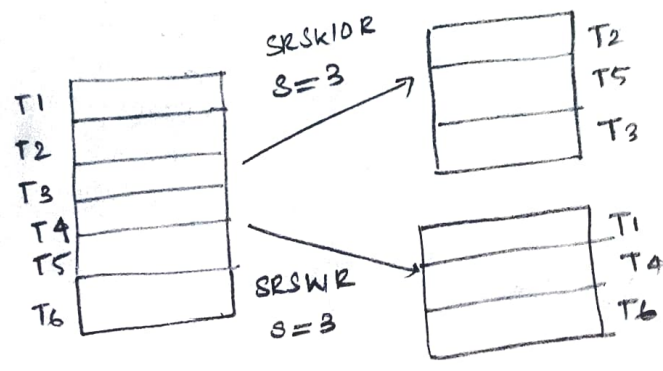| Date id |
|---|
| Day |
| Month |
| Year |

Star schema.

② 
ⓐ Sampling can be used as a data reduction technique because it allows a large data set to be represented by a much smaller random sample. Suppose that a large data set D, contains N tuples.

(i) Simple random sample without replacement (SRSWOR) of size s:
This is created by drawing s of the N tuples from (D≤N) where the probability of drawing any tuple in D is $1/N$, i.e; after all tuples are equally likely to be sampled.

(ii) Simple random sample with replacement (SRSWR) of size s.
This is similar to SRSWOR, except that each time a tuple is drawn from D, it is recorded and then replaced.
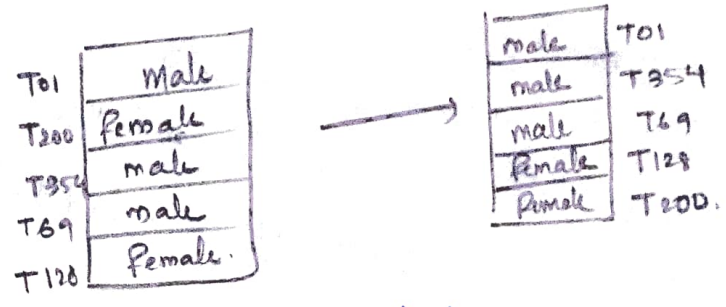i.e; after a ~~tex~~ tuple is drawn, it is placed back in D. so that it may be drawn again.

(iii) cluster sample : If the tuples in D are grouped into M mutually disjoint "clusters", then an SRS of s cluster can be obtained, where $s \leq M$.

(iv) strafied sample : If D is divided into mutually disjoint parts called strata, a straified sample of D is generated by obtaining an SRS at each stratum. This helps ensure a representative sample, especially when data are skewed.



SRSklOR
s=3

SRSWR
s=3

T1 T2 T3 T4 T5 T6

T2 T5 T3

T1 T4 T6



T200
T1
T2
T3
T100

T201
T202
:
T300

T701
T702
$\}$
T800.

cluster
sample (s=2)



| T01 | Male |
| T200 | Female |
| T354 | male |
| T69 | male |
| T128 | female |

| male | T01 |
| male | T354 |
| male | T69 |
| female | T128 |
| female | T200 |

straatified
sample
decording to male, female..

⑦

(b)  $A = (22, 1, 42, 10)$  $B = (20, 0, 36, 8)$.

(i) ~~Hamill~~ Euclidean distance.

$$D_E(A, B) = \sqrt{(22-20)^2 + (1)^2 + (42-36)^2 + (10-8)^2}$$

$$= \sqrt{4 + 1 + 36 + 4}$$

$$= \sqrt{45}$$

$$= 6.708$$

(ii) ~~Ha~~ Manhattan distance

$$D_M(A, B) = |22-20| + |1-0| + |42-36| + |10-8|$$

$$= 2 + 1 + 6 + 2$$

$$= 11.$$

⑦  (a)  minimum support = 3       confidence = 80%

| TID | items_bought |
|-----|-------------|
| T100 | {M; O; N, K; E, Y} |
| T200 | {D, O, N, K; E, Y}. |
| T300 | { MA, K, E} |
| T400 | { M, U, C, K, Y}. |
| T500 | { C, O, O, K, I, E}. |

$c_1 \rightarrow$

| Itemset | count | | |
|---------|-------|---|---|
| M | 3. | D | 1 |
| O | @3. | A | 1 |
| N | 2 | U | 1 |
| K | 5 | C | 2 |
| E | 4 | | |
| Y | 3 | | |

Since the minimum count is 3.

∴ L will be ⟹

| Itemset | Support count |
|---------|---------------|
| M | 3 |
| O | 3 |
| K | 5 |
| E | 4 |
| Y | 3 |

cluster 1.

{M, O, K, E, Y}.

$C_2$ ⟹

| itemset | support count |
|---------|---------------|
| (M, O) | 1 |
| (M, K) | 3 |
| (M, E) | 2 |
| (M, Y) | 2 |
| (O, K) | 3 |
| (O, E) | 3 |
| (O, Y) | 2. |
| (K, E) | 4 |
| (K, Y) | 3 |
| (E, Y). | 2 |

$L_2$ ⟹

| itemset | support count |
|---------|---------------|
| (M, K) | 3 |
| (O, K) | 3 |
| (O, E) | 3 |
| (K, E) | 4 |
| (K, Y) | 3. |

$C_3$ ⟹

| itemset | support count |
|---------|---------------|
| (M, O, K) | 1 |
| (M, K, E) | 2 |
| (M, K, Y). | 2 |
| (O, K, E) | 3 |
| (O, K, Y) | 2. |

c.Christy

$L_3 \Rightarrow$ | itemset | support count
| $(O, K, E)$ | 3.

∴ The frequent itemset $(O, K, E)$ is obtained by Apriori Algorithm.

min. confidence = 80%

consider $L_3 = (O; K, E)$

Association rule formed from $(O, K, E)$.

$\{O, K\} \Rightarrow E$  confidence $= \dfrac{3}{4} = 100\%$

$\{O, E\} \Rightarrow K$  confidence $= \dfrac{3}{3} = 100\%$

$\{E, K\} \Rightarrow O$  confidence $= \dfrac{3}{4} = 75\%$.

i. Strong association $\Rightarrow$

$E \Rightarrow \{O, K\}$  confidence $= \dfrac{3}{4}$

$K \Rightarrow \{O, E\}$  confidence $= \dfrac{3}{5}$

$O \Rightarrow \{K, E\}$.  confidence $= \dfrac{3}{3} = 100\%$.

∴ strong association $\Rightarrow$  $\{O, K\} \Rightarrow E$
$\{O, E\} \Rightarrow K$
$O \Rightarrow \{K, E\}$

④

⑧  P(headache = Y | runny nose = N | fever = Y) = ?

$$P(flu = Y) = \frac{1}{2}$$
$$P(flu = N) = \frac{1}{2}$$

$$P(headache = Y | flu = Y) = \frac{2}{3}$$
$$P(headache = Y | flu = N) = \frac{1}{3}.$$
$$P(runnynose = N | flu = Y) = \frac{1}{3}.$$
$$P(runnynose = N | flu = N) = \frac{2}{3}.$$
$$P(fever = Y | flu = Y) = \frac{2}{3}.$$
$$P(fever = Y | flu = N) = \frac{1}{3}.$$

$$P(headache = Y | runny nose = N | fever = Y | flu = Y)$$
$$= \frac{2}{3} \times \frac{1}{3} \times \frac{2}{3} \times \frac{1}{7} = .074$$

$$P(headache = Y | runnynose = N | fever = Y | flu = N)$$
$$= \frac{1}{3} \times \frac{2}{3} \times \frac{1}{3} \times \frac{1}{7} = 0.037.$$

max (0.074, 0.037) = 0.074

So the patient is likely to have flu = Yes.

|   | −ve | +ve |
|---|-----|-----|
| −ve | TN | FP |
| +ve | FN | TP |

④ ⑩.
⑤ ⑥  $$Recall = \frac{TP}{TP + FP}$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}.$$

| | +ve | ~ve |
|---|---|---|
| +ve | 100 TP | 5 FN |
| -ve | 10 FP | 50 TN |

① Accuracy = $\dfrac{100 + 50}{100 + 50 + 5 + 10}$

Accuracy = 0.909

② Recall = $\dfrac{TP}{TP + FN}$ = $\dfrac{100}{\cancel{110}\,105}$ = 0.952

④ ⓐ **Working of classification**

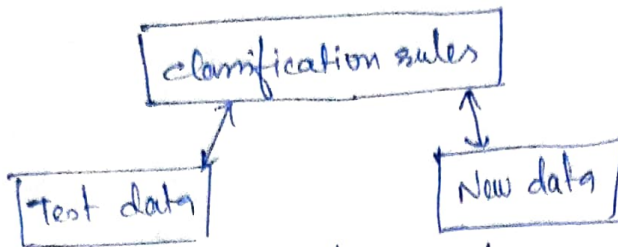Data classification is a two-step process.

**step 1:** In the first step, a classifier is built describing a predetermined set of data classes or concept. This is the learning step, where a classification algorithm builds the classifier by analysing or "learning" from a training set.

**Step 2:** In the second step the model is used for classification.

first, the predictive accuracy of the classifier is estimated. If we were to use the training set to measure the accuracy of the classifier, this estimate would likely be optimistic. It is also known as the classification step

The accuracy of a classifier on a given test set is the % of the test tuple that are correctly classified by the classifier.

```
   Training          →  | classification |
   data set  |---·------→ | Algorithm      |
                                |
                                ↓
                          | classification |
                          | Rules          |
```

step1: learning step

```
                    ┌─────────────────────┐
                    │ clanification rules │
                    └─────────────────────┘
                      ↗                  ↑↓
          ┌───────────┐          ┌───────────┐
          │ Test data │          │ New data  │
          └───────────┘          └───────────┘
```
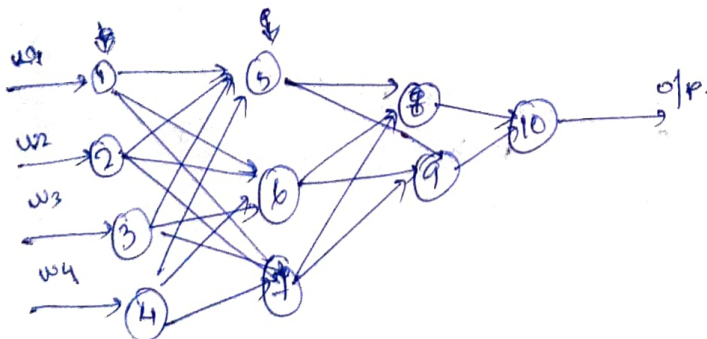
step 2: classification step.

(5) (a) Backpropagation Algorithm

Backpropogation learns by iteratively processing a data set of training tuples, comparing the network's prediction for each tuple with the actual known target value. The target value may be known class label of training tuple. For each training tuple, the weights are modified so as to minimize the mean squared error between the network's prediction and the actual target value. These modifications are made in "backwards" direction. The algo steps involved are expressed in terms of $y_{ps}$, o/ps and errors and may seem awkward if this is your first look at neural network learning.
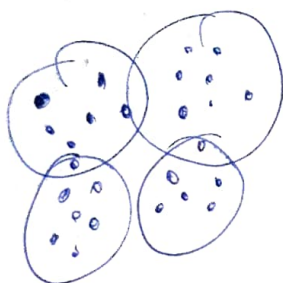
⑧ ⓐ

Ⓐ **Agglomerative Method**

Agglomerative methods work by starting with singleton sets and then merging them untill S is covered.
The agglomerative methods cannot be used directly as it scales quadratically with the number of data points. However Hierarchical methods usually generate spherical clusters clusters and not of arbitrary shapes.

Each objects initially represents a cluster of its own.
Integrating hierarchical clustering with other integrating hierarchical clustering with other techniques are BIRCH, CURE. etc.
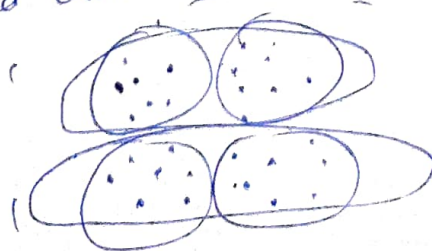


It is a bottom up approach

Diagram

Ⓑ **Divisive Method :** It works by recursively partitioning the set of data points S until singleton sets are obtained.

All objects intially belong to one cluster. Then the cluster is divided into sub-clusters which are sucsievely divided into their own sub-clusters. This process continues untill the desired cluster structure is obtained.
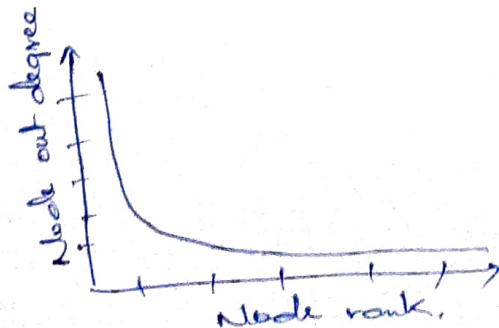
It is a top-down approach.



Diagram

**⑧ ⓑ Characteristics of social networks**

Social networks are rarely static. Their graph representation evolve as nodes and edges are added or deleted over time. In general, social networks tend to establish the following performance.

(i) **Densification power law :** It was believed that as a n/w evolves, the number of degree grows linearly in the number of nodes. However, extensive experiments have shown the n/w become increasingly dense over time with the average degree increasing.

(ii) **Heavy tailed out degree & indegree distributions :** The no. of out-degrees for a node tends to follow a heavy-tailed distribution by observing the power law, $1/n^a$, $a < 2$. The smaller the value of heavier the tail. The in degree also follows a heavy tailed distribution although it tends to more skewed than the out degrees distribution..

(iii) **Shrinking diameter :** It has been experimentally shown that the effective diameter tends to decrease as the network grows. This contradicts an earlier belief than the diameter slowly increase as a function of n/w size.



*y-axis: Node out degree, x-axis: Node rank.*

*(signature)*