

1/7/19  
Data mining - use old data to find some (knowledge discovery from data) pattern among them.

Data - Structured (RDBMS)

Semi-structured (Social media sites)

Non-structured (Big data)

(Unstructured)

any file (audio, video, text etc uploaded).

3/7/19

## Data PreProcessing

Collected data can be:

- Inconsistent
- Noisy
- Incomplete

Can occur due to the tendency of humans to keep some information private. This can lead to missing the incomplete data in most cases.

Noisy data refer to the irrelevant data punched in. Ex- Entering 'abc' in salary field. Hence, data validation is a must.

Inconsistent data refers to the values of fields which are inconsistent with each other. Ex. Age is 23 and DOB is 19/7/2018 which is impossible.

So we can say that it is not the quality of data or quality of mining result.

Thus, to avoid redundant entries,  
data must be preprocessed.

### Bon Ternouï Principle

Some events are very rare. Even then, if the event has to be considered to increase its probability, the dataset should be enhanced.

Based on this idea, Bon Ternouï principle states that more the data, higher are the chance to detect some rare patterns.

$\Rightarrow \begin{array}{l} A - 10 \\ \text{Apple} - 20 \\ AP - 15 \\ B - 21 \end{array}$  } represent same item. Hence, they need to be transformed after integration from various source

If transformation is not applied, the answer comes to be B, which is incorrect.

Data reduction - keeping only the relevant data which is related to the analysis to be performed but without affecting the mining result

### Dimensionality reduction

Ex: Text summarization deals with extraction & abstraction

to take out the lead & understand the entire text then prepare a

thus, to achieve an optimal solution, taking a subset of the dimensions (subset of data) is required. Like - not considering attributes such as age, DOB, name, address for analysing the sales trend is appropriate and in fact, required.

- Wavelet transformation

- Data cube aggregation

$\Rightarrow$  PCA deals with eigen values

$$AX = \lambda X$$

matrix  $\uparrow$        $\lambda$   $\uparrow$  eigen value. reduction  
eigen vector

data compression (part of)

### Data Preprocessing

- Data cleaning / Inconsistency removal.
- Data integration
- Data transformation
- Data reduction

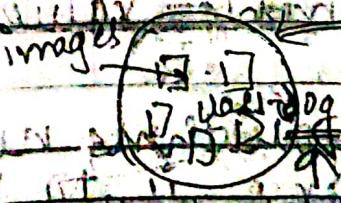
8/7/19

### Data Cleaning:

- Handling missing values :

- i) ignore the tuple

Not ideal but, suitable for problem of class label.



Here, in case of missing data. The tuple can be ignored if we don't know what it is.

ii) fill data manually  
possible for very small dataset.  
the missing data is filled manually  
based on heuristic predictions

iii) use a global value constant to  
fill missing value.  
Ex - putting in (keeping) global constant  
as: NaN, unknown, ..

Some interesting pattern may be  
lost if freq. of this constant  
increases.

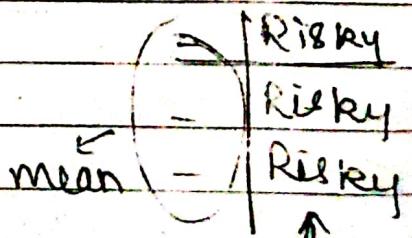
iv) Replace missing value by mean  
(mode can also be used)

v) Use attribute mean of same class

Ex:  

Good	Good
Bad	Good
if some value is missing, it	Good

  
Good from same class  
if some value is missing, it  
can be replaced by the  
mean of the values  
of same class.



v) Replace most probable value in the  
missing tuple.

Classification algorithms need to be  
trained and inference is to be carried  
out.

⇒ (iii) to (vi) can lead to biasing or biased data set, because they involve mean

### - Noisy data:

Ex - Age / Salary having -ve values

Random noise by variance obtained if a variable is noise.

## ii) Binning

$$\text{L.C.M.} = 2^3 \times 3^2 \times 5^1 \times 7^1 = 2520$$

Sort the data.

(Already sorted)

Partition into (equi depth) bins:

-- Bin 1 : 4, 8, 15

Bin 2: 21, 21, 24

Bin 3: 25, 26, 34

## Smoothing by mean

Bin 1: 9, 9, 9

Bin 2: 22, 22, 22

Bin 3: 29, 29, 29

## Smoothening by bin boundary

Bin 1: 4, 4, 15

Bin 2: 21, 21, 24

Bin 3: 25, 25, 34

If bin 1: 4,710

then, it becomes

bin 1: 4, 4, 15, 11

## ii) clustering

grouping the similar items.

~~Passing of unsupervised data~~

The ~~frequency~~ which mostly is considered as noise) and is simply dropped.

iii) Combining computer and human inspections.

It involves clustering at the first step and then combining the outliers. Human heuristic is used to classify them further.

#### iv) Regression

Regression  
Selecting the best fit line which covers most of the points.

The points which are far away from this line are considered as outliers

## Data Inconsistency

Ex: DOB Age  
1911/2018 23

91719

## Data Integration

Problem may arise due to data inconsistency

Ex: Mr. Zuckerberg I refer to same  
Mr. Mark person but, when  
combined, they may  
be considered as  
separate entity.

Henie, early recognition  
of his genius

## Introduction to relation analysis

If any info. is redundant, it can be easily ignored.

It may be possible that the other

value

attribute will not derived directly exactly. But, if both the attr. have the same trend, one of them can be dropped.

Mean (Average) : Biased towards the smallest or largest data values. Might not lead to correct prediction for a individual data.

$$\bar{Y} = \frac{\sum_{i=1}^n x_i}{n}$$

Median : gives the middle value for a dataset.

Mode : gives the data value having the highest frequency.

Mean, median, mode are all measures of central tendency.

$\Rightarrow$  Weighted average :  $\bar{Y} = \frac{\sum_{i=1}^n w_i x_i}{n}$

Variance : Deviation from mean value.  
(rough measure) can be given by :

$$\sigma_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

Standard deviation :  $\sigma_x = \sqrt{\text{variance}} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}}$   
(more exact measure)

$$Y_{A,B} = \frac{\sum (A - \bar{A})(B - \bar{B})}{(n-1) \Delta_A \Delta_B}$$

$\rightarrow$  Co-variance

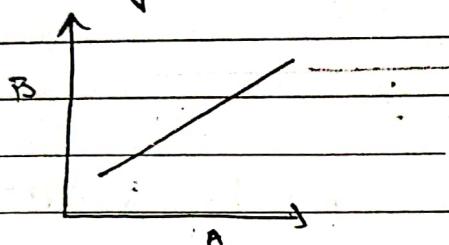
$r_{A,B} \rightarrow$  Pearson's correlation coefficient

If  $r_{A,B} > 0$  +ve co-relation

$= 0$  Independent of each other

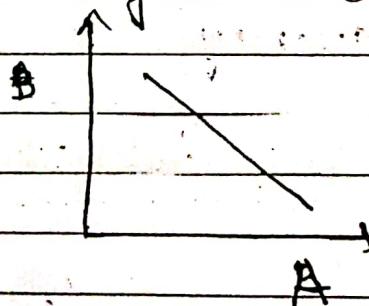
$< 0$  -ve co-relation

+ve-ly co-related :  $A \uparrow B \uparrow$



$r_{A,B} = 1$ , then they are strongly co-related

-ve-ly co-related :  $A \uparrow B \downarrow$



Chi-square : (Given by Pearson)

$$\chi^2 = \sum_{i=1}^n \sum_{j=1}^n \left( \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \right)$$

( $O_{ij}$ ): Observed value

$E_{ij}$ : Expected value

## Data Transformation

- Smoothening
- Aggregation
- Generalization (Based on granularity level  
the data is put into  
generalized labels)
- Normalization (in case of numeric data)
  - MIN-MAX Normalization

Ex: Age  $\rightarrow$  normalize in range  
50  $\rightarrow$   $[0, 1]$

$$\text{find } \frac{24}{26} = \frac{\text{new\_max}}{\text{new\_min}_A}$$

$\therefore$  for every value,  $v$ :

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

$$\Rightarrow \text{min} = \$12,000$$

$$\text{max} = \$98,000$$

$$N.R = [0.0, 1.0]$$

Normalize value:  $\$73,600$

$$v' = \frac{73600 - 12000}{98000 - 12000} (1.0) + 0.0$$

$$= \frac{61600}{86000} = 0.716$$

~~•  $Z = \frac{x - \mu}{\sigma}$  (Standardization or Z-mean)~~

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

~~Helpful when min & max. value  
are given with no explicit~~

~~range is not known. E.g. data given~~

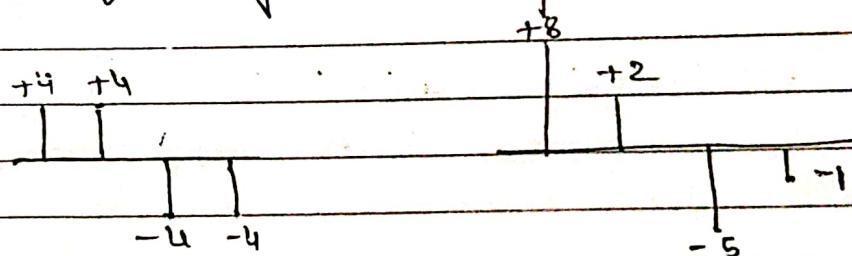
11/11/29

$$\text{Variance} = \frac{\sum (x - \bar{x})^2}{n} \quad \leftarrow \text{entire population}$$

$$S.P = \sqrt{\frac{\sum (x - \bar{x})^2}{(n-1)}}$$

When we are considering a subset of the population.

If we have uniform data, variance works perfectly correct. But, when the data is skewed, standard deviation give a more exact answer, whereas variance is quite fluctuating.



$$\text{variance} = +4$$

$$\text{variance} = +4 \quad \text{which is not very accurate}$$

Decimal scaling

$$v' = \frac{v}{10^j}; \quad j \text{ is smallest Int such that } (\max |v'|) < 1 \quad (\max |v|) < 1$$

Ex- Range (A) : -986, 917

max. absolute value of A = 986  
new - 0.70 (minimizer,  $v' = 2.3$ )

$\therefore -986$  becomes  $-0.986$

③ Standardization

In addition reduced variation of the original

dataset, without compromising the integrity of the data.

$$\text{Linear regression : } \hat{y} = w_1 x + b$$

weight  $\downarrow$  bias  
 original

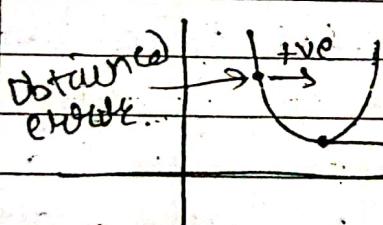
$w$  &  $b$  are simple parameters here.  
To find hyper-parameters,  $w$  &  $b$  we need to find:

$$\min \left[ \frac{\sum (y_i - \hat{y}_i)^2}{n} \right]$$

so as to (almost) correctly predict the trend in  $x$  &  $y$ .

Hence, to find one parameter, we need another parameter.

$$w' = w - \alpha \frac{\partial}{\partial w} (\text{adjusted error})$$



hyper-parameter / learning rate

Thus, we need to min. error. add some value to adjust the error.

15/7/19

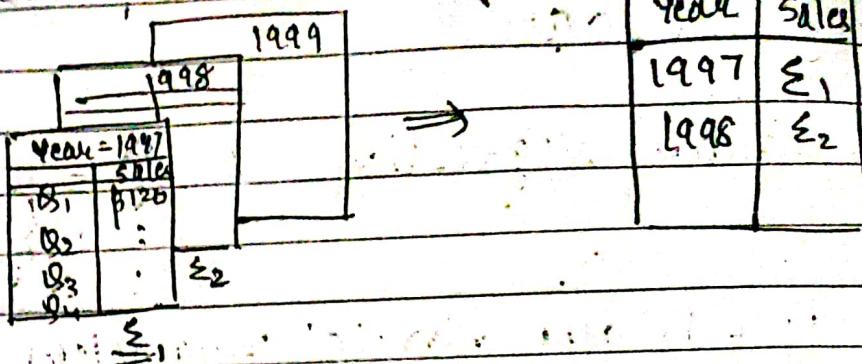
Techniques for data reduction :

i) Data Cube Aggregation

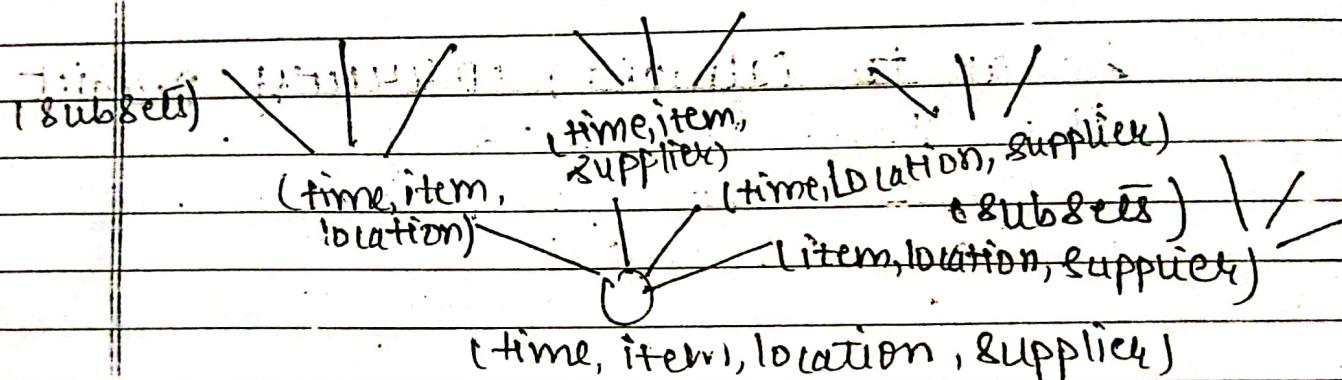
	Branch A	Branch B	Branch C	Branch D
Employees	500	450	400	350
Computers	120	100	80	60
Phones	600	500	400	300
Scanners	720	600	480	360
Total	3000	2500	2000	1500

this slice represents data for branch A, made 8 years)

## Compress dimensionality:



## Concept hierarchy:



From these various levels of abstraction, we can abstract out the info relevant to us and thus assist in mining mining the big. data.

ii) Dimensionality Reduction: speed  
The computational power and storage required to maintain high dimension data is quite large

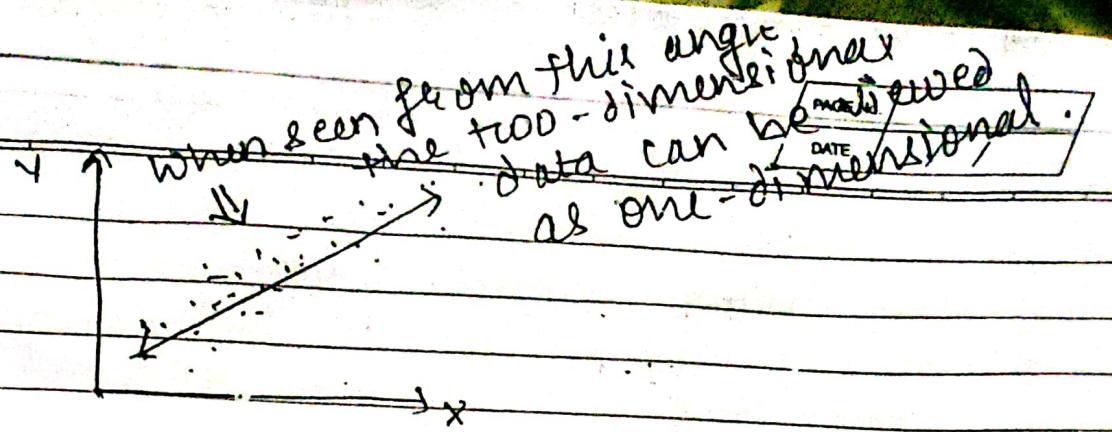
High dimensionality data also contains irrelevant and redundant attributes

Irres.

Red.

A dimensional data is also

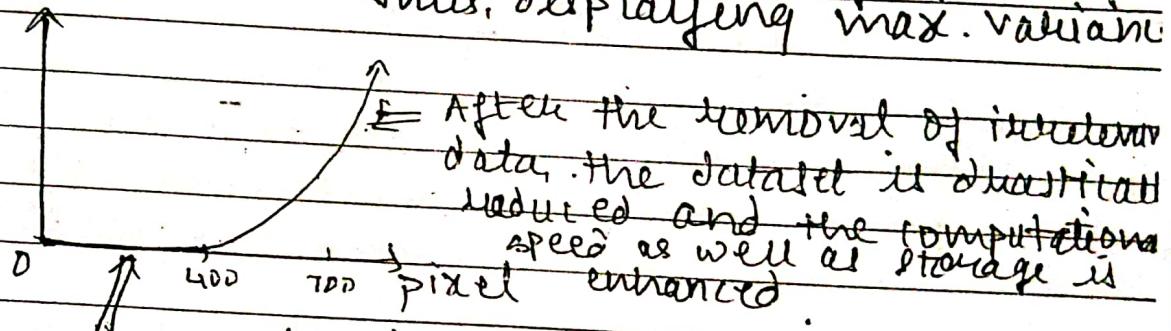
8.2907 MB



⇒ The dimension having the maximum variance is the most important.

white.
Black

This line of separation represents the highest frequency of change in intensity for images. Thus, displaying max. variance.



This portion has almost negligible variance and thus, removing it from the data set does not affect the quality of the image. The mining result stays unaffected.

⇒ The data distribution (in the probability space) of the original dataset and the data obtained by taking a subset of attributes should remain same. Then we can say that the original property of the data is preserved.

⇒ Exhaustive search (go on creating subsets until it's feasible).

A, B, C  
Take subset of 2

A	B
C	D

Find co-variance:

Combine all the co-variance values to form a co-variance matrix. Finding the matrix's determinant, we can know how related two attributes are.

Ex: 
$$\begin{vmatrix} 1 & 1 & 3 \\ 1 & 2 & 2 \\ 2 & 2 & 5 \end{vmatrix} = D$$

(pair)  
same column. Thus, these two attributes are said to be strongly dependent.

iii) Step wise forward Selection:

- $\{\}$       ii) Start with empty set
- $\{A_1\}$     ii) Add the attribute having highest variance
- $\{A_1, A_2\}$     ii) repeat (2)

Results are checked after every addition and when the desired results are obtained, the process is stopped and we get the subset of attributes which are most relevant for our need.

iv) Step wise backward Selection:  
-  $\{A_1, A_2, \dots, A_n\}$

- $\{A_1, A_2, A_3, \dots, A_n\}$   $A_2 X$
- $\{A_1, A_3, \dots, A_n\}$   $A_4 X$

Reverse of forward selection

v) Combination of iii) & iv)

$$\{A_1, A_2, A_3, \dots, A_n\}$$

$A_2$  best  
subset  $X$   $\{A_3\}$

Keep removing irrelevant attributes to another set (at every step)  
Keep adding the best attribute to another set (at every step)

vi) Decision Tree induction.

Wrapper Approach

Age < 32

Tree is built based on the existing data and the new data is run through the tree to predict the probability of a certain event

Filter Approach

Income < 10K

Income > 20K

Yes | No

Ex:  $\{A_1, A_2, A_3, \dots, A_6\}$

Populate a tree

Individuals with  $A_4 = 1$  or  $A_5 = 1$  or  $A_6 = 1$  can be predicted.  $A_1, A_2, A_3, A_4, A_5, A_6$  can be predicted.



15/7/19

## Data Compression :

depends on whether the uncompressible data is lossy or loss-less.

Ex: PCA, Wavelet

Ex: Huffman

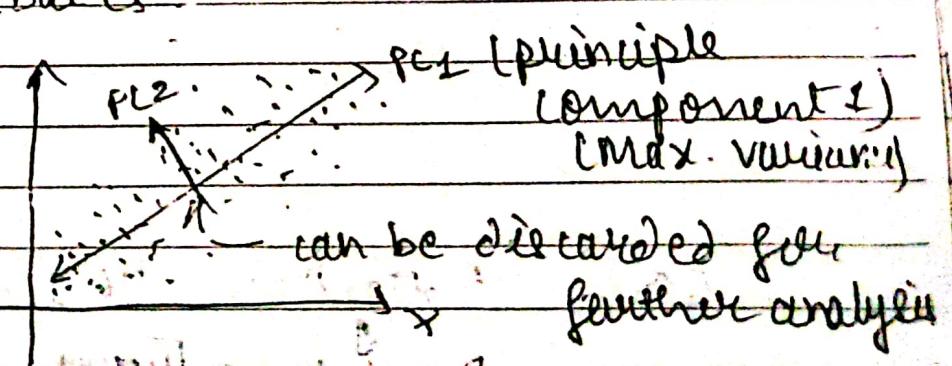
PCA (Principle component analysis)  
To reduce multi-dimensional data into a set with only relevant attributes

can be carried out by considering:

- i) Max. variance
- ii) Least square approximation
- iii) Matrix factorization (spectral decomposition)

• Select principle components.

PCA transforms the entire data set in some other coordinate system.  
2 dimensional rather than taking a subset of the attributes



STEPS:

- i) For PCA, normalizing the attributes plays an important role. Mostly, Z-score normalization is applied.

STEPS FOR SIGNIFICANT

iii) construct co-variance matrix

Ex: For three attributes:

A	$\text{cov}(A,A)$	$\text{cov}(A,B)$	$\text{cov}(A,C)$
B	$\text{cov}(B,A)$	$\text{cov}(B,B)$	$\text{cov}(B,C)$
C	$\text{cov}(C,A)$	$\text{cov}(C,B)$	$\text{cov}(C,C)$

Where,  $\text{cov}(A,A) = \sigma^2(A)$  [variance of A]  
and  $\text{cov}(A,B) = \text{cov}(B,A)$ .

iii) find eigen value & corresponding eigen vector

iv) sort the eigen values in descending order

v) The eigen vector corresponding to the highest eigen value becomes the first principle component ( $PC_1$ ).

$\frac{\lambda_i}{\sum \lambda_j} \times 100$  represents the percentage that  $\lambda_i$  occupies as a part of whole. (importance of  $\lambda_i$ )

Orthogonal: Sum of all vectors should sum up to 1.

Orthonormal:  $A^T$

$$\begin{bmatrix} v_1 & v_2 & v_3 & \dots & v_n \end{bmatrix}$$

When  $A^T A = I$ , then vectors are orthogonal & linearly independent.

PCA is used when the data is sparse  
(better than wavelet in this case)

$$\text{new Data} = [\text{PCA-matrix}]_{3 \times 5} \times [\text{orig-data}]_{5 \times ?}$$

no. of  
Observation

⇒ In this entire process, the data is modified only once [during normalization]

22/1/19

### Numerosity reduction

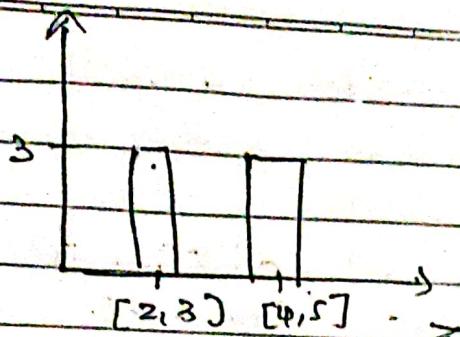
$$y = w^T x + b$$

for any input  $x$ , we need to train  $w$  &  $b$  such that the line represented by it gives minimum outliers.

- Log-linear model
- useful discrete multi-dimensional probability distribution
- Make use of logit function  
 $[ \exp(w^T x + b) ]$
- Sigmoid func. (converts values into probability)  
 $\Delta(x) = \frac{1}{1 + e^{-(w^T x + b)}}$
- useful for sparse data
- works correctly upto 1D dimensions
- Histogram

Values are plotted against their frequency. OR Range of values are plotted against their freq.

{ 2, 2, 3, 5, 2, 4, 5 }



Equi-width histogram:  
as the bucket size is uniform (range).

Equi-depth histogram (uniform freq):  
The range is divided such that each range has the same freq.

- V-optimal histogram (v stands for variable) Divided based on data having least variance.

Max. difference histogram represents data such based on the difference between two adjacent values/range

### clustering

- Quality of cluster is measured in terms of diameter.

Diameter is the max. distance between two entities in a cluster.

- Centroid distance is the avg. distan from each obj. in a cluster to its centroid.

Data can be reduced by discarding values from a cluster. Ex- taking only 10-15 nearest/ farthest points

from the centroid based on the requirements.

B<sup>+</sup> tree provides abstraction from bottom to top. (Complete abstraction to PMSF level). The data values are present in leaf nodes only. Thus, we can reduce the dataset by taking only certain keys based on the level of abstraction key. (As a key leads only to a fixed range of data).

### sampling

- simple Random sample without replacement (SRSWR)

T <sub>1</sub>		T <sub>3</sub>		} cannot be replaced Ex - a tuple sample cannot be considered twice)
T <sub>2</sub>			T <sub>5</sub>	
T <sub>3</sub>			T <sub>4</sub>	
T <sub>4</sub>				
T <sub>5</sub>				
T <sub>6</sub>				

- simple Random sample with replacement (SRWSR)

T <sub>1</sub>		T <sub>3</sub>		} can be replaced
T <sub>2</sub>			T <sub>1</sub>	
T <sub>3</sub>			T <sub>5</sub>	
T <sub>4</sub>			T <sub>1</sub>	
T <sub>5</sub>			T <sub>1</sub>	
T <sub>6</sub>			T <sub>1</sub>	

- Cluster sample

Divide the entire dataset,  $D$  having  $N$  tuples into mutually exclusive equal size clusters.

Depends on requirement  
(not a compulsion)

- Stratified sample

It will divide the dataset into mutually exclusive strata.

Young	Young
M-A	M-A
M-A	M-A
Aged	Aged
Aged	Aged

Based on probability distribution, pick a representative [sub-set of] data from every stratum.

### Discretization & Concept hierarchy:

- Segmentation by natural partition (3-4-5 rule)

- Concept hierarchy generation for categorical data.

→ What value should be converted into a probability distribution graph?

→ When we need to take some decision based on a threshold value, we can write the point by a few decimal

places. But, when the entire set is represented as a probability, we can predict the chance based on the overall set. (Ex- determining whether a cell is contagious or not)

25/7/19

Discritization & concept hierarchy generation by natural partition for numerical data.

- Binning
- Histogram Analysis
- Clustering
- Segmentation by natural Partition (3-4-5 rule)
  - i) If interval covers 3, 6, 7, 9 values at MSB, then divide your range into 3 partitions.
  - ii) If interval covers 2, 4, 8 values at MSB, then divide your range into 4 partitions.
  - iii) If interval covers 1, 5, 10 values at MSB, then divide your range into 5 partitions.

Q 5. ~~With equal width~~ with unequal width

Ex:

Starting with \$ 351,976.00

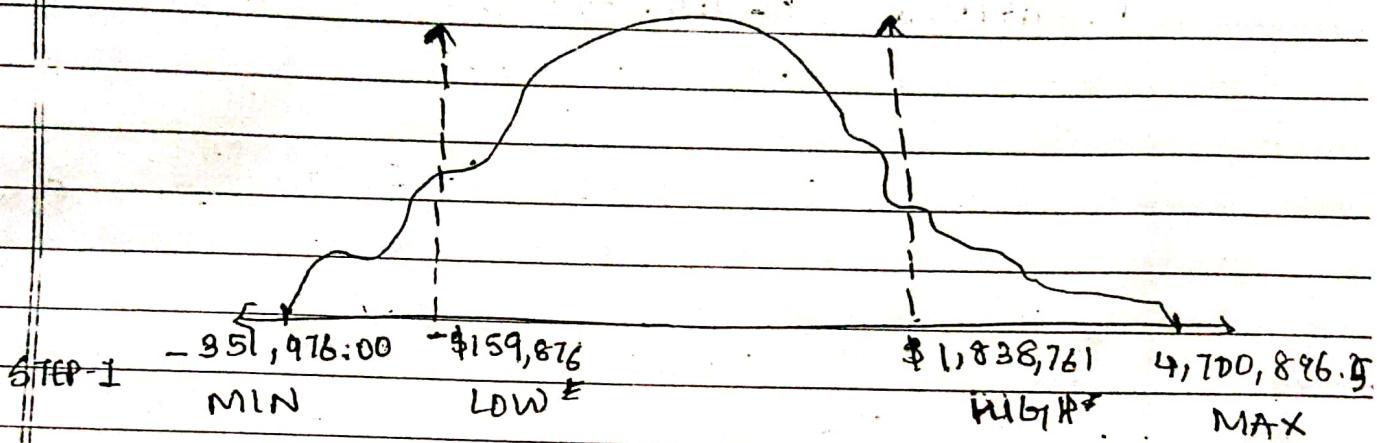
37

to \$ 4,700,896.50

with the given range bins

Partitions are expressed as  $(l, r]$  to avoid overlapping.

Discard 5 percentile and over 95 percentile data (Clipping)  
 Suppose:  $-\$159,876$  to  $+\$1838,761$   
 (low) (high)



STEP-2

$\rightarrow -159,876 \xrightarrow{\text{round down}} \text{nearest 2 digit no.}$

$1,838,761 \xrightarrow{\text{round up}} +2,000,000 \xrightarrow{\text{AIEEE}} (\text{a})$

$$\Rightarrow \underline{2,000,000} - \underline{(-1,000,000)} \\ 1,000,000$$

$\Rightarrow 3$

$$(-1,000,000, +2,000,000)$$

$$(-1,000,000, 0) \quad [0, +1,000,000] \quad (+1,000,000, +2,000,000)$$

If we round up  $-159,876$  to  $+1,000,000$ ,

then spread of the data is 1081.

Since,  $\text{LOW} \leq \text{MIN}$

Partitions:

~~$[-4,000,000, 0]$~~

~~$[0, +1,000,000]$~~

~~$[+1,000,000]$~~

Range:  $(-\$4,000,000,$

$+\$5,000,000)$

~~$1,000,000 + 1,000,000 = 2,000,000$~~

~~$1,000,000 + 1,000,000 = 2,000,000$~~

Min. diversification at lower values.

Partitions:  $[-4,000,000, 0]$

$[0, +1,000,000]$

$[+1,000,000, +2,000,000]$

$[+2,000,000, +5,000,000]$

$\$5,000,000 - \$2,000,000$

$= 3,000,000$

$= 3$

$0 - (-4,000,000)$

$100,000$

$= 4$

$1,000,000 - 0$

$1000000$

$= 10$

$2,000,000 - 1,000,000$

$1000000$

$= 1$

(4 partitions) (5 partitions) (5 partitions)

Hence, we see that there are more partitions in the range where max data resides.

$[-4,000,000, -3,000,000]$

$[-3,000,000, -2,000,000]$

$[-2,000,000, -1,000,000]$

$[-1,000,000, 0]$

$(\$1,000,000, \$1,000,000) \quad [0,000,000, 100,000,000] \quad [\$1,000,000, \$1,000,000]$

Data mining primitive:

1. Task-relevant data

Some of the tuples in the dataset can be discarded based on the task.

and accordingly query is fired. This leads to a correct pattern analysis.

2) The kind of knowledge to be mined  
we need to determine whether we need to classify, cluster the data or find an association rule.

3) Background knowledge  
Required for processing the data and concept hierarchy generation.

4) Interestingness measure  
How unique and informative a mined pattern is.

Terms: support, confidence, simplicity, novelty

Support - (Utility)

$\text{age}(x, "30..39") \wedge$   
 $\text{income}(x, "40K..44K")$  } Association Rule  
 $\Rightarrow \text{buye}(x, "VCR")$   
 $[2.2\%, 60\%]$   
Support      Confidence

and

$\text{occupation}(x, "Student") \wedge$

$\text{age}(x, "20..29")$

$\Rightarrow \text{buye}(x, "computer")$

$[1.4\%, 70\%]$  Lifetime

Support

Example:

Support(A  $\rightarrow$  B) = tuples containing A & B  
total no. of tuples.

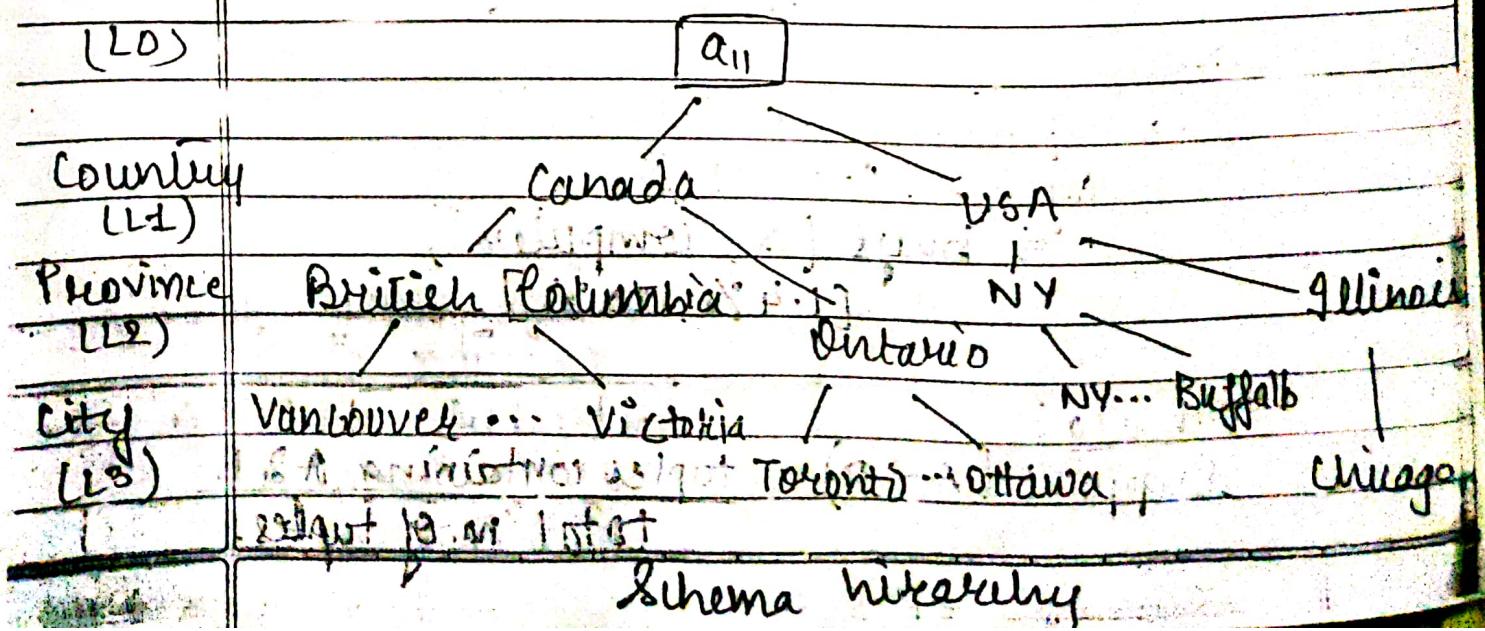
confidence ( $A \rightarrow B$ ) =  $\frac{\text{number of tuples containing } A \text{ & } B}{\text{total } A \text{ tuples}}$   
(certainty)

gives a measure of trustworthiness of a rule.

Simplicity - The rule shouldn't be too complex.

It also defines the support and confidence. A threshold is set and if the support/confidence comes out to be higher than that, the rule is considered to be an important one or the pattern is considered interesting.

Novelty - If a certain rule does not yield any new info. (or yields info. which can be mined by some other, already established rule) then it is not said to be novel.



i) Schema hierarchy is based on total or partial ordering of the dataset.

Higher the level in the schema hierarchy lesser is its frequency (generally).

location (x, "Canada")  $\Rightarrow$  buys(x, "SONY-TV")  $[8\%, 70\%]$

location (x, "Montreal")  $\Rightarrow$  buys(x, "SONY-TV")  $[2\%, 71\%]$

The second rule is redundant (in term of confidence) as Montreal is a part of Canada and becomes a subset of the first rule. (Provides no new info.  $\rightarrow$  not novel).

ii) Set grouping.

{ young, middle-aged, senior } c all(age)

{ 20..39 } c young

{ 40..59 } c middle-aged

{ 60..89 } c senior

iii) Operation derived hierarchy.

book @ cs.stu.ca

login-name < department < university < country

iv) Rule based hierarchy

low-profit margin (x)  $\Leftarrow$  price(x, p<sub>1</sub>)  $\wedge$  cost(x, p<sub>2</sub>)  $\wedge$  ((p<sub>1</sub> - p<sub>2</sub>)  $< \$50$ )

medium-profit margin (x)  $\Leftarrow$  price(x, p<sub>1</sub>)  $\wedge$  net(x, p<sub>2</sub>)  $\wedge$  ((p<sub>1</sub> - p<sub>2</sub>)  $\geq \$50$ )  $\wedge$  ((p<sub>1</sub> - p<sub>2</sub>)  $< \$250$ )

## Representations

- Rules
  - Cross Table
  - Bar Charts
  - Tables
  - Pie Charts
  - Decision tree

[coupling 4.3, 4.4.] (E.1)