

Data warehouse and

Knowledge mining

Assignment - 1

Name: R. R. Dewadasa

USN: EN1420DS0033

SEM: 5

SEC: H

COURSE CODE: 20 DS3501

Module - 1

(i) choose a real time domain and apply all the OLAP operations on the cube.

(ans) Dataset on sales of phone and earphones.

	Unnamed	month	year	location	Product	Value	quarter	country
0	18	1	2017	CA	PS5	420	Q1	USA
1	39	1	2017	NY	PS5	510	Q1	USA
2	77	1	2017	CA	VR	400	Q1	USA
3	80	1	2017	WA	VR	560	Q1	USA
4	114	1	2017	NY	PS5	600	Q1	USA

(i) Dice:

each dimension to a certain range of values, while keeping the numbers of dimensions the same is the resulting cube.

we can focus on sales happening in [Jan/Feb/Mar]

year	Product	month	amount
2017	PS5	1	7100
		2	6890
		3	6622
	VR	1	5600
		2	4960
		3	6920

(ii) Roll up:

applying an aggregation function to collapse a number of dimension, we want to focus in the annual revenue for each product and collapse the location dimension.

year	amount
2018	200801
2019	199867

(iii) slice:

fixing certain dimensions to analyse the remaining dimensions, we can focus in the sales happening in 2019 or can focus on sales happening in 2019

	month	quarter	year	location	Product	amount
2 6 3	2	Q1	2019	CA	VR	2100
2 6 4	2	Q1	2019	QU	PS5	620
2 6 5	2	Q1	2019	WA	VR	2560
2 6 6	2	Q1	2019	NY	PS5	590
2 6 7	2	Q1	2019	CA	PS5	660

(iv) Drilldown:

in the reverse of rollup and applying an aggregation function to a fines level of granularity. we want to focus in the annual and monthly revenue for each product

year	amount
2018	12340
2019	210129

(v) Pivot:

analyzing the combination of a pair of selected dimensions. we want to analyze the revenue by year and month.

Product	CA	NY	QU	WA
VR	1219.1111	2291.3076	1410.000	1436.1020
PS5	629.4769	6612.419	2169.2120	1120.1211

→ applying this on code:

(i) Dice:

CA	395		
NY			
Q1	605		
Q1			

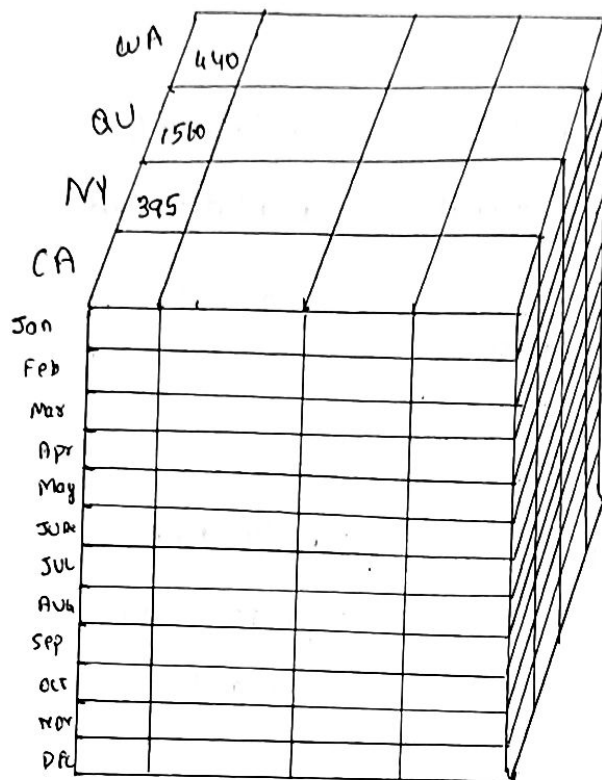
(ii) Rollup:

NY	2000			
CA				
Q1	1000			
Q2				
Q3				
Q4				

(iii) Slice:

WA	440			
QU	1560			
NY	395			
CA				
Q1	605	825	14	400
Q2				
Q3				
Q4				

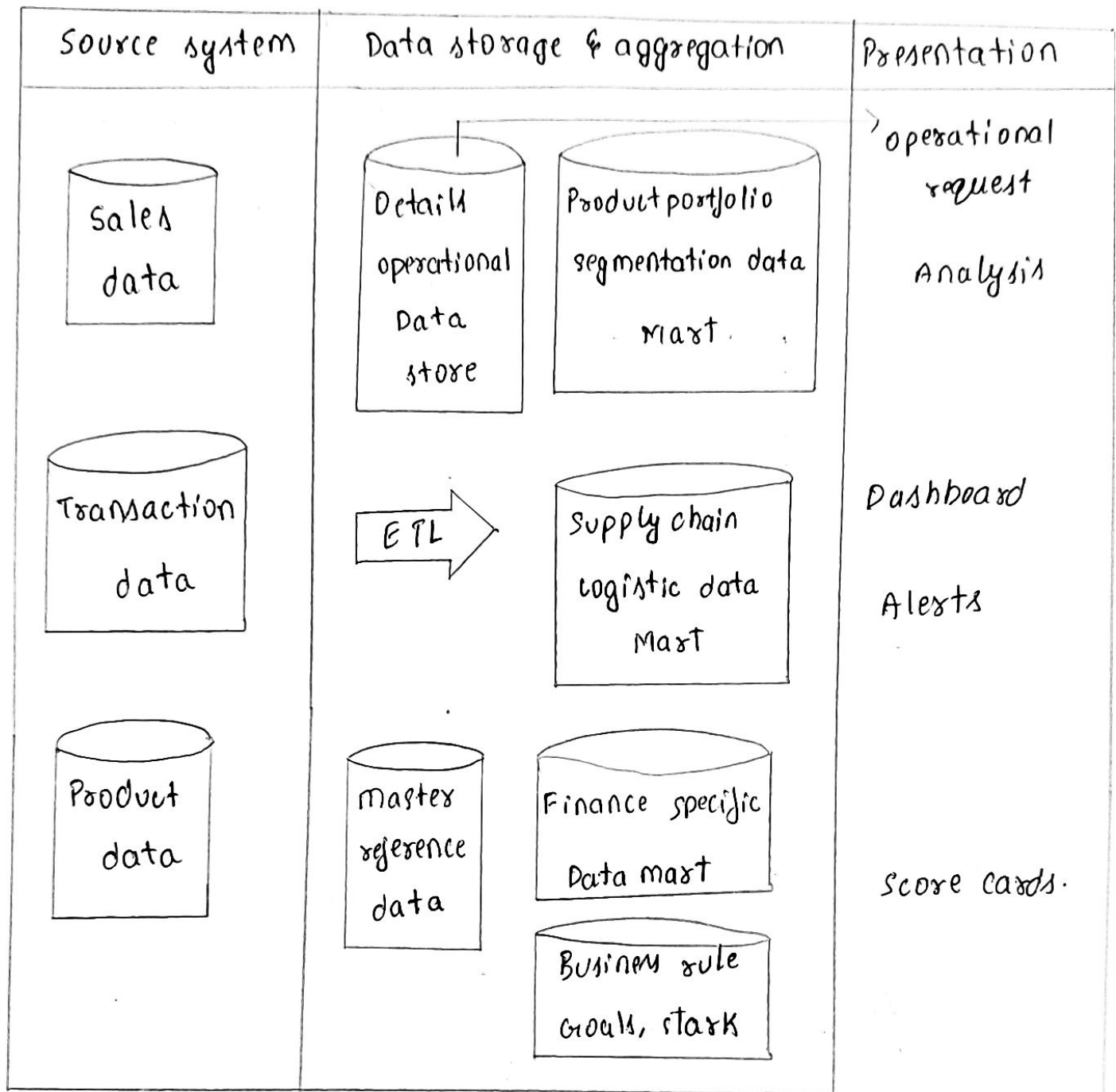
(iv) Drill down:



(v) Pivot:

YR			605
PSS			825

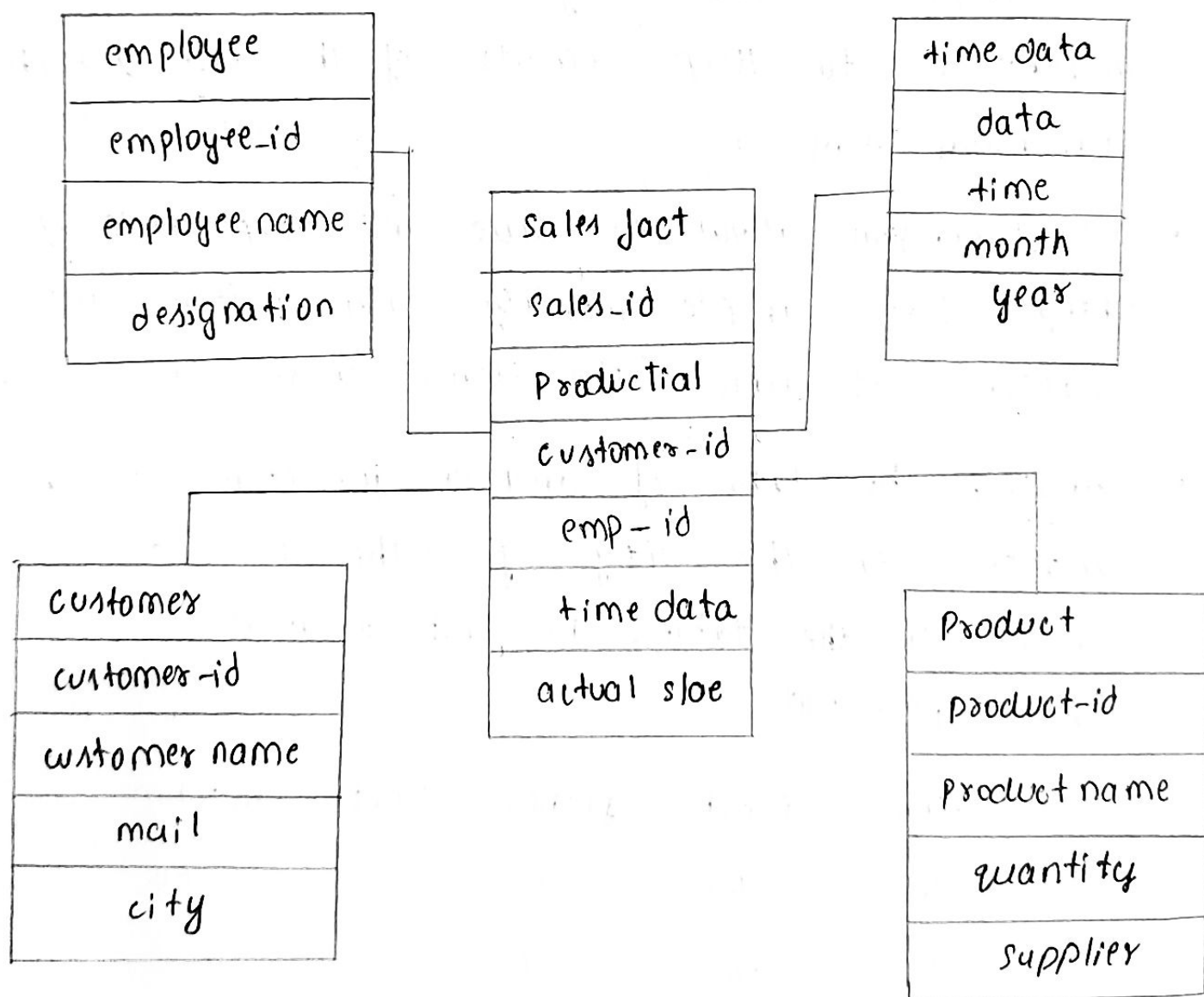
<2> Design a data warehouse for parallel processing.
(Domain- Retail).



- Data warehouse architecture of retail industry helps to understand the basic overview which can be optimised for retail stores
- Data is obtained from multiple sources for example sales data, transaction data.
- The data is extracted from the sources and data cleaning operations are implemented on this data

- This transformed data is fed to data storage and aggregation layer which consists of data warehouse and several data mode
- The data is then presented to the decision making in the form of operational reports, analysis, dashboard.

Schema diagram



<3> Multidimensional data model for amazon

<ans>

- a multi dimensional model views data in the form of data cube. A data cube enables data to be modeled and viewed in multiple dimensions. It is defined by dimensions and fact
- The dimensions are the perspective or entities concerning which an organisation keeps record. For example Amazon may create a sales data warehouse to keep records of the sales for the dimension time.
- The dimension allow to save and keep track of things for example monthly sales of items and location at which the items were sold.
- consider the data of amazon for items sold per quarter in the city of delhi. The data is shown in the table. The fact or measure displayed in rupee-sold.

Time	Fresh	fashion	Electro	Kitchen
Q ₁	260	508	15	60
Q ₂	390	256	20	90
Q ₃	436	396	50	40
Q ₄	528	483	35	50

- If we want to view the sales data with a third dimension. Suppose the data according to time and item, as well as the location is considered for the cities Chennai.

The 3D data are shown below.

Chennai				Kolkata		
Time	Fresh	fashion	electronic	Fresh	fashion	electronic
Q ₁	340	360	20	435	460	20
Q ₂	490	490	16	389	385	45
Q ₃	680	583	46	684	490	39
Q ₄	535	694	39	335	365	83

Time

Chennai	340	360	20	10
Kolkata	435	460	20	15
Mumbai	390	385	20	39
Delhi				
Q ₁	260	508	15	60
Q ₂	390	256	20	90
Q ₃	436	396	50	40
Q ₄	528	483	85	50

item (types)

(4) Design a data cube for market basket analysis.

Ans

- First let understand market basket analysis
- a data mining technique that gives the careful study of purchase done by a customer is a supermarket.
- This concept identifies the pattern of frequent purchase items by customers. This analysis can help to promote deals, offer by the company and data mining technique.
- example:
Data mining concept are in use for sales and marketing to provide better customer service, to improve cross-selling opportunity to increase direct mail response rate
- Customer retention in the form of pattern identification and prediction to likely detection is possible is possible by data mining.
- Risk averted and Fraud area also use the data mining concept for identifying inappropriate or unusual behaviours etc.
- 2 D data:

Bangalore					Pune			
Time	egg	milk	bread	Biscuit	egg	milk	bread	Biscuit
Q ₁	400	360	20	10	500	460	20	15
Q ₂	300	490	16	50	200	385	45	35
Q ₃	200	583	46	43	100	490	39	48
Q ₄	100	694	39	38	600	365	83	35

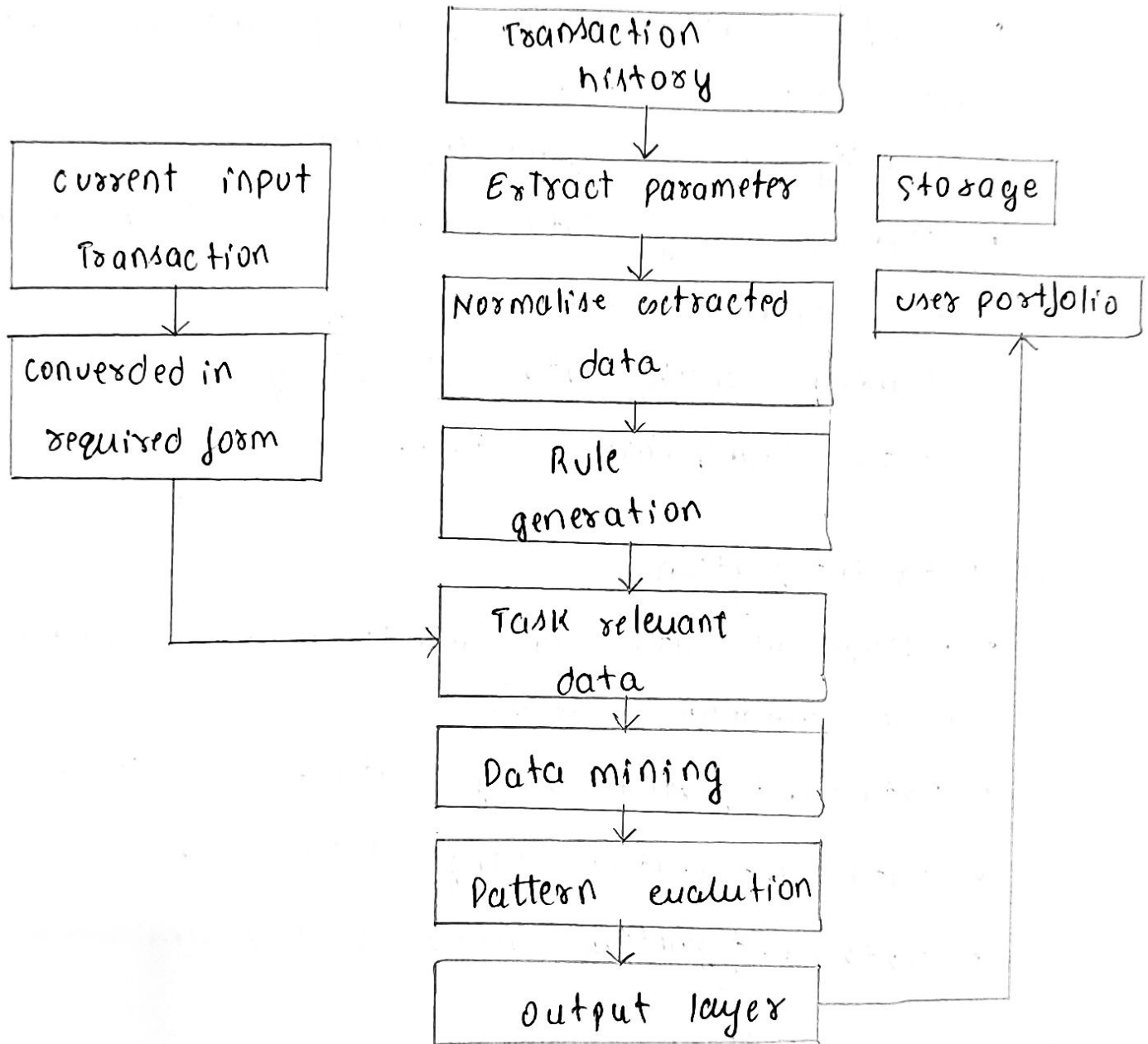
• 3D data:

Location	Bangalore	400	360	20	10
	Pune	500	460	20	15
	Salapur	600	385	20	39
	Vijaypur	700			
Q ₁		260	508	15	60
Q ₂		390	256	20	90
Q ₃		436	396	50	40
Q ₄		528	483	35	56

item.

Module - 2

(1) Design a KDD model for credit card fraud detection



Knowledge

(2) Explain the data preprocessing for information retrieval application.

Ans)

(i) Data cleaning:

- is defined as removal of noisy and irrelevant data from collection
- cleaning in case of missing value
- cleaning noisy data, where noise is a random or variance error.
- cleaning with data discrepancy detection and data transformation tools.

(ii) Data integration:

- heterogeneous data from multiple source combined in a common source.
- Data integration using data mining tools
- Data integration using data synchronization
- Data integration using ETL process

(iii) Data selection:

- process where data relevant to the analysis is decided and retained from the data collection
- Data collection using neural network

- Data selection using decision tree
- Data selection using naive bayes
- Data selection using clustering, regression.

(iv) Data Transformation

- process of transforming data into appropriate required form by miming
- data mapping
- code generation

(v) Data mining:

clever techniques that are applied to extract pattern potentially useful

- Transforms task relevant data into pattern
- Decides purpose of model using classification or characteristics.

(vi) Pattern evaluation:

- identifying strictly increasing pattern representative knowledge based measures
- Find interesting score of each pattern

(vii) Knowledge representation:

- utilizes visualization tools to represent data mining results.
- generate report
- generate table

<3> Evaluate the statistical description for stock market analysis with data visualization.

(ans) GitHub.