

MODULE I

CHAPTER 1

Data Warehouse (DWH) Fundamentals with Introduction to Data Mining

University Prescribed Syllabus w.e.f Academic Year 2021-2022

DWH characteristics, Dimensional modeling : Star, Snowflakes, OLAP operation, OLTP vs OLAP Data Mining as a step in KDD, Kind of patterns to be mined, Technologies used, Data Mining applications.

Self-learning Topics : Data Marts, Major issues in Data Mining.

1.1	Introduction to Data Warehouse	1-3
1.1.1	Features of Data Warehouse	1-3
GQ.	Why is data integration required in a data warehouse more so than in an operational application?	1-3
1.1.2	Need of Data Warehouse.....	1-4
1.1.3	Applications of Data Warehouse.....	1-4
1.1.4	Benefits of Data Warehouse	1-5
1.1.5	Approaches to Build Data Warehouse	1-5
1.1.5(A)	Top-Down Approach	1-5
1.1.5(B)	Bottom-Up Approach.....	1-6
1.2	Data Warehouse Architecture	1-6
GQ.	Consider Metadata as an equivalent of Amazon book store, where each data element is book. What this metadata will contain? Explain.....	1-6
1.2.1	Source Data Component.....	1-7
1.2.2	Data Staging Component.....	1-7
1.2.3	Data Storage Component	1-8
1.2.4	Information Delivery Component	1-8
1.2.5	Metadata Component.....	1-8
1.2.5(A)	Types of Metadata	1-8
1.2.5(B)	Examples of Metadata	1-8
GQ.	What is Metadata? Why do we need metadata when search engines like Google seem so effective?	1-9
1.2.6	Management and Control Component.....	1-9
1.3	Data Warehouse Vs Data Marts	1-10
1.4	E-R Modelling Vs Dimensional Modelling	1-10
GQ.	Why is entity-relationship modeling technique not suitable for the data warehouse? How is dimensional modeling different?.....	1-10
1.4.1	Elements of Dimensional Data Model.....	1-10
1.4.2	Steps of Dimensional Modelling.....	1-11
1.5	Information Package Diagram.....	1-13
1.6	Data Warehouse Schemas	1-13
1.6.1	Star Schema	1-13
1.6.1(A)	Characteristics of Star Schema	1-14
1.6.1(B)	Keys in Star Schema.....	1-14
1.6.1(C)	Advantages of Star Schema	1-14

Data Mining & Business Intelligence (MU-Sem 6-IT) (Data Warehouse Fund)

(Data Warehouse Fund)

Data Mining & Business Intelligence (MU-Sem 6-IT) (Data Warehouse Fund with Intro to Data Mining) Page no. (1-3)

Module

1.6.1(D) Disadvantages of Star Schema	1-15
1.6.1(E) Example	1-15
1.6.2 Snowflake Schema	1-15
1.6.2(A) Characteristics of Snowflake Schema	1-16
1.6.2(B) Advantages of Snowflake Schema	1-16
1.6.2(C) Disadvantages of Snowflake Schema	1-16
1.6.3 Star Schema Vs Snowflake Schema	1-17
1.6.4 Factless Fact Table	1-18
1.6.5 Fact Constellation Schema	1-18
1.6.5(A) Advantages of Fact Constellation Schema	1-18
1.6.5(B) Disadvantages of Fact Constellation Schema	1-19
1.6.6 Schema Definition	1-19
1.7 Update to the Dimension Tables	1-24
1.7.1 Slowly Changing Dimensions	1-25
1.7.2 Rapidly Changing Dimension (RCD)	1-26
1.7.3 Conformed Dimension	1-27
1.7.4 Junk Dimension	1-27
1.7.5 Degenerated Dimension	1-27
1.7.6 Role Playing Dimension	1-27
1.8 Major Steps in ETL Process	1-27
1.9 OLTP Vs OLAP	1-29
1.10 OLAP Operations	1-31
1.11 OLAP Servers	1-37
1.12 Applications of OLAP	1-39
1.13 Hypercube	1-40
1.14 Aggregate Fact Tables	1-40
1.15 Introduction to Data Mining	1-40
1.15.1 Sources of Data that can be Mined	1-41
1.15.2 Kind of Patterns to be mined (Data Mining Techniques)	1-42
1.15.3 Data Mining Technologies	1-43
1.16 Data Mining Task Primitives	1-44
1.17 Data Mining Architecture	1-45
Q. Suppose your task as a software engineer at DBU University is to design a data mining system to examine their university course database, which contains the following information: the name, address, and status (e.g., undergraduate or graduate) of each student, the courses taken, and their cumulative grade point average (GPA). Describe the architecture you would choose. What is the purpose of each component of this architecture?	1-46
1.18 KDD Process	1-47
Q. Explain KDD process with diagram. MU - Dec. 2019	1-47
Q. What is data mining? Explain KDD process with diagram. MU - May 2019	1-47
1.19 Issues in Data Mining	1-48
Q. What are the major issues in data mining? MU - May 2019	1-48
20 Applications of Data Mining	1-49
Chater Ends	1-51

Module	1
1.1 INTRODUCTION TO DATA WAREHOUSE	1. Data Warehouse
	<p>• Data Warehouse is a relational database management system (RDBMS) constructed to meet the requirement of transaction processing systems.</p> <p>(3) It can be loosely described as any centralized data repository which can be queried for business benefits.</p> <p>• It is a database that stores information oriented to satisfy decision-making requests. It is a group of decision support technologies that targets to enabling the knowledge worker (executive, manager, and analyst) to make superior and higher decisions.</p> <p>So, data warehousing support architectures and tool for business executives to systematically organize, understand and use their information to make strategic decisions.</p> <p>It includes historical data derived from transaction data from single and multiple sources.</p> <p>• According to William H. Inmon, "A data warehouse is a subject-oriented, integrated, time-variant, and non-volatile collection of data in support of management's decision making process."</p> <p>(1) The four keywords subject-oriented, integrated, time-variant and non-volatile distinguish data warehouses from other data repository systems, such as relational database systems, transaction processing systems, and file systems.</p> <p>(5) Ultimately a lot of heterogeneous data needs to be integrated and combined. Data historization in the data warehouse requires surrogate keys or artificial hash keys.</p>
	<p>• A data warehouse integrates various heterogeneous data sources like RDBMS, flat files, and online transaction records.</p> <p>(2) It requires performing data cleaning and integration during data warehousing to ensure consistency in naming conventions, encoding structures, attributes measures, etc., among different data sources.</p> <p>• Operational back-ends do have integrations with other systems but generally speaking the integration breadth compared to data warehouse is much smaller because of the limited scope in particular with micro services.</p> <p>• In contrast the scope of a data warehouse is very wide and encompasses practically all important operational systems. Thus data warehouse systems ingest data from practically all important operational systems to power analytics with a broad and complete picture of all enterprise data plus other data sources beyond the scope of the enterprise.</p>
	<p>(3) A data warehouse system ingests data from practically all important operational systems to power analytics with a broad and complete picture of all enterprise data plus other data sources beyond the scope of the enterprise.</p> <p>(4) Ultimately a lot of heterogeneous data needs to be integrated and combined. Data historization in the data warehouse requires surrogate keys or artificial hash keys.</p> <p>• Data across separate sources needs to be aligned and harmonized, and standardized. The need for data cleansing and data quality control is significant. This is why integration work and scope in data warehousing is higher than in operational systems.</p>
	<p>3. Time-variant</p> <p>(1) Historical information is kept in a data warehouse. For example, one can retrieve files from 3 months, 6 months, 12 months, or even previous data from a data warehouse.</p> <p>(2) These varies with transactions systems, where often only the most current file is kept.</p> <p>(3) Every key structure in the data warehouse contains either implicitly or explicitly a time element.</p>
	<p>4. Non-volatile</p> <p>(1) The data warehouse is a physically separate data storage, which is transformed from the source operational RDBMS.</p>

(5N) x 2	Database Systems, Transaction Processing Systems, and File Systems.
1.1.1 Features of Data Warehouse	
GO.	Why is data integration required in a data warehouse more so than in an operational application?
The Key Features of a Data Warehouse are discussed below:	The key features of a data warehouse are discussed below:
1. Subject-Oriented	<p>1. Subject-Oriented</p> <ul style="list-style-type: none"> A data warehouse target on the modeling and analysis of data for decision-makers. Therefore, data warehouses typically provide a concise and straightforward view around a particular subject, such as customer, product, or sales, instead of the global organization's ongoing operations.

GO.	Why is data integration required in a data warehouse more so than in an operational application?
The Key Features of a Data Warehouse are discussed below:	The key features of a data warehouse are discussed below:
1. Subject-Oriented	<p>1. Subject-Oriented</p> <ul style="list-style-type: none"> A data warehouse target on the modeling and analysis of data for decision-makers. Therefore, data warehouses typically provide a concise and straightforward view around a particular subject, such as customer, product, or sales, instead of the global organization's ongoing operations.
2. Integrated	<p>2. Integrated</p> <ul style="list-style-type: none"> A data warehouse integrates various heterogeneous data sources like RDBMS, flat files, and online transaction records.
3. Time-variant	<p>3. Time-variant</p> <ul style="list-style-type: none"> (1) Historical information is kept in a data warehouse. For example, one can retrieve files from 3 months, 6 months, 12 months, or even previous data from a data warehouse. (2) These varies with transactions systems, where often only the most current file is kept. (3) Every key structure in the data warehouse contains either implicitly or explicitly a time element.
4. Non-volatile	<p>4. Non-volatile</p> <ul style="list-style-type: none"> (1) The data warehouse is a physically separate data storage, which is transformed from the source operational RDBMS.

- The operational updates of data do not occur in the data warehouse, i.e., update, insert, and delete operations are not performed.
- It usually requires only two procedures in data accessing: initial loading of data and access to data.
- Therefore, the data warehouse does not require transaction processing, recovery, and concurrency capabilities, which allows for substantial speedup of data retrieval.
- Non-Volatile defines that once entered into the warehouse, and data should not change.

1.1.2 Need of Data Warehouse

Data Warehouse is needed for the following reasons:

- Business User
 - Store historical data
 - Make strategic decisions
 - For data consistency and quality
 - High response time
- **1. Business User :** Business users require a data warehouse to view summarized data from the past. Since these people are non-technical, the data may be presented to them in an elementary form.
- **2. Store historical data :** Data Warehouse is required to store the time variable data from the past. This input is made to be used for various purposes.
- **3. Make strategic decisions :** Some strategies may be depending upon the data in the data warehouse. So, data warehouse contributes to making strategic decisions.
- **4. For data consistency and quality :** Bringing the data from different sources at a commonplace, the user can effectively undertake to bring the uniformity and consistency in data.
- **5. High response time :** Data warehouse has to be ready for somewhat unexpected loads and types of queries, which demands a significant degree of flexibility and quick response time.
- **6. Complete advantage**

is used:
Here, are most common sectors where data warehouse

- Airline
- Banking
- Healthcare
- Public sector
- Investment and Insurance sector
- Retain chain
- Telecommunication
- Hospitality Industry

► **8. Hospitality Industry :** This Industry utilizes warehouse services to design as well as estimate their advertising and promotion campaigns where they want to target clients based on their feedback and travel patterns.

1.1.3 Applications of Data Warehouse

1. Airline

In the Airline system, it is used for operation purpose like crew assignments, analyses of route profitability, frequent flyer program promotions, etc.

► **2. Banking :** It is widely used in the banking sector to manage the resources available on desk effectively. Few banks also use for the market research, performance analysis of the product and operations.

► **3. Healthcare :** Healthcare sector also use data warehouse to strategize and predict outcomes, generate patient's treatment reports, share data with tie-in insurance companies, medical aid services, etc.

► **4. Public sector :** In the public sector, data warehouse is used for intelligence gathering. It helps government agencies to maintain and analyze tax records, health policy records, for every individual.

► **5. Investment and Insurance sector :** In this sector, the warehouses are primarily used to analyze data patterns, customer trends, and to track market movements.

► **6. Retail chain :** In retail chains, data warehouse is widely used for distribution and marketing. It also helps to track items, customer buying pattern, promotions and also used for determining pricing policy.

► **7. Telecommunication :** A data warehouse is used in this sector for product promotions, sales decisions and to make distribution decisions.

► **8. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

1.1.4 Benefits of Data Warehouse

1. Delivers enhanced business intelligence

By having access to information from various sources from a single platform, decision makers will no longer need to rely on limited data or their instinct.

Additionally, data warehouses can effortlessly be applied to a business's processes, for instance, market segmentation, sales, risk, inventory, and financial management.

► **2. Saves times :** A data warehouse standardizes, preserves, and stores data from distinct sources, aiding the consolidation and integration of all the data. Since critical data is available to all users, it allows them to make informed decisions on key aspects. In addition, executives can query the data themselves with little to no IT support, saving more time and money.

► **3. Enhances data quality and consistency :** A data warehouse converts data from multiple sources into a consistent format. Since the data from across the organization is standardized, each department will produce results that are consistent. This will lead to more accurate data, which will become the basis for solid decisions.

► **4. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **5. Improves the decision-making process :** Data warehousing provides better insights to decision makers by maintaining a cohesive database of current and historical data. By transforming data into purposeful information, decision makers can perform more functional, precise, and reliable analysis and create more useful reports with ease.

► **6. Enables organizations to forecast with confidence :** Data professionals can analyze business data to make market forecasts, identify potential KPIs, and gauge predicated results, allowing key personnel to plan accordingly.

► **7. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **8. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **9. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **10. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **11. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **12. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **13. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **14. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **15. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **16. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **17. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **18. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **19. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **20. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **21. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **22. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **23. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **24. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **25. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **26. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **27. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **28. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **29. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **30. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **31. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **32. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **33. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **34. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **35. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **36. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **37. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **38. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **39. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **40. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **41. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **42. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **43. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **44. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **45. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **46. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **47. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **48. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **49. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **50. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **51. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **52. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **53. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **54. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **55. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **56. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **57. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **58. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **59. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **60. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **61. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **62. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **63. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **64. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **65. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **66. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **67. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **68. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **69. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **70. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **71. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **72. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **73. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **74. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **75. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **76. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **77. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **78. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **79. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **80. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **81. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **82. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **83. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **84. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **85. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **86. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **87. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **88. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **89. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **90. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **91. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **92. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **93. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **94. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **95. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **96. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **97. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **98. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **99. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **100. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **101. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **102. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **103. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **104. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **105. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **106. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **107. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **108. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **109. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **110. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **111. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **112. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **113. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **114. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **115. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **116. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **117. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **118. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **119. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **120. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **121. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **122. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **123. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **124. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

► **125. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

► **126. Improves the decision-making process :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

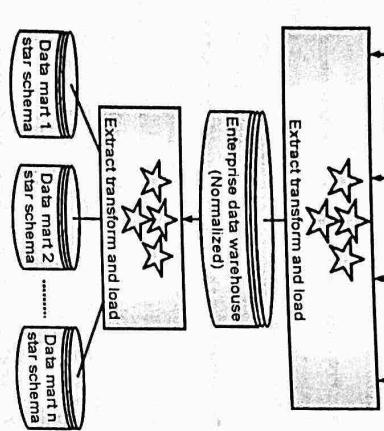
► **127. Provides competitive advantage :** Data warehouses help get a holistic view of their current standing and evaluate opportunities and risks, thus providing companies with a competitive advantage.

► **128. Streamlines the flow of information :** Data warehousing facilitates the flow of information through a network connecting all related or non-related parties.

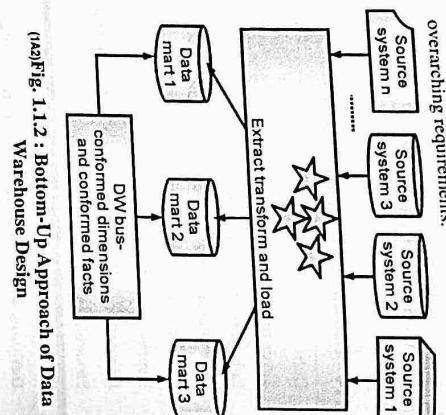
► **129. Generates a high Return on Investment (ROI) :** Companies experience higher revenues and cost savings than those that haven't invested in a data warehouse.

Data Mining & Business Intelligence (MU-Sem 6-IT)

(Data Warehouse Fund with Intro. to Data Mining)...Page no. (1-6)



(1a)Fig. 1.1.1 : Top-Down Approach of Data Warehouse Design



(1a)Fig. 1.1.2 : Bottom-Up Approach of Data Warehouse Design

Advantages

- (1) Represents a data view from the perspective of the enterprise.
- (2) Inherently designed—not a mash-up of disparate data marts.
- (3) Data about the content is stored in a single, central location.
- (4) Centralized control and rules.

Disadvantages

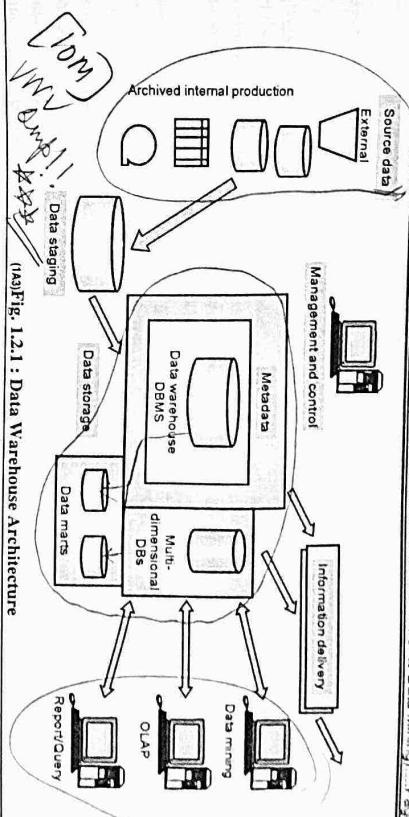
- (1) Even with an iterative strategy, building takes longer.
- (2) High failure risk/exposure
- (3) Requires a high level of cross-functional expertise
- (4) Expenses are high without proof of concept.

1.1.5(B) Bottom-Up Approach

- You create departmental data marts one by one using this bottom-up method.
- To figure out which data marts to build first, you'd create a priority list.
- The most serious disadvantage of this method is data fragmentation.

Data Mining & Business Intelligence (MU-Sem 6-IT)

(Data Warehouse Fund with Intro. to Data Mining)...Page no. (1-7)



(1a)Fig. 1.2.1 : Data Warehouse Architecture

Advantages

- (1) Implementation of small portions is faster and easier.
- (2) Favourable return on investment and proof of concept.
- (3) There is a lower chance of failure.
- (4) Inherently incremental; significant data marts can be scheduled first.
- (5) Allows the project team to grow and develop.

Disadvantages

- (1) Each data mart has its own skewed perspective on information.
- (2) Every data mart is flooded with redundant information.
- (3) Perpetuates data that is inconsistent and irreconcilable.
- (4) Increases the number of unmanageable interfaces.

1.2 DATA WAREHOUSE ARCHITECTURE

1.2.1 Source Data Component

Source data coming into the data warehouse may be grouped into four broad categories, as discussed here.

- (1) Production data
- (2) Internal data
- (3) Archived data
- (4) External data

Q. Consider Metadata as an equivalent of Amazon book store, where each data element is book. What this metadata will contain? Explain.

We put together numerous components to make up an operational system like order entry, claims processing, or a savings account. The GUI (graphical user interface) is used to interact with users for data entry in the front-end component.

- To get the data ready, three major functions must be completed. You must first extract the data, then transform it and at last load it into the data warehouse storage.
- A staging area is where the three major functions of extraction, transformation, and loading take place.
- Data staging is a place with a set of functions that clean, change, combine, convert, deduplicate, and prepare source data for storage and use in the data warehouse.

1.2.3 Data Storage Component

- The data warehouse's data storage is kept in a separate repository. Large volumes of historical data must be kept in a data warehouse's data repository for analysis. Furthermore, the data in the data warehouse must be kept in structures suitable for analysis rather than for quick retrieval of individual pieces of information.
- As a result, the data warehouse's data storage is kept separate from the data storage for operational systems.

1.2.4 Information Delivery Component

- The information delivery component includes various methods of information delivery in order to provide information to the large community of data warehouse users.
- Information can be delivered in the form of ad hoc reports, complex queries, multidimensional (MD) analysis, statistical analysis, enterprise information system (EIS) feeds, etc.
- This information can be delivered through online, intranet, internet or email mode.

1.2.5 Metadata Component

- Metadata is "data that describes other data". Metadata in a data warehouse is similar to the data dictionary or the data catalog in a database management system.
- In the data dictionary, you keep the information about the logical data structures, the information about the files and addresses, the information about the indexes, and so on. The data dictionary contains data about the data in the database.
- Similarly, the metadata component is the data about the data in the data warehouse.

1.2.5(A) Types of Metadata

Metadata in a data warehouse fall into three major parts:

(a) Operational Metadata

- As we know, data for the data warehouse comes from various operational systems of the enterprise. These source systems include different data structures. The data elements selected for the data warehouse have various fields lengths and data types.
- In selecting information from the source systems for the data warehouses, we divide records, combine factors of documents from different source files, and deal with multiple coding schemes and field lengths. When we deliver information to the end-users, we must be able to tie that back to the source data sets. Operational metadata contains all of this information about the operational data sources.

(b) Extraction and Transformation Metadata

- Extraction and transformation metadata include data about the removal of data from the source systems, namely, the extraction frequencies, extraction methods, and business rules for the data extraction.
- Also, this category of metadata contains information about all the data transformation that takes place in the data staging area.

(c) End-User Metadata

- The end-user metadata is the navigational map of the data warehouses. It enables the end-users to find data from the data warehouses.
- The end-user metadata allows the end-users to use their business terminology and look for the information in those ways in which they usually think of the business.

1.2.5(B) Examples of Metadata

- A library catalog may be considered metadata. The directory metadata consists of several predefined components representing specific attributes of a resource, and each item can have one or more values. These components could be the name of the author, the name of the document, the publisher's name, the publication date, and the methods to which it belongs.

1.2.6 Management and Control Component

- This component of the data warehouse architecture sits on top of all the other components. The management and control component coordinates the services and activities within the data warehouse.
- This component controls the data transformation and the data transfer into the data warehouse storage.
- On the other hand, it moderates the information delivery to the users. It works with the database management systems and enables data to be properly stored in the repositories.
- It monitors the movement of data into the staging area and from there into the data warehouse storage itself.
- The management and control component interacts with the metadata component to perform the management and control functions.

- As the metadata component contains information about the data warehouse itself, the metadata is the source of information for the management module.

1.3 DATA WAREHOUSE VS DATA MARTS

- A data mart is a small, single-subject data warehouse subset that provides decision support to a small group of people.
- Data Mart can serve as a test vehicle for companies exploring the potential benefits of Data Warehouses.
- Data Mart address local or departmental problems, while a Data Warehouse involves a company-wide effort to support decision making at all levels in the organization.

Table 1.3.1 : Data Warehouse Vs Data Mart

Sr. No.	Data Warehouse	Data Mart
1.	Data warehouse is a centralized system.	Data mart is a decentralized system.
2.	In data warehouse, In Data mart, highly denormalization takes place.	In Data mart, highly denormalization takes place.
3.	Data warehouse is top-down model.	Data mart is a bottom-up model.
4.	To build a warehouse is difficult.	To build a mart is easy.
5.	In data warehouse, Fact constellation schema is used.	In data mart, Star schema and snowflake schema are used.
6.	Data warehouse is flexible.	Data mart is not flexible.
7.	Data warehouse is data-oriented in nature.	Data mart is the project-oriented in nature.
8.	Data warehouse has long life.	Data-mart has short life than warehouse.
9.	In data warehouse, data are contained in detail form.	In data mart, data are summarized form.
10.	Data Warehouse is vast in size.	Data mart is smaller than warehouse.

- Q. Why is entity-relationship modeling technique not suitable for the data warehouse? How is dimensional modeling different?

1.4 E-R MODELLING VS DIMENSIONAL MODELLING

- Dimensional Modeling (DM) is a data structure technique optimized for data storage in a Data warehouse.
- The purpose of dimensional modeling is to optimize the database for faster retrieval of data. The concept of Dimensional Modelling was developed by Ralph Kimball and consists of "fact" and "dimension" tables. Dimensional table records information on each dimension, and fact table records all the "fact", or measures.
- A dimensional model in data warehouse is designed to read, summarize, analyze numeric information like values, balances, counts, weights, etc. in a data warehouse.
- In contrast, relation models are optimized for addition, updating and deletion of data in a real-time Online Transaction System.
- These dimensional and relational models have their unique way of data storage that has specific advantages.
- For instance, in the relational model, normalization and ER models reduce redundancy in data. On the contrary, dimensional model in data warehouse arranges data in such a way that it is easier to retrieve information and generate reports.
- Hence, Dimensional models are used in data warehouse systems and not a good fit for relational systems.

1.4.1 Elements of Dimensional Data Model

- Q. What are the measurement/metrics or facts from your business process. For a Sales business process, a measurement would be quarterly sales number.
- Q. Who - Customer Names, Product Names, Who - Location
- Fact : Facts are the measurements/metrics or facts from your business process. For a Sales business process, a measurement would be quarterly sales number.
 - Dimensions : Dimensions provide the context surrounding a business process event. In simple terms, they give who, what, where of a fact. In the Sales business process, for the fact quarterly sales number, dimensions would be : ③ **Product Name**, **Customer Name**, **Who - Customer Names**, **Product Name**, **Who - Location**.
 - Where - Location

- What – Product Name
- In other words, a dimension is a window to view information in the facts.

1.4.2 Steps of Dimensional Modeling

- The accuracy in creating your dimensional modeling determines the success of your data warehouse implementation. The model should describe the Why, How much, When/Where/Who and What of your business process. Here are the steps to create dimension model.
- Q. Who - Customer Names, Product Names, Who - Location
1. Identify Business Process
- Dimensions are nouns like date, store, inventory, etc. These dimensions are where all the data should be stored. For example, the date dimension may contain data like a year, month and weekday.
 - Example of Dimensions: The CEO at an MNC wants to find the sales for specific products in different locations on a daily basis. So, the grain is "product sale information by location by the day".
2. Identify Dimensions and Attributes
- Dimensions are nouns like date, store, inventory, etc. These dimensions are where all the data should be stored. For example, the date dimension may contain data like a year, month and weekday.
 - Example of Dimensions: The CEO at an MNC wants to find the sales for specific products in different locations on a daily basis.
 - (i) Dimensions : Product, Location and Time
 - (ii) Attributes : For Product: Product key (Foreign Key), Name, Type, Specifications
 - (iii) Hierarchies : For Location: Country, State, City, Street Address, Name

- This step is co-associated with the business users of the system because this is where they get access to data stored in the data warehouse. Most of the fact table rows are numerical values like price or cost per unit, etc.

- Example of Facts : The CEO at an MNC wants to find the sales for specific products in different locations on a daily basis. The fact here is Sum of Sales by product by location by time.

- In this step, you implement the Dimension Model. A schema is nothing but the database structure (arrangement of tables). There are two popular schemas

- Star Schema
- Snowflake Schema

Table 1.4.1 : E-R Modeling Vs Dimensional Modeling

Sr. No.	E-R Modeling	Dimensional Modeling
1.	E-R modeling have logical and physical model.	Dimensional modeling have physical model.
2.	E-R modeling is used for normalizing the OLTP design.	Dimensional modeling is used for de-normalizing the OLAP design.
3.	E-R modeling revolves around the entities and their relationships to capture the overall process of the system.	Dimensional modeling revolves around dimensions (point of analysis) for decision making and dimensions by time, product, location, and customer demographics.
4.	In E-R modeling the data is in normalized form, so more number of joins are required, which may adversely affect the system performance.	In dimensional modeling the data is de-normalized, so less number of joins are required, by which system performance will improve.
5.	In E-R modeling, a view of data is from data processing.	In dimensional modeling, a view of data is from business processing.
6.	It is not mapped for creating schemas.	It is mapped for creating schemas.
7.	It uses the current data.	It uses the historical data.
8.	Size of data varies from MBs to GBs.	Size of data varies from GBs to TBs.
9.	Data storage is volatile.	Data storage is non-volatile.
10.	High Create/Read/Update/Delete activity.	High Select activity only insert & access.
11.	Advantages : Removes data redundancy. Ensures data consistency. Expresses the relationship between the entities.	Advantages : Captures critical measures. Views along dimensions. Useful to business users.

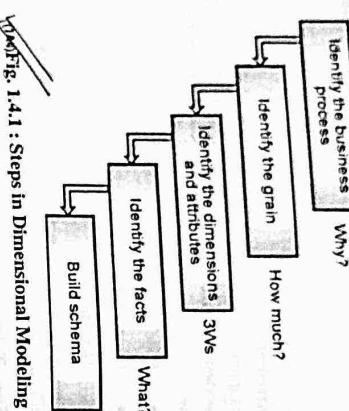


Fig. 1.4.1 : Steps in Dimensional Modeling

- The first and most generalized level of an information model is its information package diagram. This model focuses on the data gathering activities for the user's information packaging requirements.

- An information package diagram defines the relationships between subject matter and key performance measures.

- The information package diagram has a highly targeted purpose, providing a focused scope for user requirements.

- Because information package diagrams target what the users want, they are effective in facilitating communication between the technical staff and the users, indicating any inconsistencies between the requirements and what the data warehouse will deliver.

Example : Information package for analyzing sales for a certain business. It allows users to evaluate sales metrics by time, product, location, and customer demographics.

The subject here is sales. The measured facts or the measurements that are of interest for analysis are shown in the bottom section of the package diagram. In this case, the measurements are actual sales, forecast sales, and budget sales.

1.5 INFORMATION PACKAGE DIAGRAM

Fig.

1.5

Information Package Diagram

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

1

Fact Tables

- It is a table in a star schema which contains facts and connected to dimensions.

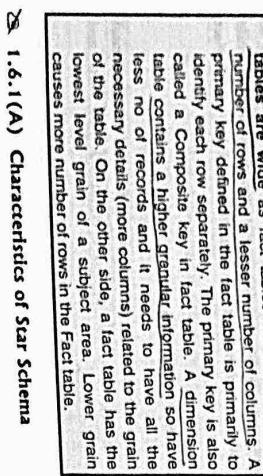
A fact table has two types of columns: those that include fact and those that are foreign keys to the dimension tables.

The primary key of the fact tables is generally a composite key that is made up of all of its foreign keys.

Dimension Tables

- A dimension is an architecture usually composed of one or more hierarchies that categorize data.
- If a dimension has not got hierarchies and levels, it is called a flat dimension or list.
- The primary keys of each of the dimension table are part of the composite primary keys of the fact table.
- Dimensional attributes help to define the dimensional value. They are generally descriptive, textual values.
- Dimensional tables are usually small in size than fact table.

Note : Fact tables are deep whereas dimension tables are wide as fact tables will have a higher number of rows and a lesser number of columns. A primary key defined in the fact table is primarily to identify each row separately. The primary key is also called a Composite key in fact table. A dimension table contains a higher granular information so have less no of records and it needs to have all the necessary details (more columns) related to the grain of the table. On the other side, a fact table has the lowest level grain of a subject area. Lower grain causes more number of rows in the Fact table.



(FIG. 1.6.1 : Star Schema for Digital Electronics Sale

1.6.1(A) Characteristics of Star Schema

- It creates a de-normalized database that can quickly provide query responses.
- It provides a flexible design that can be changed easily or added to throughout the development cycle, and as the database grows.
- It provides a parallel in design to how end-users typically think of and use the data.
- It reduces the complexity of metadata for both developers and end-users.

1.6.1(B) Keys In Star Schema

- | | |
|-----------------|-------------------|
| 1. Primary Keys | 2. Surrogate Keys |
|-----------------|-------------------|

- Primary Keys :** The primary key of the dimension table identifies each row in a dimension table. Example: In a student dimension table, student_id is the primary key which identifies each student uniquely.
- Surrogate Keys :** System generated sequence numbers are called surrogate keys. They do not have any built-in meanings.

- Foreign Keys :** Every dimension table has one-to-one relationship with the fact table. The primary key in the dimension table acts as a foreign key in the fact table.

- Easily Understood**: A star schema is simple to understand and navigate, with dimensions joined only through the fact table.
- These joins are more significant to the end-user because they represent the fundamental relationship between parts of the underlying business. Customer can be retrieved.

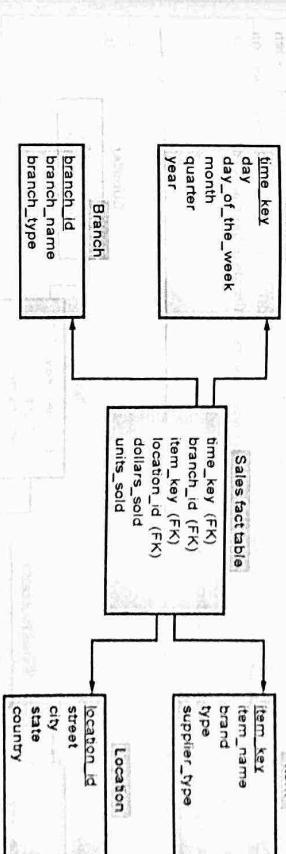
- 1.6.1(E) Example**
- A star schema for Digital Electronics sales is shown. Sales are considered along four dimensions: time, item, branch, and location. The schema contains a central fact table for sales that contains keys to each of the four dimensions, along with two measures: dollars sold and units sold.

1.6.1(D) Disadvantages of Star Schema

- Dimension table can be populated once and occasionally refreshed. We can add new facts regularly and selectively by appending records to a fact table.

- Built-in referential integrity**

- Date integrity is not enforced at all since in a highly normalized schema state.
- Not flexible in terms of analytical needs as a normalized data model.
- Star schemas don't reinforce many-to-many relationships within business entities, at least not frequently.



(FIG. 1.6.2 : Snowflake Schema

- The snowflake schema is a variant of the star schema. Here, the centralized fact table is connected to multiple dimensions. In the snowflake schema, dimensions are present in a normalized form in multiple related tables.
- The snowflake structure is materialized when the dimensions of a star schema are detailed and highly structured, having several levels of relationship, and the child tables have multiple parent table.
- The snowflake schema affects only the dimension tables and does not affect the fact tables.
- In other words, a dimension table is said to be snowflaked if the low-cardinality attribute of the dimensions has been divided into separate normalized tables. These tables are then joined to the original dimension table with referential constraints (foreign key constraint).

1.6.2(A) Characteristics of Snowflake Schema

- (1) The snowflake schema uses small disk space.
- (2) It is easy to implement dimension that is added to schema.
- (3) There are multiple tables, so performance is reduced.
- (4) The dimension table consist of two or more sets of attributes which define information at different grains.
- (5) The sets of attributes of the same dimension table are being populated by different source systems.

1.6.2(B) Advantages of Snowflake Schema

- (1) It provides structured data which reduces the problem of data integrity.
- (2) It uses small disk space because data are highly structured.

1.6.2(C) Disadvantages of Snowflake Schema

- (1) Snowflaking reduces space consumed by dimension tables, but compared with the entire data warehouse the saving is usually insignificant.
- (2) Avoid snowflaking or normalization of a dimension table unless required and appropriate.

1.6.2(D) Example

- A snowflake schema for Digi1 Electronics sales is given. Here, the sales fact table is identical to that of the star schema in Fig. 1.6.1. The main difference between the two schemas is in the definition of dimension tables. The single dimension table for item in the star schema is normalized in the snowflake schema, resulting in new item and supplier tables.
- For example, the item dimension table now contains the attributes item_key, item_name, brand, type, and supplier_key, where supplier_key is linked to the supplier dimension table, containing supplier_key and supplier_type information.
- Similarly, the single dimension table for location in the star schema can be normalized into two new tables: location and city. The city_key in the new location table links to the city dimension. Notice that, when desirable, further normalization can be performed on state and country in the snowflake schema.

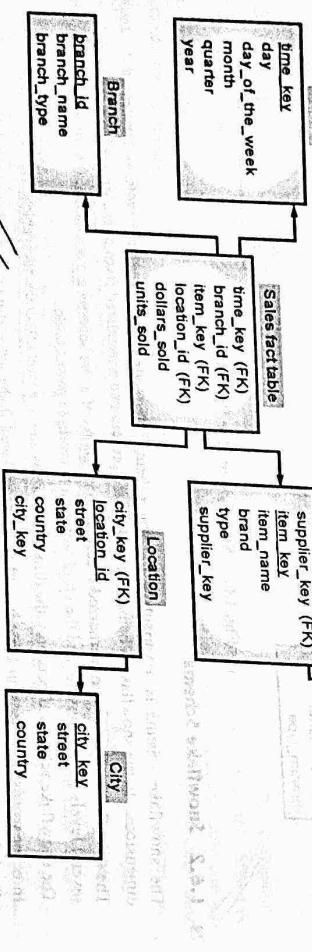


Fig. 1.6.2 : Snowflake Schema for Digi1 Electronics Sale

1.6.3 Star Schema Vs Snowflake Schema

Table 1.6.1 : Star Schema Vs Snowflake Schema

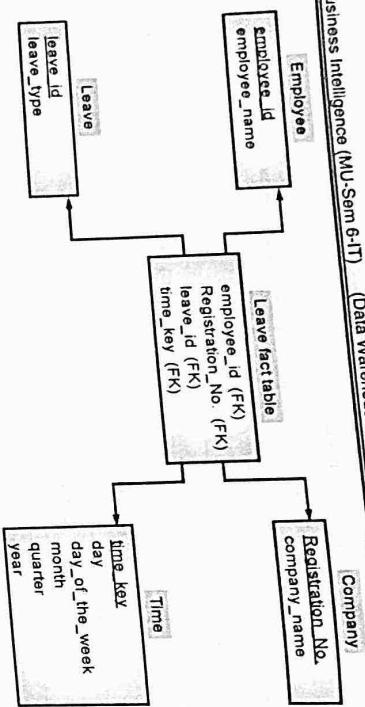
Sr. No.	Basis	Star Schema	Snowflake Schema
1.	Ease of Maintenance/change	It has <u>redundant data</u> and hence less easy to maintain/change.	No redundancy and therefore more easy to maintain and change.
2.	Ease of Use	<u>Less complex queries</u> and simple to understand.	More complex queries and therefore less easy to understand.
3.	Parent table	In a star schema, a dimension table will not have any parent table.	In a snowflake schema, a dimension table will have one or more parent tables.
4.	Query Performance	Less number of <u>foreign keys</u> and hence <u>lesser query execution time</u> .	More foreign keys and thus <u>more query execution time</u> .
5.	Normalization	It has <u>De-normalized tables</u> .	It has <u>normalized tables</u> .
6.	Type of Data	Good for data marts with simple relationships (one to one or one to many).	Good to use for data warehouse core to simplify complex relationships (many to many).
7.	Joins	Execution	Execution
8.	Dimension Table	It contains only a single dimension table for each dimension.	Hierarchies for the dimension are stored in the dimensional table itself in a star schema.
9.	Hierarchies	Dimensional table	Hierarchies are broken into <u>separate tables</u> in a snowflake schema. These hierarchies help to drill down the information from topmost hierarchies to the lowermost hierarchies.
10.	When to use	When the dimensional table contains less number of rows, we can go for Star schema.	When dimensional table store a huge number of rows with redundancy information and space is such an issue, we can choose snowflake schema to save space.
11.	Data Warehouse	Work best in any data warehouse/ data mart.	Better for <u>small data warehouse</u> mart.

1.6.4 Factless Fact Table

- A data warehouse factless fact table is a fact that does not have any measures stored in it. This table will only contain keys from different dimension tables. The fact-less fact is often used to resolve a many-to-many cardinality issue.

- There are two types of factless fact tables:
1. Event capturing factless fact
 2. Coverage Table – Describing condition

- | <u>Data Mining & Business Intelligence (Mu-Sem 6-T)</u> | (Data Warehouse Fund. with Intro. to Data Mining)... Page no. [1-19] |
|--|---|
| <p>1.6.5(C) Example</p> <ul style="list-style-type: none"> A fact constellation schema is shown in Fig. 1.6.4 below. This schema specifies two fact tables, sales and shipping. The sales table definition is identical to that of the star schema above. The shipping table has four dimensions, | <ul style="list-style-type: none"> or keys_item_key, time_key, shipper_key, location_id and two measures dollars_cost and units_shipped. A fact constellation schema allows dimension tables to be shared between fact tables. For example, the dimension tables for time, item, and location are shared between the sales and shipping fact tables. |



(a) Fig. 1.6.3 : Factless Fact Table for Leave

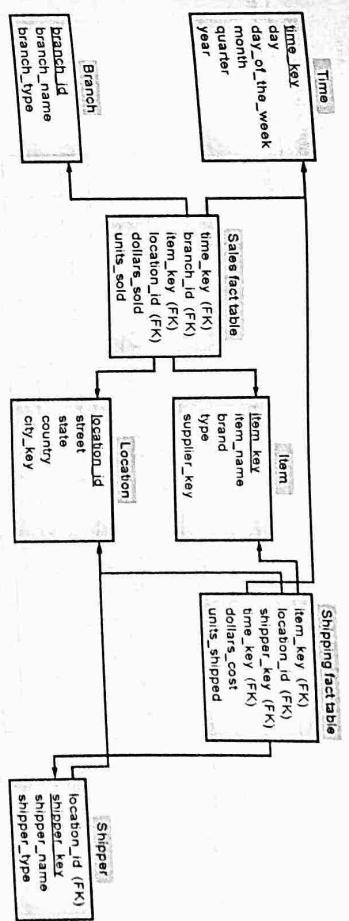
1. Event Capturing

- This type of fact table captures dimension members from various dimension tables without any measured value.

For example, Student attendance (student-teacher relation table) capturing table is the factless fact. Table can have entries into it whenever student attend class.

will have thus ...

- attendance table .
 1. Which student is taught by the maximum number of teachers?
 2. Which class has maximum number of attendance?
 3. Which teacher teaches maximum number of students?
 - All the above queries are based on the COUNT(), MAX() with GROUP BY.
 - **2. Coverage Table-Describing Condition**
 - This is another kind of factless fact. A factless fact table can only answer ‘optimistic’ queries (positive query) but cannot answer a negative query.
 - Coverage fact is used to support negative analysis reports. For example, an electronic store did not sell any product for given period of time.
 - If you consider the student-teacher relation table, the event capturing fact table cannot answer ‘which teacher did not teach any student?’
 - Coverage fact attempts to answer this question by



(1a8)Fig. 1.6.4 : Fact Constellation Schema for Digi1 Electronics Sale

1.6.6 Schema Definition

Multidimensional schema is defined using Data mining [\(see Fig. 1.5\)](#), and dimension definition, can be used for defining the data warehouses and data marts.

Syntax for Cube Definition

Syntax for Dimension Definition

```
define cube < cube_name > [ < dimension_name > ] [ < dimension_name > ] ...
```

Star Schema Definition

The star schema that we have discussed can be defined in the Mining Query Language (DMQL) as follows:

using Data Mining to define cube sales star [time,

`dollars_sold = sum(sales in dollars), units_sold = sum(units sold),
define dimension time as (time_key, day, day_of_week,`

```
define dimension  
supplier_type  
define dimension branch as (branch_id, branch_name,
```

define dimension location
country)

Fact Constellation Schema Definition

Fact constellation schema can be defined using DML as follows:

```
define cube sales [time, item, branch, location] dollars_sold =  
sum(sales in dollars), units_sold = count(*)  
define dimension time as (time_key, day, day_of_week,  
month, quarter, year)  
define dimension item as (item_key, item_name, brand, type,  
supplier_type)  
define dimension branch as (branch_id, branch_name,  
branch_type)
```

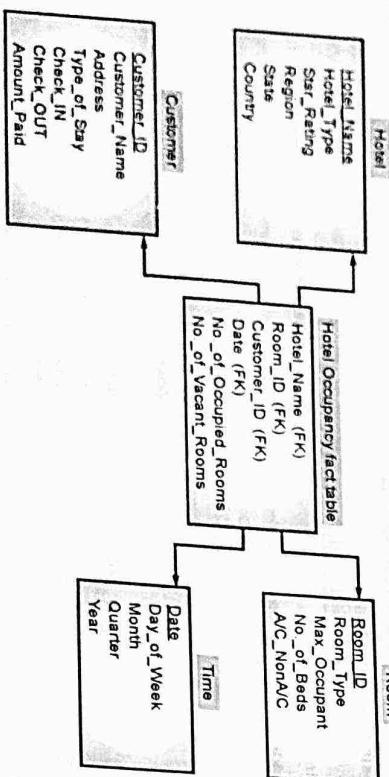
- Ex 1.6.1 :** Consider a data warehouse for hotel occupancy, where there are four dimensions namely (a) Hotel (b) Room (c) Time (d) Customer and two measures (i) Occupied rooms (ii) Vacant rooms.
- Draw information package diagram, star schema and snowflake schema.

Soln : (a) Information Package Diagram

Dimensions					
Hotel	Room	Time	Customer
Hotel Name	Room ID	Date	Customer ID		
Hotel Type	Room Type	Day of Week	Customer Name		
Star Rating	No. of Beds	Month	Address		
Region	A/C Non A/C	Quarter	Type of Stay		
State	Year	Check IN			
Country		Check OUT			
		Amount Paid			

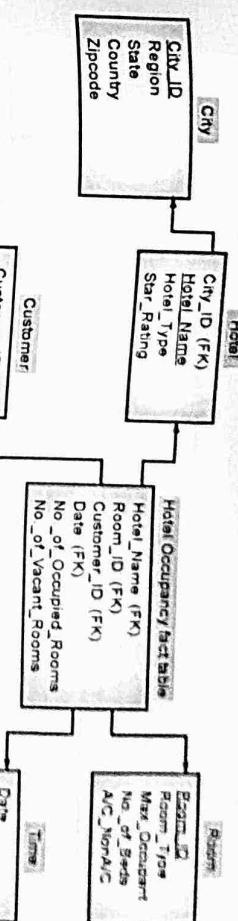
Facts : Occupied Rooms, Vacant Rooms

(b) Star Schema



(i)(a)Fig. P. 1.6.1(a) : Star Schema for Hotel Occupancy

(c) Snowflake Schema



(i)(b)Fig. P. 1.6.1(b) : Snowflake Schema for Hotel Occupancy

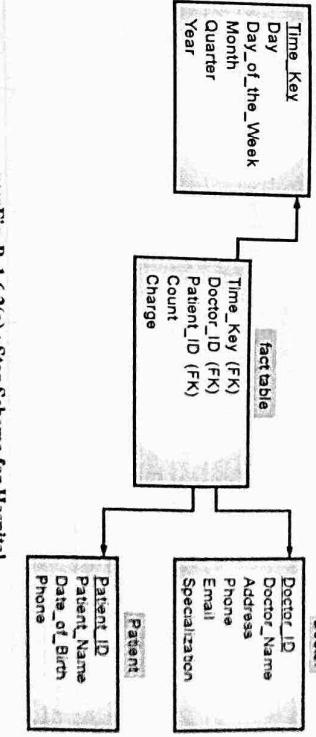
Ex 1.6.2 : Consider a data warehouse for a hospital where there are three dimensions namely (a) Doctor (b) Patient (c) Time and two measures (i) count (ii) charge where charge is the fee that the doctor charges a patient for a visit.

- (i) Draw star and snowflake schema.
(ii) Starting with the base cuboid (day, doctor, patient), what specific OLAP operations should be performed in order to list the total fee collected by each doctor in 2010?
(iii) To obtain the same list, write an SQL query assuming the data are stored in a relational database with the schema free (day, month, year, doctor, hospital, patient, count, charge).

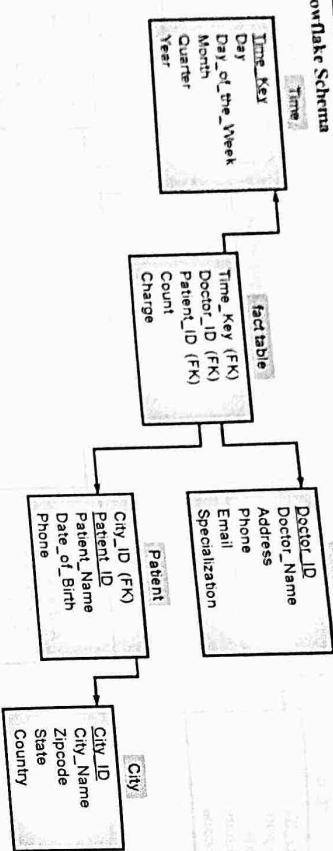
Soln. :

(i) star Schema and Snowflake Schema

(a) Star Schema



(i)(a)Fig. P. 1.6.2(a) : Star Schema for Hospital

(b) Star Schema

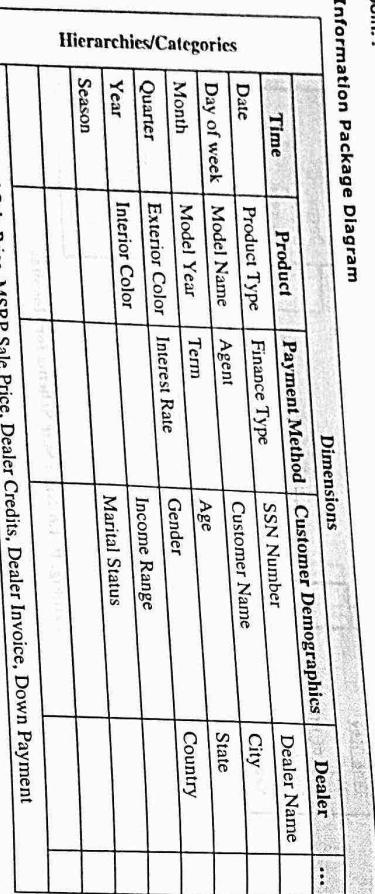
(1a12)Fig. P. 1.6.2(b) : Snowflake Schema for Hospital

(ii) First, we should use roll-up operation to get the year 2010 (rolling-up from day then month to year). After getting that, we need to use slice operation to select (2010). Finally, we get list the total fee collected by each doctor in 2010. So,

1. roll up from day to month to year
2. slice for year = "2010"
3. roll up on patient from individual patient to all
4. slice for patient = "all"
5. get the list of total fee collected by each doctor in 2010.

(iii) SELECT doctor, sum(charge) FROM tee WHERE year = 2004 GROUP BY doctor;

(c) Payment Method (d) Customer Demographics (e) Dealer and facts/measures like Actual Sale Price, MSRP Sale Price, Dealer Credits, Dealer Invoice, Down Payment. Draw the information package diagram and star schema.

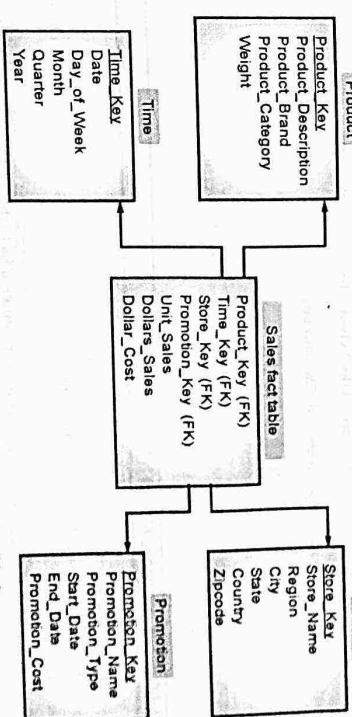
Soln. :**(a) Information Package Diagram**

(1a13)Fig. P. 1.6.3 : Star Schema for Automaker Sales

Ex. 1.6.4 : For a supermarket chain, consider the following dimensions, namely product, store, time and promotion. The schema contains a central fact table for sales with three measures unit_sales, dollars_sales and dollar_cost.

- i. Draw star schema.
- ii. Calculate the maximum number of base fact table records for warehouse with the following values given below:

 - Time period-5 years
 - Store - 300 stores reporting daily sales
 - Product - 40,000 products in each store (about 4000 sell in each store daily).
 - Promotion: a sold item may be in only one promotion in a store on a given day.

Soln. :**i. Star Schema**

(1a14)Fig. P. 1.6.4 : Star Schema for Supermarket Chain

Facts : Actual Sale Price, MSRP Sale Price, Dealer Credits, Dealer Invoice, Down Payment

11. Fact table records

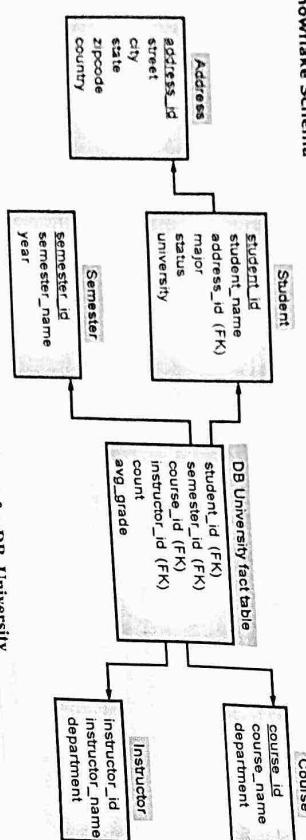
- Time period = 5 years = 5×365 days = 1825 days
 - Number of Stores = 300
 - Product sell in each store daily = 4000
 - Promotion = 1
- Therefore, maximum number of fact table records = $1825 \times 300 \times 4000 \times 1 = 2190000000$ records

Ex. 1.6.5 : Suppose that a data warehouse for DB-University consists of the following four dimensions: student, course, semester, and instructor, and two measures, count and avg_grade. When at the lowest conceptual level (e.g., for a given student, course, semester, and instructor combination), the avg_grade measure stores the average grade for the given combination.

- At higher conceptual levels, avg_grade stores the average grade for the data warehouse.
- Draw a snowflake schema diagram for the data warehouse.
 - Starting with the base cuboid [student, course, semester, instructor], what specific OLAP operations (e.g., roll-up from semester to year) should one perform in order to list the average grade of CS courses for each DB_University student.

Soln. :

(a) Snowflake Schema



(145)Fig. P. 1.6.5 : Snowflake Schema for DB_University

(b) OLAP Operations to list the average grade of CS courses

- Roll-up on course from course_id to department.
- Roll-up on student from student_id to university.
- Dice on course, student with department = "CS" and university = "DB_University".
- Drill-down on student from university to student_name.

1.7 UPDATE TO THE DIMENSION TABLES

- The number of rows in the fact table continues to increase over time. Rows in a fact table are rarely updated with changes. Even if there are adjustments to the prior numbers, these are processed as additional adjustment rows and added to the fact table.
- Dimension tables are more stable and less volatile. However, unlike a fact table, which changes as the number of rows increases, a dimension table changes its attributes themselves change.

1.7.1 Slowly Changing Dimensions

- Slowly changing dimensions are the dimensions that change slowly over time, rather than changing on regular schedule. In time-base. In data warehouse there is a need to track changes in dimension attributes in order to report historical data. In other words, implementing one of the slowly changing dimension types should enable users assigning proper dimension's attribute value for given date. Example of such dimensions could be: customer, geography, employee.
- There are many approaches how to deal with slowly changing dimension. The most popular are:
 - Type 0 - The passive method
 - Type 2 - Creating a new additional record
 - Type 3 - Adding a new column
 - Type 4 - Using historical table
 - Type 6 - Combine approaches of types 1+2+3 (1+2+3=6)
- Type 0 - The passive method. In this method no special action is performed upon dimensional changes. Some dimension data can remain the same as it was first time inserted, others may be overwritten.
- Type 1 - Overwriting the old value. In this method no history of dimension changes is kept in the database. The old dimension value is simply overwritten by the new one. This type is easy to maintain and is often use for data in which changes are caused by processing corrections(e.g. removal of special characters, correcting spelling errors).

Before the change

Customer_ID	Customer_Name	Customer_Type
1	Cust_1	Retail

After the change

Customer_ID	Customer_Name	Customer_Type
1	Cust_1	Retail
1	Cust_2	Corporate

- Type 2 - Creating a new additional record. In this methodology all history of dimension changes is kept in the database. You capture attribute change by adding a new row with a new surrogate key to the dimension table. Both the prior and new rows contain as attributes the natural key(or other durable identifier). Also effective_date' and 'current_indicator' columns are used in this method. There could be only one record with current indicator set to 'Y'. For 'effective_date' columns, i.e. start_date and end_date, the end_date for current record usually is set to value 31-12-9999. Introducing changes to the dimensional model in type 2 could be very expensive database operation so it is not recommended to use it in dimensions where a new attribute could be added in the future.

Before the change

Customer_ID	Customer_Name	Customer_Type	Start Date	End Date	Current_Flag
1	Cust_1	Corporate	22-07-2010	31-12-9999	Y

After the change:

Customer_ID	Customer_Name	Customer_Type	Start Date	End Date	Current_Flag
1	Cust_1	Corporate	22-07-2010	17-05-2012	N
2	Cust_1	Retail	18-05-2012	31-12-9999	Y

- Type 3 - Adding a new column. In this type usually only the current and previous value of dimension is kept in the database. The new value is loaded into 'current/new' column and the old one into 'old/previous' column. Generally speaking, the history is limited to the number of column created for storing historical data. This is the least commonly needed technique.

Customer_ID	Customer_Name	Current_Type	Previous_Type
1	Cust_1	Corporate	

Customer_ID	Customer_Name	Current_Type	Previous_Type
1	Cust_1	Retail	Corporate

- Type 4 - Using historical table.** In this method a separate historical table is used to track all dimension's attribute historical changes for each of the dimension. The 'main' dimension table keeps only the current data e.g. customer and customer_history tables.

Customer_ID	Customer_Name	Customer_Type
1	Cust_1	Corporate

Customer_ID	Customer_Name	Customer_Type	Start_Date	End_Date
1	Cust_1	Retail	01-01-2010	21-07-2010
1	Cust_1	Other	22-07-2010	17-05-2012
1	Cust_1	Corporate	18-05-2012	31-12-9999

- Type 5 - Track changes in dimension table.** In this type we have in dimension table such additional columns as:

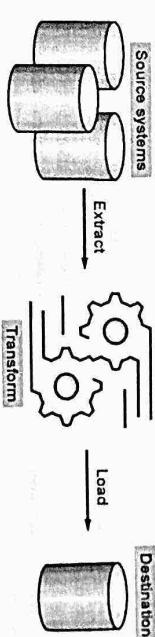
- current_type - for keeping current value of the attribute. All history records for given item of attribute have the same current value.
- historical_type - for keeping historical value of the attribute. All history records for given item of attribute could have different values.
- start_date - for keeping start date of 'effective date' of attribute's history.
- end_date - for keeping end date of 'effective date' of attribute's history.
- current_flag - for keeping information about the most recent record.

- In this method to capture attribute change we add a new record as in type 2. The current_type information is overwritten with the new one as in type 1. We store the history in a historical_column as in type 3.

Customer_ID	Customer_Name	Current_Type	Historical_Type	Start_Date	End_Date	Current_Flag
1	Cust_1	Corporate	Retail	01-01-2010	21-07-2010	N
2	Cust_1	Corporate	Other	22-07-2010	17-05-2012	N
3	Cust_1	Corporate	Corporate	18-05-2012	31-12-9999	Y

1.7.2 Rapidly Changing Dimension (RCD)

- A dimension is a fast changing or rapidly changing dimension if one or more of its attributes in the table changes very fast and in many rows. Handling rapidly changing dimension in data warehouse is very difficult because of many performance implications.



(Ans)Fig. 1.8.1 : ETL Process

- decision to implement the rapidly changing dimensions,
- For example : Consider patient dimension where there are 1000 rows in it. On average basis, each patient changes the 10 of attributes in a year. If you use the type 2 to manage this scenario, there will be $1000 * 10 = 10000$ rows. Imagine if the table has 1 million rows, it will become very hard to handle the situation with type 2. Hence we use rapidly changing dimension approach.

1.7.3 Conformed Dimension

- A conformed dimension is the dimension that is shared across multiple data mart or subject area. Company may use the same dimension table across different projects without making any changes to the dimension table.
- Conformed dimension example would be Customer dimension, i.e. both marketing and sales department can use Customer dimension for their reporting purpose.

1.7.4 Junk Dimension

- A junk dimension is a grouping of typically low cardinality attributes, so you can remove them from main dimension.
- You can use junk dimensions to implement the rapidly changing dimension where you can use it to stores the attribute that changes rapidly. For example, attributes such as flags, weights, BMI (body mass index), etc.

- The ETL process requires active inputs from various stakeholders, including developers, analysts, testers, top executives, and is technically challenging.
- To maintain its value as a tool for decision-makers, Data warehouse technique needs to change with business changes.

1.7.5 Degenerated Dimension

- A degenerated dimension is a dimension that is derived from fact table and does not have its own dimension

1.7.6 Role Playing Dimension

- Dimensions which are often used for multiple purposes within the same database are called role-playing dimensions. For example, you can use a date dimension for "date of sale", as well as "date of delivery", or "date of hire".

1.8 MAJOR STEPS IN ETL PROCESS

- The process of extracting data from source systems and bringing it into the data warehouse is commonly called ETL, which stands for extraction, transformation, and loading. Note that ETL refers to a broad process, and not three well-defined steps.
- It is a process in which an ETL tool extracts the data from various data source systems, transforms it in the staging area and then finally, loads it into the Data Warehouse system.
- The ETL process requires active inputs from various stakeholders, including developers, analysts, testers, top executives, and is technically challenging.
- ETL is a recurring method (daily, weekly, monthly) of a Data warehouse system and needs to be agile, automated, and well documented.
- ETL Tools : Most commonly used ETL tools are Sybase, Oracle Warehouse builder, Clover ETL and MarkLogic.

Let us understand each step of the ETL process in depth:

Extraction

- The first step of the ETL process is extraction. In this step, data from various source systems is extracted which can be in various formats like relational databases, NoSQL, XML and flat files into the staging area.

It is important to extract the data from various source systems and store it into the staging area first and not directly into the data warehouse because the extracted data is in various formats and can be corrupted also.

Hence loading it directly into the data warehouse may damage it and rollback will be much more difficult. Therefore, this is one of the most important steps of ETL process.

ETL Data Extraction Techniques

There are two types of data warehouse extraction methods: Logical and Physical extraction methods.

(A) Logical Extraction

Logical Extraction method in-turn has two methods:

(i) Full Extraction

- In this method, data is completely extracted from the source system. The source data will be provided as-is and no additional logical information is necessary on the source system. Since it is complete extraction, so no need to track source system for changes.
- For example, exporting complete table in the form of flat file.

(ii) Incremental Extraction

- In incremental extraction, the changes in source data need to be tracked since the last successful extraction. Only these changes in data will be extracted and then loaded. Identifying the last changed data itself is the complex process and involve many logics.
- You can detect the changes in the source system from the specific column in the source system that has the last changed timestamp. You can also create a change table in the source system, which keeps track of the changes in the source data.

(B) Physical Extraction

- Physical extraction has two methods: Online and Offline extraction.

(i) Online Extraction

In this process, extraction process directly connects to the source system and extract the source data.

(ii) Offline Extraction

- The data is not extracted directly from the source system but is staged explicitly outside the original source system.
- You can consider the following common structure in offline extraction:

- Flat file : Generic format
- Dump file : Database specific file

Transformation

- The second step of the ETL process is transformation. In this step, a set of rules or functions are applied on the extracted data to convert it into a single standard format.

It may involve following processes/tasks:

- Filtering – loading only certain attributes into the data warehouse.
- Cleaning – filling up the NULL values with some default values, mapping U.S.A, United States and America into USA, etc.
- Joining – joining multiple attributes into one.
- Splitting – splitting a single attribute into multiple attributes.
- Sorting – sorting tuples on the basis of some attribute (generally key-attribute).

Data Transformation Techniques

- The third and final step of the ETL process is loading. In this step, the transformed data is finally loaded into the data warehouse.

Sometimes the data is updated by loading into the data warehouse very frequently and sometimes it is done after longer but regular intervals.

The rate and period of loading solely depends on the requirements and varies from system to system.

Loading can be carried in two ways:

- Data Smoothing :** This method is used for removing the noise from a dataset. Noise is referred to as the distorted and meaningless data within a dataset. Smoothing uses algorithms to highlight the special features in the data. After removing noise, the process can detect any small changes to the data to detect special patterns.
- Data Aggregation :** Aggregation is the process of collecting data from a variety of sources and storing it in a single format. Here, data is collected, stored, analyzed and presented in a report or summary format. It helps in gathering more information about a particular data cluster. The method helps in collecting vast amounts of data.

Data Mining & Business Intelligence (MU-Sem 6-IT)

(Data Warehouses Fund. with Intro. to Data Mining)...Page no. (1-29)

Discretization : This is a process of converting continuous data into a set of data intervals. Continuous attribute values are substituted by small interval labels. This makes the data easier to study and analyze.

Generalization : In this process, low-level data attributes are transformed into high-level data attributes using concept hierarchies. For example, age data can be in the form of (20, 30) in a dataset. It is transformed into a higher conceptual level into a categorical value (young, old).

Attribute construction : In the attribute construction method, new attributes are created from an existing set of attributes. For example, in a dataset of employee information, the attributes can be employee name, employee ID and address. These attributes can be used to construct another dataset that contains information about the employees who have joined in the year 2019 only. This method of reconstruction makes mining more efficient and helps in creating new datasets quickly.

Normalization : Also called data pre-processing, this is one of the crucial techniques for data transformation in data mining. Here, the data is transformed so that it falls under a given range. When attributes are on different ranges or scales, data modelling and mining can be difficult. Normalization helps in applying data mining algorithms and extracting data faster.

OLAP (On-line Analytical Processing) is characterized by relatively low volume of transactions.

Queries are often very complex and involve aggregations. For OLAP systems a response time is an effectiveness measure. OLAP applications are widely used by Data Mining techniques. In OLAP database there is aggregated historical data, stored in multi-dimensional schemas (usually star schema). For example, a bank storing years of historical records of check deposits could use an OLAP database to provide reporting to business users. OLAP databases are divided into one or more cubes. The cubes are designed in such a way that creating and viewing reports become easy. At the core of the OLAP concept is an OLAP Cube. The OLAP cube is a data structure optimized for very quick data analysis. The OLAP Cube consists of numeric facts called measures which are categorized by dimensions. OLAP Cube is also called the hypercube.

Module 1.9 OLTP VS OLAP

We can divide IT systems as transactional (OLTP) and analytical (OLAP). In general, we can assume that OLTP systems provide source data to data warehouses, whereas OLAP systems help to analyze it.

OLTP (On-line Transaction Processing) is characterized by a large number of short on-line transactions (INSERT, UPDATE, DELETE). The main emphasis for OLTP systems is put on very fast query processing, maintaining data integrity in multi-access environments and an effectiveness measured by number of transactions per second. In OLTP database there is detailed and current data, and schema used to store transactional databases is the entity model (usually 3NF), *denormalized*.

OLAP (On-line Analytical Processing) is characterized by relatively low volume of transactions. Queries are often very complex and involve aggregations. For OLAP systems a response time is an effectiveness measure. OLAP applications are widely used by Data Mining techniques. In OLAP database there is aggregated historical data, stored in multi-dimensional schemas (usually star schema). For example, a bank storing years of historical records of check deposits could use an OLAP database to provide reporting to business users. OLAP databases are divided into one or more cubes. The cubes are designed in such a way that creating and viewing reports become easy. At the core of the OLAP concept is an OLAP Cube. The OLAP cube is a data structure optimized for very quick data analysis. The OLAP Cube consists of numeric facts called measures which are categorized by dimensions. OLAP Cube is also called the hypercube.

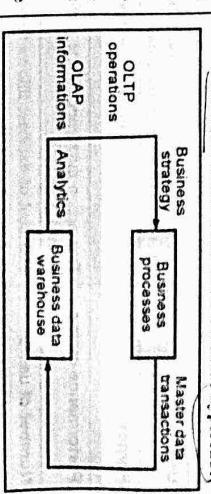


Diagram 1.9.1 : OLTP Vs OLAP Operations

- The following table summarizes the major differences between OLTP and OLAP system design.

Table 1.9.1 : OLTP vs OLAP

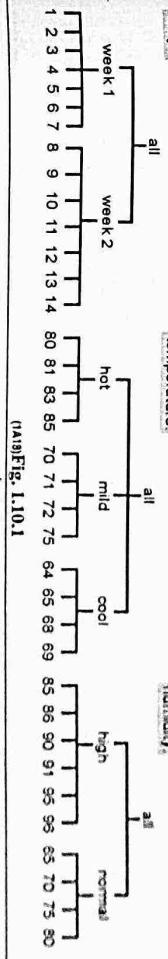
Parameters	OLTP	OLAP
<u>Process</u>	It is an online transactional system. It manages database modification.	OLAP is an online analysis and data retrieving process.
<u>Characteristic</u>	It is characterized by large numbers of short online transactions.	It is characterized by a large volume of data.
<u>Functionality</u>	OLTP is an online database modifying system.	OLAP is an online database query management system.
<u>Method</u>	OLTP uses traditional DBMS.	OLAP uses the data warehouse.
<u>Query</u>	Insert, Update, and Delete information from the database.	Mostly Select operations
<u>Table</u>	Tables in OLTP database are normalized.	Tables in OLAP database are not normalized.
<u>Source</u>	OLTP and its transactions are the sources of data.	Different OLTP databases become the source of data for OLAP.
<u>Storage</u>	The size of the data is relatively small as the historical data is archived. For e.g. MB, GB.	Large amount of data is stored typically in TB, PB.
<u>Data Integrity</u>	OLTP database must maintain data integrity constraint.	OLAP database does not get frequently modified. Hence, data integrity is not an issue.
<u>Response time</u>	It's response time is in millisecond.	Response time in seconds to minutes.
<u>Data quality</u>	The data in the OLTP database is always detailed and organized.	The data in OLAP process might not be organized.
<u>Usefulness</u>	It helps to control and run fundamental business tasks.	It helps with planning, problem-solving, and decision support.
<u>Operation</u>	Allow read/write operations.	Only read and rarely write.
<u>Audience</u>	It is a market oriented process.	It is a customer orientated process.
<u>Category Type</u>	Queries in this process are standardized and simple.	Complex queries involving aggregations.
<u>Back-up</u>	Complete backup of the data combined with incremental backups.	OLAP only need a backup from time to time. Backup is not important compared to OLTP
<u>Design</u>	DB design is application oriented. Example: Database design changes with industry like Retail, Airline, Banking, etc.	DB design is subject oriented. Example: Database design changes with subjects like sales, marketing, purchasing, etc.
<u>User type</u>	It is used by Data critical users like clerk, DBA & Data Base professionals.	Used by Data knowledge users like workers, managers, and CEO.
<u>Purpose</u>	Designed for real time business operations.	Designed for analysis of business measures by category and attributes.
<u>Performance metric</u>	Transaction throughput is the performance metric.	Query throughput is the performance metric.
<u>Number of users</u>	This kind of database allows thousands of users.	This kind of database allows only hundreds of users.
<u>Productivity</u>	It helps to increase user's self-service and productivity.	Help to increase productivity of the business analysis.

Parameters	OLTP	OLAP
<u>Challenge</u>	Data Warehouses historically have been a costly to build.	An OLAP cube is not an open SQL server data warehouse. Therefore, technical knowledge and experience is essential to manage the OLAP server.
<u>Process</u>	It provides fast result for daily used data.	It ensures that response to the query is quicker consistently.
<u>Characteristic</u>	It is easy to create and maintain.	It lets the user create a view with the help of a spreadsheet.
<u>Style</u>	OLTP is designed to have fast response time, low data redundancy and is normalized.	A data warehouse is created uniquely so that it can integrate different data sources for building a consolidated database

1.10 OLAP OPERATIONS

(S.M.) In the multidimensional model, data are organized into multiple dimensions, and each dimension contains multiple levels of abstraction defined by concept hierarchies. This organization provides users with the flexibility to view data from different perspectives.

For example, we have attributes as day, temperature and humidity, we can group values in subsets and name these subsets, thus obtaining a set of hierarchies as shown in Fig. 1.10.1.



- OLAP provides a user-friendly environment for interactive data analysis. A number of OLAP data cube operations exist to materialize different views of data, allowing interactive querying and analysis of the data.
- The most popular end user operations on dimensional data are:

1. Roll-Up

- The roll-up operation (also called drill-up or aggregation operation) performs aggregation on a data cube, either by climbing up a concept hierarchy for a dimension or by climbing down a concept hierarchy, i.e. dimension reduction. Let us explain roll up with an example:
- Consider the following cubes illustrating temperature of certain days recorded weekly:

Temperature	Cool	Mild	Hot
Week1	2	1	1
Week2	2	1	1

2. Drill-Down

- The drill-down operation (also called roll-down) is the reverse operation of roll-up. Drill-down is like zooming-in on the data cube. It navigates from less

- Detailed record to more detailed data. Drill-down can be performed by either stepping down a concept hierarchy for a dimension or adding additional dimensions.
- Figure shows a drill-down operation performed on the dimension time by stepping down a concept hierarchy which is defined as day, month, quarter, and year. Drill-down appears by descending the time hierarchy from the level of the quarter to a more detailed level of the month.
- Because a drill-down adds more details to the given data, it can also be performed by adding a new dimension to a cube. For example, a drill-down on the central cubes of the figure can occur by introducing an additional dimension, such as a customer group. Drill-down adds more details to the given data.

dimension. For example, a slice operation is executed when the customer wants a selection on one dimension of a three-dimensional cube resulting in a two-dimensional site. So, the slice operations perform a selection on one dimension of the given cube, thus resulting in a sub-cube. It will form a new subcubes by selecting one or more dimensions.

For example, if we make the selection, temperature = cool we will obtain the following cube:

Temperature	Cool	mild	Hot
Day 1	0	0	0
Day 2	0	0	0
Day 3	0	1	0
Day 4	0	1	0
Day 5	1	0	0
Day 6	0	0	0
Day 7	1	0	0
Day 8	1	0	0
Day 9	1	0	0
Day 10	0	0	0
Day 11	0	0	0
Day 12	0	0	0
Day 13	0	0	0
Day 14	0	0	0

Temperature	Cool
Day 1	0
Day 2	0
Day 3	0
Day 4	0
Day 5	1
Day 6	0
Day 7	1
Day 8	1
Day 9	1
Day 10	0
Day 11	0
Day 12	0
Day 13	0
Day 14	0

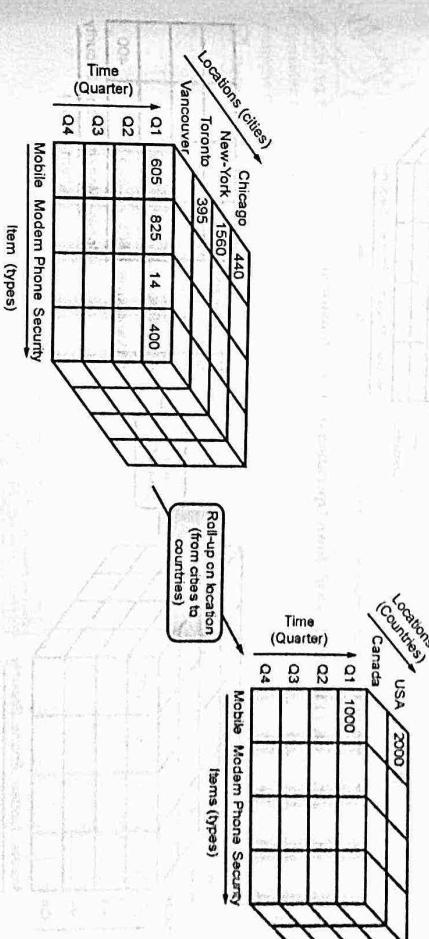
Dice

- The dice operation describes a sub-cube by operating a selection on two or more dimension
- For example, Implement the selection (time = day 3 OR time = day 4) AND (temperature = cool OR temperature = hot) to the original cubes we get the following sub-cube (still two-dimensional)

Temperature	Cool	hot
Day 3	0	1
Day 4	0	0

Pivot

- A slice is a subset of the cubes corresponding to a single value for one or more members of the row-dimensions into the column dimensions.

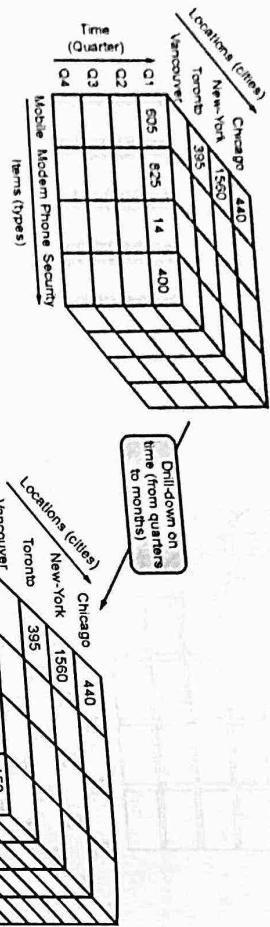


(142) Fig. 1.10.3 : Roll-up Operation on Location Dimension

NOTES

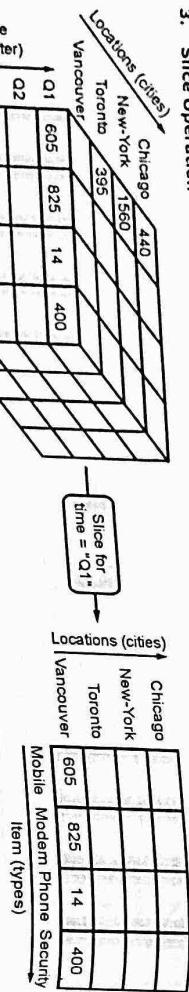
- The pivot operation is also called a rotation. Pivot is a visualization operation which rotates the data axes in view to provide an alternative presentation of the data. It may contain swapping the rows and columns or moving one of the row-dimensions into the column dimensions.

2. Drill-down Operation



(1421) Fig. 1.10.4 : Drill-down Operation on Time Dimension

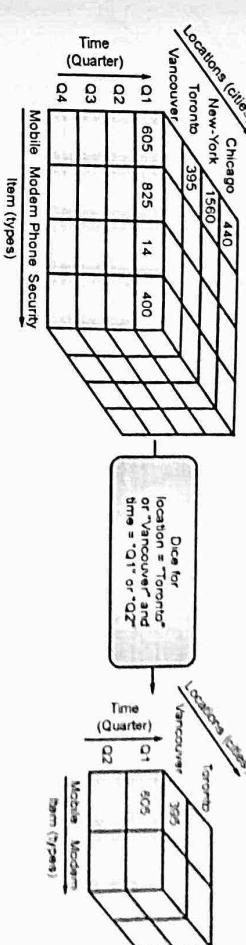
3. Slice Operation



(1422) Fig. 1.10.5 : Slice Operation for Time Dimension

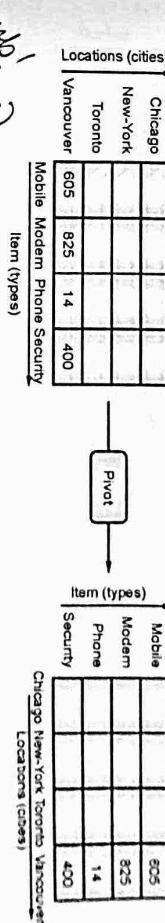
NOTES

4. Dice Operation



(1423) Fig. 1.10.6 : Dice Operation for Location and Time Dimensions

5. Pivot Operation

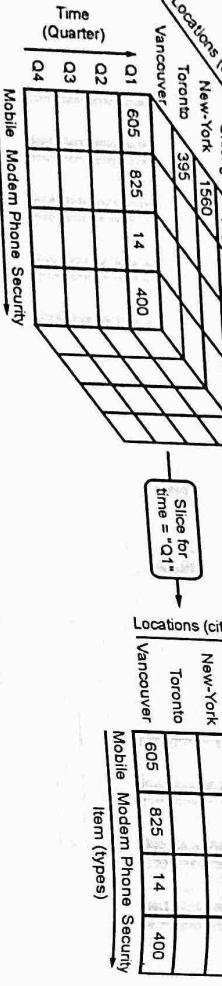


(1424) Fig. 1.10.7 : Pivot Operation for Location and Item Dimensions

Wanted (DM)
Ex 1.10.1: Consider a data warehouse for a hospital where there are three dimensions: a) Doctor b) Patient c) Time. Consider two measures(i) Count(ii) Charge where charge is the fee that the doctor charges a patient for a visit. For the above example create a cube and illustrate the following OLAP operations. 1) Rollup 2) Drill down 3) Slice 4) Dice 5) Pivot.

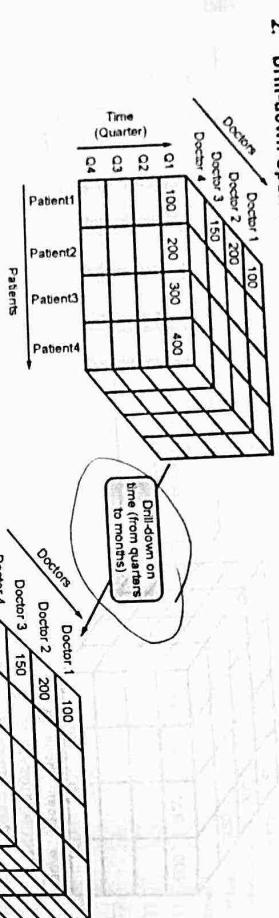
Soln. 1. Roll-up Operation

see one more such number from Pg 1-32



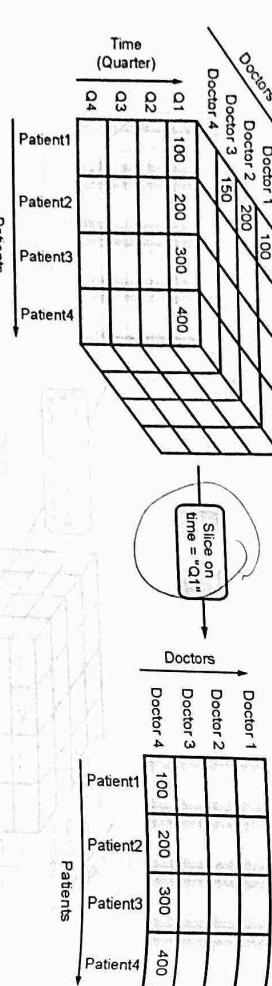
(1425) Fig. P. 1.10.1(a) : Roll-up Operation on Doctors Dimension

2. Drill-down Operation



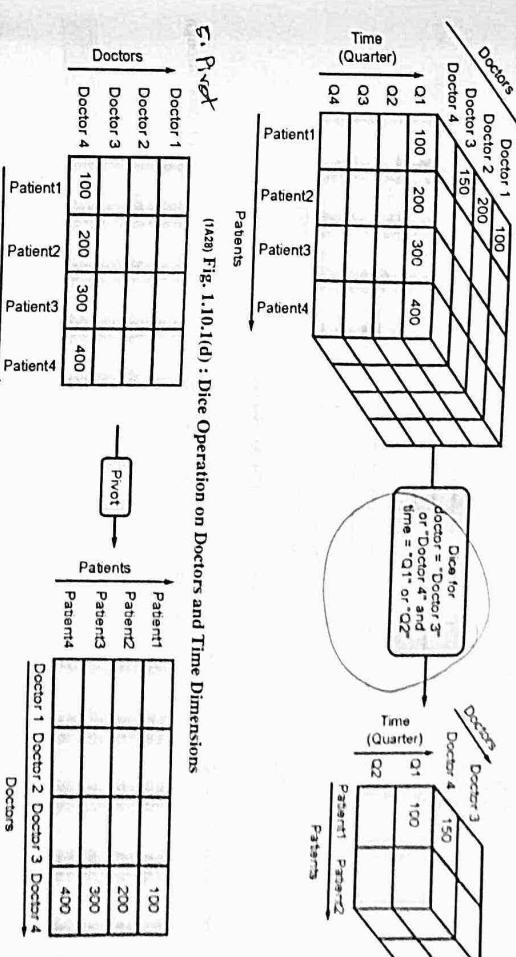
(1426) Fig. 1.10.1(b) : Drill-down Operation

3. Slice Operation

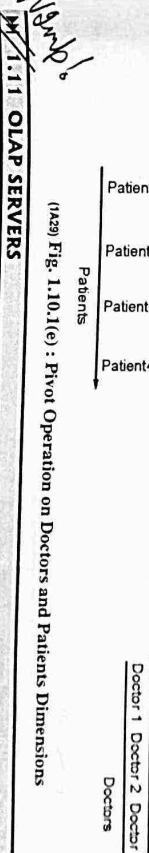


(1427) Fig. 1.10.1(c) : Slice Operation for Time Dimension

4. Dice Operation



(1428) Fig. 1.10.1(d) : Dice Operation on Doctors and Time Dimensions



(1429) Fig. 1.10.1(e) : Pivot Operation on Doctors and Patients Dimensions

1.1 OLAP SERVERS

There are three types of OLAP servers, namely, Relational OLAP (ROLAP), Multidimensional OLAP (MOLAP) and Hybrid OLAP (HOLAP).

1. Relational OLAP (ROLAP)

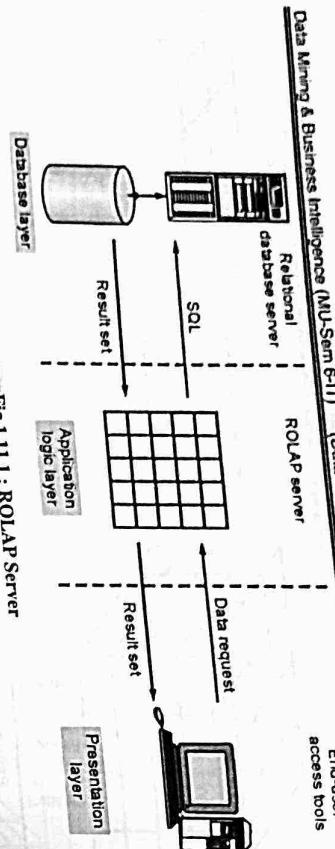
- Relational On-Line Analytical Processing (ROLAP) work mainly for the data that resides in a relational database, where the base data and dimension tables are stored as relational tables.
- ROLAP servers are placed between the relational back-end server and client front-end tools.
- ROLAP servers use RDBMS to store and manage warehouse data, and OLAP middleware to support missing pieces.
- Example : DSS Server of Microstrategy

2. Advantages of ROLAP

- ROLAP can handle large amounts of data.
- Can be used with data warehouse and OLTP systems.

3. Disadvantages of ROLAP

- Limited by SQL functionalities.
- Hard to maintain aggregate tables.

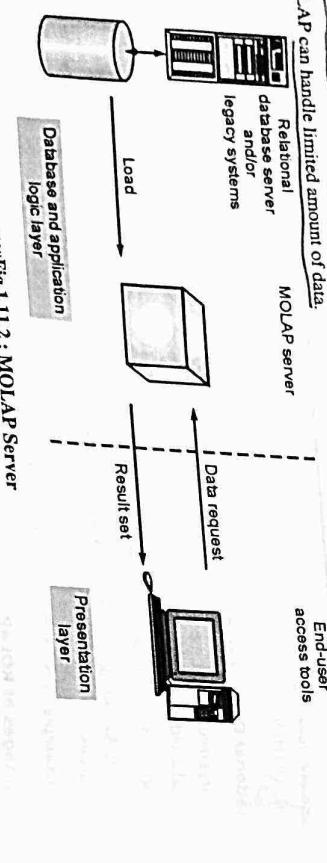


(11.1) Fig.1.11.1 : ROLAP Server

2. Multidimensional OLAP (MOLAP)

- Multidimensional On-Line Analytical Processing (MOLAP) support multidimensional views of data through array-based multidimensional storage engines.
- With multidimensional data stores, the storage utilization may be low if the data set is sparse.

Example : Oracle Essbase



(11.2) Fig.1.11.2 : MOLAP Server

1.1.12 APPLICATIONS OF OLAP

OLAP system is to analyze the business which helps in decision-making, forecasting, planning, problem solving. Some of the applications of OLAP include:

1. Financial Applications

- Resource (man-power, raw material) allocation
- Budgeting.

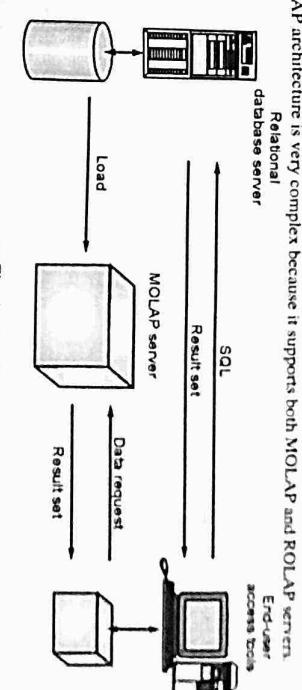
2. Sales Applications

- Research on market analysis
- Forecasting sales
- Analyzing sales promotions
- Analyzing customer requirements
- Dividing market based on customer.

3. Business Modelling

- Understanding and simulating the market trend and business behavior
- Decision support system for managers, executives, CEO, data scientists.
- Dividing market based on customer.

3. Hybrid OLAP (HOLAP)
- Hybrid On-Line Analytical Processing (HOLAP) is a combination of ROLAP and MOLAP.
 - HOLAP provide greater scalability of ROLAP and the faster computation of MOLAP.
 - Example: Microsoft SQL Server 2000



(11.3) Fig. 1.11.3 : HOLAP Server

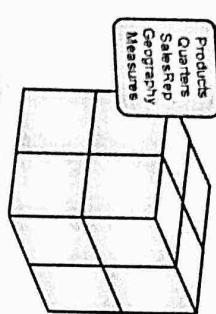
1.1.13 HYPERCUBE

- Multidimensional databases can present their data for an application using two types of cubes: hypercube and multi-cubes. In a hypercube, as shown in Fig. 1.13.1, all data appears logically as a single cube. All parts of the manifold represented by this hypercube have identical dimensionality. Each dimension belongs to one cube only. A dimension is owned by the hypercube. This simplicity makes it easy for users to understand.

- Designing a hypercube model is a top-down process with three major steps.

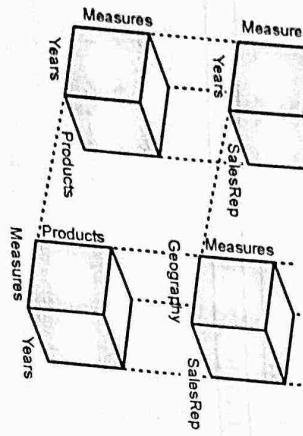
- You decide which process of the business you want to capture in the model, such as sales activity.
- You identify the values that you want to capture, such as sales amounts. This information is always numeric.
- You identify the granularity of the data, meaning the lowest level of detail at which you want to capture. These elements are the dimensions. And time, geography, product, and customer are some common dimensions. For example, a single cell in a cube could refer to the sales amount of Sony TVs in the first quarter of the year, in PA, USA.

1.14 AGGREGATE FACT TABLES



(1.13.1) Fig. Hypercube

- In the multi-cube model, data is segmented into a set of smaller cubes, each of which is composed of a subset of the available dimensions, as shown in Fig. 1.13.2. They are used to handle multiple fact tables, each with different dimensionality.



(1.13.2) Fig. Multi-cube

- The drawback is that this is less straightforward than hypercube and can carry steeper learning curves. Some systems use the combined approach of hypercube and multi-cubes, by separating the storage, processing, and presentation layers. It stores data as multi-cubes but presents as a hypercube.

1.15 INTRODUCTION TO DATA MINING

- Aggregate fact tables are special fact tables in a data warehouse that contain new metrics derived from one or more aggregate functions (AVERAGE, COUNT, MIN, MAX, etc.) or from other specialized functions that output totals derived from a grouping of the base data.
- These new metrics, called "aggregate facts" or "summary statistics" are stored and maintained in the data warehouse database in special fact tables at the grain of the aggregation.
- Likewise, the corresponding dimensions are rolled up, and condensed to match the new grain of the fact.
- These specialized tables are used as substitutes whenever possible for returning user queries. The reason is Speed. Querying a tidy aggregate table is much faster and uses much less disk I/O than the base, atomic fact table, especially if the dimensions are large as well.
- If you want to wow your users, start adding aggregates. You can even use this "trick" in your operational systems to serve as a foundation for operational reports.
- For example, take the "Orders" business process from an online catalog company where you might have customer orders in a fact table called FactOrders with dimensions Customer, Product, and OrderDate.
- With possibly millions of orders in the transaction fact, it makes sense to start thinking about aggregates.
- To further the above example, assume that the business is interested in a report: "Monthly orders by state and product type".
- While you could generate this easily enough using the FactOrders fact table, you could likely speed up the data retrieval for the report by at least half (but likely much, much more) using an aggregate.

1.15.1 Sources of Data that can be Mined

The data from multiple sources are integrated into a common source known as Data Warehouse. Let's discuss what type of data can be mined:

1. Flat Files
2. Relational Databases
3. Data Warehouse
4. Transactional Databases
5. Multimedia Databases
6. Spatial Database
7. Time-series Databases
8. WWW

- Flat files are defined as data files in text form or binary form with a structure that can be easily extracted by data mining algorithms.
- Data stored in flat files have no relationship or path among themselves, like if a relational database is stored on flat file, then there will be no relations between the tables.
- They can be stored on Object-Oriented Databases.
- They are used to store complex information in a pre-specified formats.
- Application : Digital libraries, video-on demand, news-on demand, musical database, etc.

- ✓ Christopher J Matheus define data mining as "The non-trivial extraction of implicit, previously unknown, and potentially useful information from data."
- The process of extracting information to identify patterns, trends, and useful data that would allow the business to take the data-driven decision from huge sets of data is called Data Mining.
- The computer is responsible for finding the patterns by identifying the underlying rules and features in the data.
- A data warehouse is defined as the collection of data integrated from multiple sources that will queries and decision making.
- There are three types of data warehouse: Enterprise areas such as data warehouses, efficient analysis, data mining algorithm, helping decision making and other data requirement to eventually cost-cutting and generating revenue.
- Application: Data Mining, ROLAP model etc.

- Logical schema in Relational databases is a schema which defines the relationship among tables.
- Standard API of relational database is SQL.
- Physical schema in Relational databases is a schema which defines the structure of tables.

- 2. Relational Databases
- A Relational database is defined as the collection of data organized in tables with rows and columns.
- Application: Business decision making, Data mining, etc.

- 3. Data Warehouse
- In other words, we can say that Data Mining is the process of investigating hidden patterns of information to various perspectives for categorization into useful data, which is collected and assembled in particular areas such as data warehouses, efficient analysis, data mining algorithm, helping decision making and other data requirement to eventually cost-cutting and generating revenue.
- Two approaches can be used to update data in Data Warehouse: Query-driven Approach and Update-driven Approach.
- Application: Business decision making, Data mining, etc.

- 4. Transactional Databases
- Transactional database is a collection of data organized by time stamps, date, etc. to represent transaction in databases.
- This type of database has the capability to roll back or undo its operation when a transaction is not completed or committed.
- Highly flexible system where users can modify information without changing any sensitive information.
- Follows ACID (Atomicity, Consistency, Isolation and Durability) property of DBMS.

- Application: Banking, Distributed systems, Object databases, etc.
- 5. Multimedia Databases
- Multimedia databases consists audio, video, images and text media.

► **6. Spatial Database**

- Store geographical information.
- Stores data in the form of coordinates, topology, lines, polygons, etc.
- Application: Maps, Global positioning, etc.

► **7. Time-series Databases**

- Time series databases contains stock exchange data and user logged activities.
- Handles array of numbers indexed by time, date, etc.
- It requires real-time analysis.

► **Application: eXtremeDB, Graphite, InfluxDB, etc.**

► **8. WWW**

- WWW refers to World wide web is a collection of documents and resources like audio, video, text, etc. which are identified by Uniform Resource Locators (URLs) through web browsers, linked by HTML pages, and accessible via the Internet network.
- It is the most heterogeneous repository as it collects data from multiple resources.
- It is dynamic in nature as Volume of data is continuously increasing and changing.
- Application: Online shopping, Job search, Research, studying, etc.

► **1.15.2 Kind of Patterns to be mined**
(Data Mining Techniques)

Data mining is highly effective, so long as it draws upon one or more of these techniques :

1. Tracking patterns
2. Classification
3. Association
4. Outlier detection
5. Clustering
6. Regression
7. Prediction

- 1. **Tracking patterns** : One of the most basic techniques in data mining is learning to recognize patterns in your data sets. This is usually a recognition of some aberration in your data happening at regular intervals, or an ebb and flow of a certain variable over time. For example, you might see that your sales of a certain product seem to spike just before the holidays, or notice that warmer weather drives more people to your website.

► **2. Classification** : Classification is a more complex data mining technique that forces you to collect various attributes together into discernible categories, which you can then use to draw further conclusions, or serve some function. For example, if you're evaluating data on individual customers' financial backgrounds and purchase histories, you might be able to classify them as "low," "medium," or "high" credit risks. You could then use these classifications to learn even more about those customers.

► **3. Association** : Association is related to tracking patterns, but is more specific to dependently linked variables. In this case, you'll look for specific events or attributes that are highly correlated with another event or attribute; for example, you might notice that when your customers buy a specific item, they also often buy a second, related item. This is usually what's used to populate "people also bought" sections of online stores.

► **4. Outlier detection** : In many cases, simply recognizing the overarching pattern can't give you a clear understanding of your data set. You also need to be able to identify anomalies, or outliers in your data. For example, if your purchasers are almost exclusively male, but during one strange week in July, there's a huge spike in female purchases, you'll want to investigate the spike and see what drove it, so you can either replicate it or better understand your audience in the process.

► **5. Clustering** : Clustering is very similar to classification, but involves grouping chunks of data together based on their similarities. For example, you might choose to cluster different demographics of your audience into different packets based on how much disposable income they have, or how often they tend to shop at your store.

► **6. Regression** : Regression, used primarily as a form of planning and modeling, is used to identify the likelihood of a certain variable, given the presence of other variables. For example, you could use it to project a certain price, based on other factors like availability, consumer demand, and competition. More specifically, regression's main focus is to help you uncover the exact relationship between two (or more) variables in a given data set.

► **1.15.3 Data Mining Technologies**

Several techniques are used in the development of data mining methods. Some of them are mentioned below :

1. Statistics
2. Machine learning
3. Information retrieval
4. Database systems and data warehouse
5. Decision support system

► **2. Machine learning**

Arthur Samuel defined machine learning as "a field of study that gives computers the ability to learn without being programmed".

When the new data is entered in the computer algorithms help the data to grow or change due to machine learning.

In machine learning, an algorithm is constructed to predict the data from the available database (Predictive analysis).

It is related to computational statistics.

The four types of machine learning are:

- (i) **Supervised learning** : It is based on the classification. It is also called as inductive learning. In this method, the desired outputs are included in the training dataset.

► **5. Decision support system**

Decision support system is a category of information system. It is very useful in decision making for organizations.

- (ii) **Unsupervised learning** : Unsupervised learning is based on clustering. Clusters are formed on the basis of similarity measures, and desired outputs are not included in the training dataset.
- (iii) **Semi-supervised learning** : Semi-supervised learning includes some desired outputs to the training dataset to generate the appropriate functions. This method generally avoids the large number of labeled examples (i.e. desired outputs).
- (iv) **Semi-supervised learning** : Active learning is a powerful approach in analyzing the data efficiently. The algorithm is designed in such a way that the desired output should be decided by the algorithm itself (the user plays important role in this type).

1.15.4 Difference between Data Mining and Data Warehouse

Table 1.15.1 : Data Mining Vs Data Warehouse

Sr.No.	Data Mining	Data Warehouse
1.	Data mining is the process of analyzing unknown patterns of data.	A data warehouse is database system which is designed for analytical instead of transactional work.
2.	Data mining is a method of comparing large amounts of data to finding right patterns.	Data warehousing is a method of centralizing data from different sources into one common repository.
3.	Data mining is usually done by business users with the assistance of engineers.	Data warehousing is a process which needs to occur before any data mining can take place.
4.	Data mining is considered as a process of extracting data from large data sets.	On the other hand, Data warehousing is the process of pooling all relevant data together.
5.	One of the most important benefits of data mining techniques is the detection and identification of errors in the system.	One of the pros of Data Warehouse is its ability to update consistently. That's why it is ideal for the business owner who wants the best and latest features.
6.	Data mining helps to create suggestive patterns of important factors. Like the buying habits of customers, products, sales. So that, companies can make the necessary adjustments in operation and production.	Data Warehouse adds an extra value to operational business systems like CRM systems when the warehouse is integrated.
7.	The Data mining techniques are never 100% accurate and may cause serious consequences in certain conditions.	In the data warehouse, there is great chance that the data which was required for analysis by the organization may not be integrated into the warehouse. It can easily lead to loss of information.
8.	The information gathered based on Data Mining by organizations can be misused against a group of people.	Data warehouses are created for a huge IT project. Therefore, it involves high maintenance system which can impact the revenue of medium to small-scale organizations.
9.	After successful initial queries, users may ask more complicated queries which would increase the workload.	Data Warehouse is complicated to implement and maintain.
10.	Organisations can benefit from this analytical tool by equipping pertinent and usable knowledge-based information.	Data warehouse stores a large amount of historical data which helps users to analyze different time periods and trends for making future predictions.
11.	Organisations need to spend lots of their resources for training and Implementation purpose. Moreover, data mining tools work in different manners due to different algorithms employed in their design.	In Data warehouse, data is pooled from multiple sources. The data needs to be cleaned and transformed. This could be a challenge.
12.	The data mining methods are cost-effective and efficient compares to other statistical data applications.	Data warehouse's responsibility is to simplify every type of business data. Most of the work that will be done on user's part is inputting the raw data.
13.	Another critical benefit of data mining techniques is the identification of errors which can lead to losses. Generated data could be used to detect a drop-in sale.	Data warehouse allows users to access critical data from the number of sources in a single place. Therefore, it saves user's time of retrieving data from multiple sources.
14.	Data mining helps to generate actionable strategies built on data insights.	Once you input any information into Data warehouse system, you will unlikely to lose track of this data again. You need to conduct a quick search, helps you to find the right statistic information.

1.16 DATA MINING TASK PRIMITIVES

- Each user will have a data mining task in mind, that is, some form of data analysis that he or she would like to have performed.
- A data mining task can be specified in the form of a data mining query, which is input to the data mining system.
- A data mining query is defined in terms of data mining task primitives.
- These primitives allow the user to interactively communicate with the data mining system during discovery in order to direct the mining process, or examine the findings from different angles or depths.
- Here is the list of Data Mining Task Primitives.

1. The set of task-relevant data to be mined :

This specifies the portions of the database or the set of data in which the user is interested. This includes the database attributes or data warehouse dimensions of interest (referred to as the relevant attributes or dimensions).

2. The kind of knowledge to be mined :

This specifies the data mining functions to be performed, such as characterization, discrimination, association or correlation analysis, classification, prediction, clustering, outlier analysis, or evolution analysis.

3. The background knowledge to be used in the discovery process :

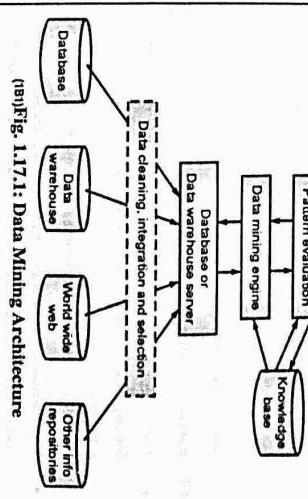
This knowledge about the domain to be mined is useful for guiding the knowledge discovery process and for evaluating the patterns found. Concept hierarchies are a popular form of background knowledge, which allow data to be mined at multiple levels of abstraction.

4. The interestingness measures and thresholds for pattern evaluation :

They may be used to guide the mining process or, after discovery, to evaluate the discovered patterns. Different kinds of knowledge may have different interestingness measures. For example, interestingness measures for association rules include support and confidence. Rules whose support and confidence

1.17 DATA MINING ARCHITECTURE

- values are below user-specified thresholds are considered uninteresting.
- 5. The expected representation for visualizing the discovered patterns : This refers to the form in which discovered patterns are to be displayed, which may include rules, tables, charts, graphs, decision trees, and cubes.



(B) Fig. 1.17: Data Mining Architecture

E3 Data Cleaning, Integration and Selection

- Before passing the data to the database or data warehouse server, the data must be cleaned, integrated and selected.
- As the information comes from various sources and in different formats, it can't be used directly for the data mining procedure because the data may not be complete and accurate. So, the first data requires to be cleaned and unified.
- More information than needed will be collected from various data sources, and only the data of interest will have to be selected and passed to the server. These procedures are not as easy as we think. Several methods may be performed on the data as part of selection, integration, and cleaning.

E3 Data Warehouse Fund. with Intro. to Data Mining)....Page no. (1-47)

- On the other hand, the pattern evaluation module might be coordinated with the mining module, depending on the implementation of the data mining techniques used. For efficient data mining, it is abnormally suggested to push the evaluation of pattern stake as much as possible into the mining procedure to confine the search to only fascinating patterns.

E3 Graphical User Interface

- The graphical user interface (GUI) module communicates between the data mining system and the user. This module helps the user to easily and efficiently use the system without knowing the complexity of the process.
- This module cooperates with the data mining system when the user specifies a query or a task and displays the results.

E3 Knowledge Base

- The knowledge base that contains the domain knowledge used to guide the search or to evaluate the interestingness of resulting patterns. For example, the knowledge base may contain metadata which describes data from multiple heterogeneous sources.
- A data mining engine, which consists of a set of functional modules for tasks such as classification, association, classification, cluster analysis, and evolution and deviation analysis.
- A pattern evaluation module that works in tandem with the data mining modules by employing interestingness measures to help focus the search towards interestingness patterns.

E3 Data Selection

- Strategies to handle missing data fields.
- 2. **Data Integration**
- Data from various sources such as databases, data warehouse, and transactional data are integrated.
- Multiple data sources may be combined into a single data format.

E3 Following steps :

1. Data Cleaning
2. Data Integration
3. Data Selection
4. Data Transformation
5. Data Mining
6. Pattern Evaluation
7. Knowledge Presentation

E3 Module

E3 1

E3 Data Mining Engine

E3 Database or Data Warehouse Server

- The database or data warehouse server consists of the original data that is ready to be processed. Hence, the server is cause for retrieving the relevant data that is based on data mining as per user request.

E3 Data Mining Engine

- The data mining engine is a major component of any data mining system. It contains several modules for operating data mining tasks, including association, characterization, classification, clustering, prediction, time-series analysis, etc.

E3 Knowledge Base

- The knowledge base is helpful in the entire process of data mining. It might be helpful to guide the search or evaluate the stake of the result patterns.
- The knowledge base may even contain user views and data from user experiences that might be helpful in the data mining process.

E3 Knowledge Base

- The pattern assessment module regularly interacts with the knowledge base to get inputs, and also update it.
- The data mining engine may receive inputs from the knowledge base to make the result more accurate and reliable.

E3 Knowledge Base

- The MU – Dec 2019 pattern assessment module receives inputs from the knowledge base to get inputs, and also update it.

E3 Knowledge Base

- What is data mining? Explain KDD process with diagram.

E3 Knowledge Base

- MU – May 2019

E3 Knowledge Base

- Knowledge discovery in the database (KDD) is the process of searching for hidden knowledge in the massive amounts of data that we are technically capable of generating and storing.

E3 Knowledge Base

- The basic task of KDD is to extract knowledge (or information) from a lower-level data (databases).

E3 Knowledge Base

- It is the non-trivial (significant) process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data.

E3 Knowledge Base

- To identify the truly interesting patterns representing knowledge based on interesting measures.

E3 Knowledge Base

- Visualizations and knowledge representation techniques are used to present mined knowledge to users.

E3 Knowledge Base

- Visualizations can be in form of graphs, charts or table.

