

Assignment - 01

Date: 01/08/2020

Title: Analyzing and Extracting data using ETL tools

Problem : For an organization of your choice, choose a set of business processes. Design Star/snow flake schemas for analyzing these processes create a fact constellation schema by combining them. Extract data from different data sources, apply suitable transformations and load data into destination table using ETL tools. For eg. Business Organization sales, ordering, Marketing process.

Objective: • Implementing of problem statement using ETL tool.

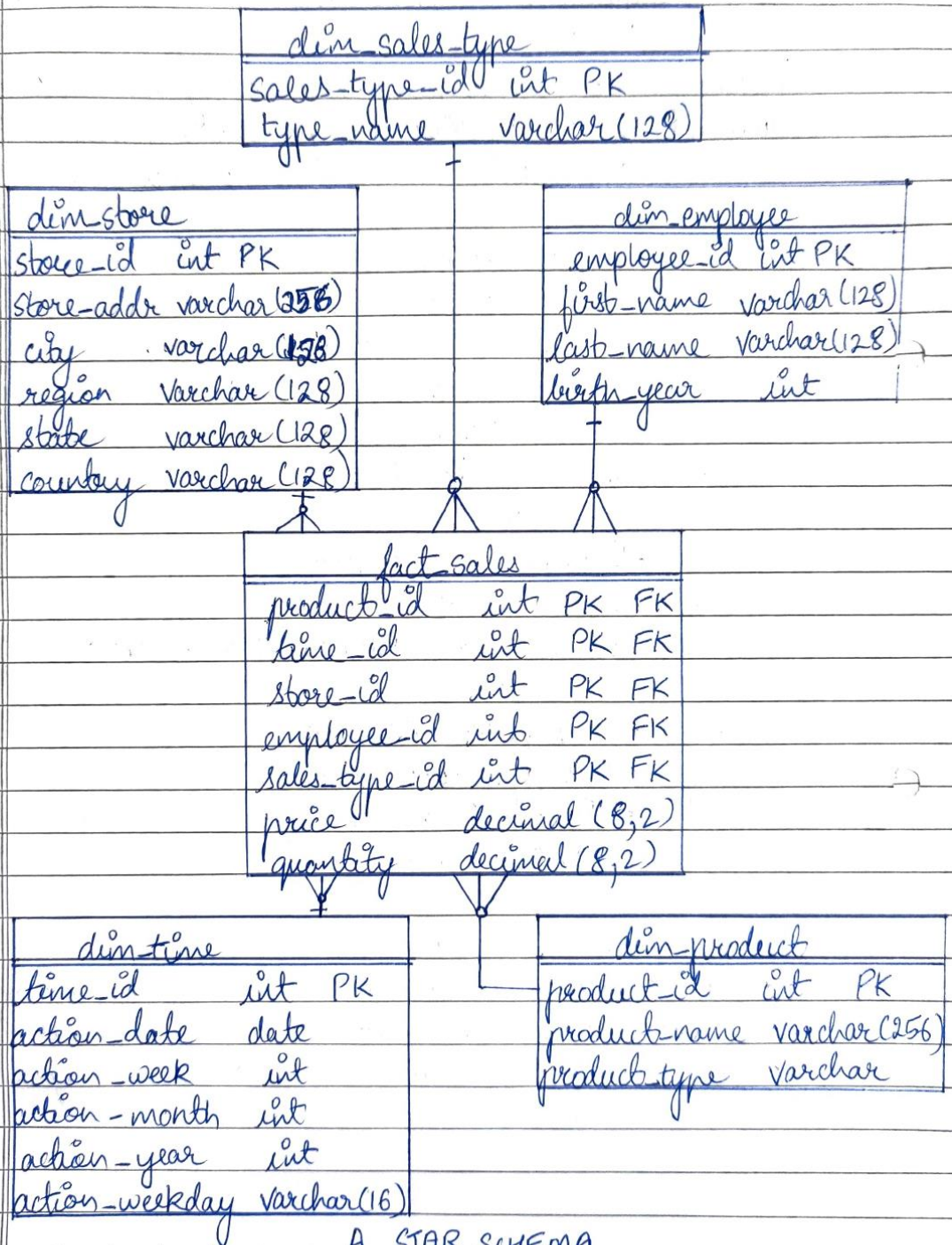
• Star/snow flake schema for analyzing process

S/W & : • PIV, 2 GB RAM, HDD

H/W req. • ETL open source tool Pentaho.
• Tomcat 8.0x with Oracle Java 8-x
• MySQL 5.6 & 5.7 (SQL 92)

Theory: A) Star Schemas.

• The schemas are a way to organize data marts or entire data warehouse using relational databases.
• Consider the following sales model represented in star schema.



A STAR SCHEMA

FOR EDUCATIONAL USE

→ Characteristics of Star schema:

- Every dimension is represented with only one-dimensional table.
- Fact table would contain key & measure.
- It is easy to understand & provides optimal disk usage.
- It is widely supported by BI tools.
- The dimension tables are not joined to each other.

B] A snowflake schema.

- It is an extension of star schema, and it adds additional dimensions.
- The dimension tables are normalized which splits data into additional table.

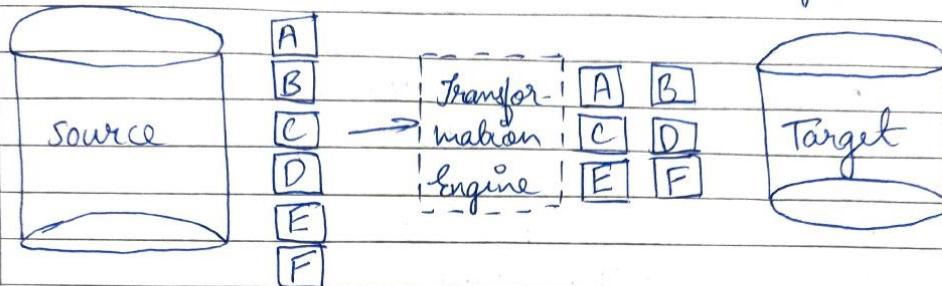
→ Characteristics of snowflake schema.

- The main benefit is that it uses smaller disk space.
- Easier to implement, a dimension is added to schema.
- Due to multiple tables, query performance is reduced.
- Need to perform more maintenance efforts because of more lookup tables.

C] ETL (Extract, Transform, Load)

- ETL is an abbreviation for Extract, Transform & Load.
- In this process an ETL tool extracts the data from different RDBMS source systems, then transforms the data by applying calculations, concatenations, etc. and load the data into data warehouse system.

- In ETL, data flows from source to target.



- ETL is a different method of looking at the tool approach to data movement.
 - Instead of transforming data before it's written, ETL lets the target system to do transformation.
 - The data is first copied to the target and then transformed in place.
 - ETL is usually used with no-sql databases like Hadoop Cluster, data appliance or cloud installation.
- List of open source ETL tools
- Clover ETL
 - Tedon
 - Pentaho
 - Talend

Test Cases :

Srno.	Description.	Expected O/P	Actual O/P
1)	Xampp/Apache Server installation	installed successfully	starts Apache & MySQL server.
2)	While installing pentaho, make sure to set PENTHO_JAVA_HOME success	success	success
2)	PENTHO-INSTALLERLICENSE.PAT environment variables		

Srno	Description.	Expected O/P	Actual O/P
3)	Perform transformation on the postal codes	successfully implements	successfully implemented.
4)	Perform transformation on missing-zipcode	Null values & invalid data is removed	success

Conclusion: Thus, we have learned to extract data from different data sources, apply suitable transformations & load into destination table using ETL tool.