

Data Loading (ETL)

DATA SET

MySQL Workbench

pw_skills_DA_Batch x

File Edit View Query Database Server Tools Scripting Help

Navigator ADVANCE SQL ASSIGNMENT FUNCTIIONS ASSIGNMENT* ETL FUNCTION* SQL File 7* SQL File 8 SQL File 6* SQL File 9 SQL File 10* SQL File 11* classess data loading as

SCHEMAS

Filter objects

- college
- company_db
- company_md
- data_loading
 - Tables
 - Views
 - Stored Procedures
 - Functions
- db_name
- employee
- etl_practice
- mitesh
- orders
- product
- pw_skills
 - Tables
 - analysis
 - classes
 - customers
 - department
 - orders
 - nrndurte

Administration Schemas

Information Schema: orders

```
10 VALUES
11 ('0101', 'C001', '4500', '12-01-2024'),
12 ('0102', 'C002', NULL, '15-01-2024'),
13 ('0103', 'C003', '3200', '2024/01/18'),
14 ('0101', 'C001', '4500', '12-01-2024'),
15 ('0104', 'C004', 'Three Thousand', '20-01-2024'),
16 ('0105', 'C005', '5100', '25-01-2024');
17 • select * from orders;
18
19
```

Result Grid

Order_ID	Customer_ID	Sales_Amount	Order_Date
O101	C001	4500	12-01-2024
O102	C002	NULL	15-01-2024
O103	C003	3200	2024/01/18
O101	C001	4500	12-01-2024
O104	C004	Three Thousand	20-01-2024
O105	C005	5100	25-01-2024

orders 3 x

Output

Action Output

#	Time	Action	Message	Duration / Fetch
✓ 13	19:19:00	CREATE TABLE orders (Order_ID VARCHAR(10), Customer_ID VARCHAR(10), Sales_Amount VA...	0 row(s) affected	0.031 sec
✓ 14	19:19:39	INSERT INTO orders (Order_ID, Customer_ID, Sales_Amount, Order_Date) VALUES ('0101', 'C001', '4500'...	6 row(s) affected Records: 6 Duplicates: 0 Warnings: 0	0.000 sec
✓ 15	19:19:53	select * from orders	6 row(s) returned	0.000 sec / 0.000 sec

QUESTION 1		Data Understanding Identify all data quality issues present in the dataset that can cause problems during data loading?					
ANSWER		Duplicate Records	O101 appears twice. Entire row is duplicated.				
		Missing Values (NULL)	Sales_Amount for C002 is Null. If column is NOT NULL, data load fails.				
		Inconsistent Date Format	Different formats: 12-01-2024 2024/01/18				
		Incorrect Data Type (Mixed Data)	Sales_Amount contains: Text value → "Three Thousand" Numeric values (4500, 3200)				
		Inconsistent Structure (Header Issue)	Order_ID separated from other columns.				
		No Defined Primary Key	If no primary key: Duplicate data possible No uniqueness enforcement				

QUESTION 2	Primary Key Validation Assume order_id is the Primary Key. a) Is the dataset violating the Primary Key rule? b) Which record(s) cause this violation?		
ANSWER	<div data-bbox="530 351 2076 509"> <p>Order_ID is the Primary Key A Primary Key must:</p> <p>Be unique</p> <p>Not be NULL</p> </div> <div data-bbox="530 517 2076 635"> <p>Is the dataset violating the Primary Key rule?</p> <p>Yes, it is violating the Primary Key rule.</p> <p>Because a Primary Key must contain unique values, and here we have a duplicate.</p> </div> <div data-bbox="530 683 2076 769"> <p>Which record(s) cause this violation?</p> <p>the duplicate values is 0101</p> </div> <div data-bbox="530 849 2076 920"> <p>Primary key cannot repeat.</p> <p>SQL will reject the second 0101 during insert if PK is applied.</p> </div>		

QUESTION 3		Missing Value Analysis								
		Which column(s) contain missing values?								
		a) List the affected records b) Explain why loading these records without handling missing values is risky								
ANSWER										
	A)	The column with missing value is:			Sales_Amount					
		O102 / C002 / NULL //15-01-2024								
	B)	NULL values may reduce total revenue.								
		Business may think sales are lower than actual.								
		Incorrect KPIs								
		If Sales_Amount is used to calculate:								
		Total Revenue								
		Average Order Value								
		Monthly Sales								

QUESTION 4

Data Type Validation

Identify records where Sales_Amount violates expected data type rules.
a) Which record(s) will fail numeric validation?
b) What would happen if this dataset is loaded into a SQL table with

ANSWER

A)

The record that violates numeric rules is:

Order_ID	Customer_ID	Sales_Amount	Order_Date
O104	C004	Three Thousand	20-01-2024

"Three Thousand" is text
It cannot be converted into INT or DECIMAL automatically

B)

What would happen if this dataset is loaded into a SQL table with:

Sales_Amount INT

4500 → Inserted
3200 → Inserted
5100 → Inserted
NULL → Inserted
"Three Thousand" → ERROR

QUESTION 5

Date Format Consistency

THE ORDERD DATE column has multiple formats.
a) List all date formats present in the dataset
b) Why is this a problem during data loading?

ANSWER

A) List all date formats present in the dataset

There are two different date formats:

Format 1 :: DD-MM-YYYY

FORMAT 2 :YYYY/MM/DD

Examples:

12-01-2024
15-01-2024
20-01-2024
25-01-2024

Example:

2024/01/18

B) Data Type Conversion Failure

Order_Date DATE

Mixed formats may cause:	Conversion errors
	Failed insert statements

12 January 2024? (DD-MM-YYYY)

Or December 1, 2024? (MM-DD-YYYY)

		5: Validate Column Constraints						
		Set Order_ID → PRIMARY KEY						
		Set Sales_Amount → NOT NULL (if required)						
		Set Customer_ID → FOREIGN KEY (if customer table exists)						
		6: Trim & Standardize Text Fields						
		Remove extra spaces						
		Ensure consistent case formatting						
		Check for hidden characters						
		7: Business Rule Validation						
		Ensure Sales_Amount > 0						
		Ensure no future dates (if rule applies)						
		Ensure each order belongs to valid customer						

QUESTION 9	<p>Loading Strategy Selection Assume this dataset represents daily sales data.</p> <p>a) Should a Full Load or Incremental Load be used?</p> <p>b) Justify your choice</p>
ANSWER	<p>Incremental Load should be used.</p> <p>Entire historical data will reload every day. It wastes time and system resources.</p> <p>Incremental load: Loads only new records for the day. Much faster and more efficient.</p> <p>Performance & Scalability</p> <p>As data grows:</p> <p>Full Load becomes slow. Database locking may increase. ETL window may exceed allowed time.</p> <p>Incremental load:</p> <p>Handles large datasets better. Scales easily over time.</p> <p>Reduces Risk of Duplicate Data</p> <p>With proper incremental logic (using Order_Date or Order_ID):</p>

		Only new records are inserted.							
		Existing records are not reloa							

QUESTION 10	BI Impact Scenario Assume this dataset was loaded without cleaning and connected to a BI dashboard. a) What incorrect results might appear in Total Sales KPI? b) Which records specifically would cause misleading insights? c) Why would BI tools not detect these issues automatically?
ANSWER	<p>1) Inflated Total Sales</p> <p>O101 is duplicated Sales amount 4500 will be counted twice So revenue will be higher than actual.</p> <p>2) Underestimated Total Sales</p> <p>One record has NULL Sales_Amount NULL is ignored in SUM() calculation</p> <p>3) Missing Sales Value</p> <p>"Three Thousand" is text</p> <p>BI tool may:</p> <p>Ignore it Treat it as 0</p> <p>c) Why would BI tools not detect these issues automatically?</p> <p>BI Tools Assume Data is Clean</p> <p>BI tools (like dashboards) focus on: Visualization Aggregation Reporting</p> <p>Duplicate Data Looks Legitimate</p> <p>Two identical rows look valid. BI tool cannot know: It is duplicate Or a real repeated transaction</p>