

Bright Light Data Analytics

Research Assignment

Due date: 30 October 2023

Section A: Database fundamentals

1. What are the main types of databases?
Relational database, NoSQL databases, object-oriented databases, time-series database

2. What is Relational Database Management System (RDBMS)?
is a program used to create, update, and manage relational database

3. What is a primary key and a foreign key in a database?

A primary key uniquely identifies a record in a table. A foreign key is a field that is used to link one table to the primary key of another table.

4. What is database normalization and why is it important?

is the process of organizing data to reduce redundancy, improve data integrity and make it more efficient

5. How is database schema?

A schema is the blueprint or structure of a database that defines its tables

fields and relationships

6. Differentiate between structured, semi-structured and unstructured.

- Structured: organized into rows and columns (e.g. SQL tables)

- Semi-structured: has some structure but not fixed

- Unstructured: No predefined format e.g. image

7. Difference between a fact table and a dimension table

Fact table: Stores numerical measures

Dimension table: Store descriptive information (e.g. product, date)

8. What is a data model and why is it important?

A data model defines how data is structured and related. It guides database

design and ensures data is consistent.

9. Difference between a database, a data warehouse and a data lake.

- Database: Store current operational data.

- Data warehouse: Store historical, structured data for analysis.

- Data lake: Stores all types of raw data - structured and unstructured.

10. What is a data mart?
It differs from

A data mart of a specific user
specific team

Section B: SQL

11. What is a query language?
SQL is the most

A query language
from a database
because it is supported by

12. What are indexes?
They improve
indexes or
searching
much data

13. What are transactions?
What are
A transaction
but as

A. Atomic
C. Consistent
I. Isolated
D. Durability

14. What is

10. what is data mart, and how does it differ from a data warehouse?

A data mart is a smaller, focused version of a data warehouse, created for a specific team or department.

Section B: SQL and data processing.

11. what is a query language, and why is SQL the most commonly used?

A query language is used to ask questions from a database. SQL is the most common because it is standardized, powerful and supported by all major db databases.

12. what are indexes in databases and how do they improve performance?

Indexes are structures that speedup searching and filtering by reducing how much data the database scans

13. what are transactions in database and what are the ACID properties?

A transaction is a group of operations that run as one:

A - Atomicity: All or nothing

C - Consistency: Must follow rules

I - Isolation: Transaction don't affect each other

D - Durability: changes are permanent

14. what is a database engine, and how does

it impact performances?

The engine is the core software that stores, reads, writes and manages data. A faster engine improves speed, reliability and efficiency.

15. What are views, stored procedures and triggers in SQL?

View: A virtual table based on a query

Stored procedure: A saved SQL script that can be run anytime

Trigger: Code that runs automatically when an event happens

16. Difference between ETL and ELT.

ETL: Extract - transform - Load (transform before loading)

ELT: Extract - load - transform (transform inside Data warehouse)

17. Difference between batch processing and stream processing.

Batch: processing large chunks of data at scheduled times

Stream: processes data continuously in real time

18. Explain what a join is in SQL. List types of joins.

A Join combines rows from 2 tables
types: INNER join, Left join, Right join, full

Outer join, cross join

19. What is referential integrity and why is it important?

It ensures relationship between tables remain valid preventing orphan records and maintaining accuracy.

20 How does data redundancy affect performance and storage?

Redundant data uses more storage slows queries, and increases inconsistency.

Section C: Data Management and Analytics Concepts

Q1. How does cloud database management differ from on-premise database.

Cloud databases are hosted online and managed by providers, while on-premise databases run on local servers and require internal maintenance.

Q2. What is data governance, and why is it important?

Data integrity means data is accurate and reliable. It is maintained through validation constraints, backups and access control.

Q3. What is data integrity, and how can it be maintained.

Data integrity means data is accurate and reliable. It is maintained through validation, constraints, backups and access control.

Q4. What is data quality and why it is critical?

Data quality means data is accurate and reliable. It is maintained through validation, constraints, backups and access control.

Q5. Explain the role of a Data Analyst.

A Data Analyst collects, cleans, organizes and analyzes to give insights to support business decisions

Q6. What are the key responsibilities of a database administrator (DBA)?

Installing database, managing users, performing backups, monitoring performance, securing data and optimizing queries

Q7. What are the main steps involved in designing a data pipeline.

Requirement - Data collection - processing - transformation - storage - Monitoring.

Q8. What are some common challenges in managing large-scale databases?

Scalability - performance slowdown, high

storage, data security, backup, and ensuring data quality.

29 popular database platforms and their use cases:

- MySQL - web apps

- Snowflake - analytics and cloud data warehousing

- PostgreSQL - advanced relational workload

- Oracle - enterprise level systems

30 main data storage formats used in analytics:

SS

CSV, JSON, Avro