Star and Snowflake Schema

1. Potential Interview Questions

Potential interview question:-

① What is star Schema?

② What is snowflake schema?

③ What is normalization ?

④ What is denormalization ?

⑤ Advantage of star schema over snow flake schema & vice versa ?

2.

fact Table → Measurement

Dim Table → Context.

3. Here we know that

4. Consider this table..here we have both dimension table and fact table column

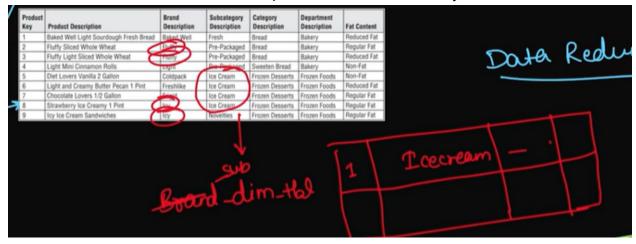| Id | Date | Product Name | Sales-quantity | Customer Name | Contact | email |
|---|---|---|---|---|---|---|
| 1 | 20-11-23 | Ice-cream | 2 | Manish kumar | 1234567 | m@gmai |
| 2 | 22-11-23 | Strawberry | 5 | Manish kumar | 1234567 | m@g |
| 3 | 22-11-23 | Bread | 2 | Manish kumar | 1234567 | m@ |
| 4 | 27-11-23 | Ice cream | 1 | Raushan Singh | 2345678 | r@ |
| 5 | 26-11-23 | Bread | 5 | Rahul Patil | 2356789 | |
| 6 | 24-11-23 | Strawberry | 7 | Rahul patil | 2356785 | |
| 7 | 24-11-23 | Cinnamon | 1 | Pritam Das | 517517 | |
| 8 | 20-11-23 | Ice-cream | 2 | Pritam Das | 5175 | |

5. Date,Product_name and customer_name can act as dim table
6. Now let us know why our transactional table..must be divided into facts and dims
7. Here we can see our data in customer_name and product name is redundant..so we'll store this data in another table
8. So here we have date dim table,customer dim and product_dim table
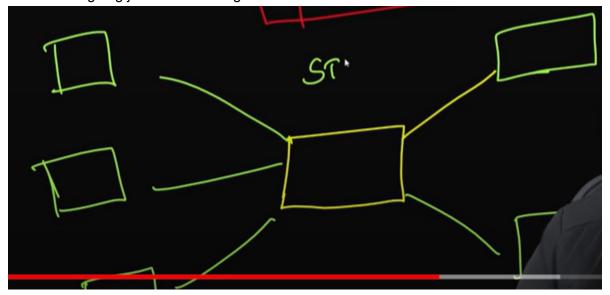9. Now in Star Schema our fact table…will be joined with dimensional tables

10.

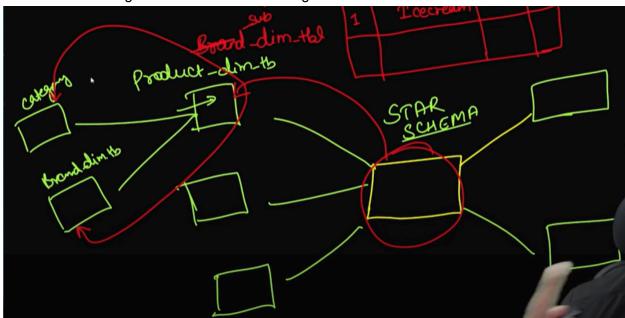11. And here the main problem data redundancy



Data Redundant

12. So we create another table for as shown pic to avoid data redundancy in a table



13. So if we are going just one level to get the data ..then it is called star schema

14. SO if we are dealing with more than 1 table to get the data..then it is snowflake schema



15. Here in the above pic..product_dim has 2 more tables attached
16. Diff bw star schema and snowflake schema

| Star Schema | Snowflake Schema | |
|---|---|---|
| It takes more Storage to keep the same info. | It takes less space than the star schema. | ⟹ Saved Storage |

Snowflake schema removes data redundancy by creating separate tables

| It takes less time in the query execution | It takes more time in the query execution. | ⟹ Wasted RAM & CPU |
|---|---|---|

In snowflake we need to use joins as there are diff tables..so it wastes RAM and CPU

| Denormalized Data | Normalized & Denormalized both Data | ⟶ Data Modeler |
|---|---|---|

in snowflake it can be both normalized and denormalized ..depends on data modeler

| | |
|---|---|
| Simple design to implement. | Complex design to implement |
| High data Redundancy | Low or No data redundancy |
| Less no. of join required. | More number of join required |


Structure of starschema and snowflake schema



17. Normalization and denormalization explained : https://g.co/gemini/share/e1bc9101c008

Primary Key and Foreign Key

1. Potential Interview Questions



Potential interview question

ⅰ) what is primary key?

ⅱ) what is foreign key?

ⅲ) Features of Primary and foreign key?

ⅳ) What is composite key?

ⅴ) what is natural / Business key?

ⅵ) what is surrogate key / fact less key?

2. Lets consider this sample table

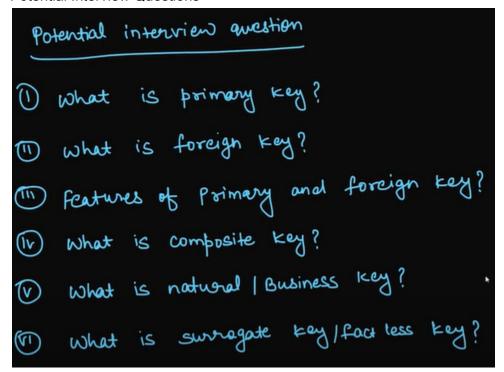| Product Table | | | | | |
|---|---|---|---|---|---|
| product_cat_id | product_category | product_sub_cat_id | product_sub_category | product_id | product_name |
| 1 | Education | 55 | Book | 101 | Manohar Pothi |
| 1 | Education | 55 | Book | 102 | Lucent GK |
| 2 | Kitchen | 201 | cookware | 103 | Stove |
| 2 | Kitchen | 201 | cookware | 104 | Microwave |
| 2 | Kitchen | 202 | Pot | 105 | Pressure Cooker |
| 2 | Kitchen | 202 | Pot | 106 | Plate |
| 3 | Grocery | 144 | cereal | 107 | Besan |
| 3 | Grocery | 144 | cereal | 108 | Maida |
| 3 | Grocery | 144 | cereal | 109 | Atta |
| 3 | Grocery | 145 | Processed | 110 | Bread |
| 3 | Grocery | 145 | Processed | 111 | cheese |

3. we have product_id…which uniquely identifies each row and there wont be any duplicates in the table it is called primary key
4. Here for product_category…there are categories which are repeating multiple times

**Category Table**

| product_cat_id | product_category |
|---|---|
| 1 | Education |
| 2 | Kitchen |
| 3 | Grocery |

5. So we'll just store this data in another table
6. And we did same for the sub categories too

**Sub Category Table**

| product_sub_cat_id(PK) | product_sub_category | category_id |
|---|---|---|
| 55 | Book | 1 |
| 201 | cookware | 2 |
| 202 | Pot | 2 |
| 144 | cereal | 3 |
| 145 | Processed | |

7. After removing this columns ..our original table will look like this…and if we need categories or sub cate…we can join the tables on id's

**Product Table**

| product_cat_id | product_sub_cat_id | product_id(PK) | product_name |
|---|---|---|---|
| 1 | 55 | 101 | Manohar Pothi |
| 1 | 55 | 102 | Lucent GK |
| 2 | 201 | 103 | Stove |
| 2 | 201 | 104 | Microwave |
| 2 | 202 | 105 | Pressure Cooker |
| 2 | 202 | 106 | Plate |
| 3 | 144 | 107 | Besan |
| 3 | 144 | 108 | Maida |
| 3 | 144 | 109 | Atta |
| 3 | 145 | 110 | Bread |
| 3 | 145 | 111 | cheese |

8. Here id are the primary columns of their table

**Category Table**

| product_cat_id(PK) | product_category |
|---|---|
| 1 | Education |
| 2 | Kitchen |
| 3 | Grocery |

**Sub Category Table**

| product_sub_cat_id(PK) | product_sub_category | category_id |
|---|---|---|
| 55 | Book | 1 |
| 201 | cookware | 2 |
| 202 | Pot | 2 |
| 144 | cereal | 3 |
| 145 | Processed | |

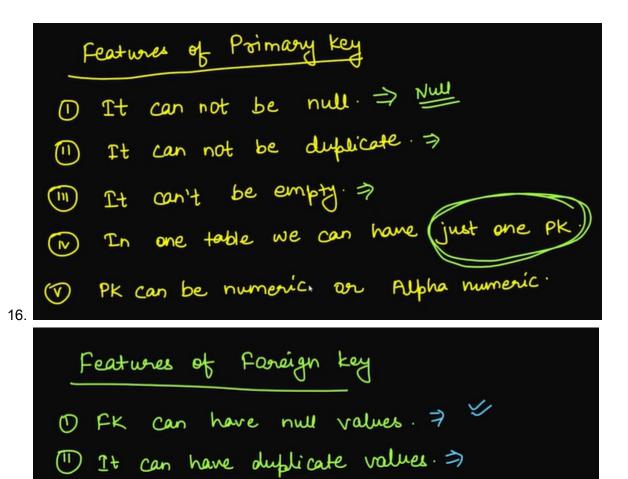9. Now in our original table…if we there are primary keys of other tables..then they are foreign keys

10. And Also the foreign key column may have the duplicates

## Category Table

| product_cat_id(PK) | product_category |
|---|---|
| 1 | Education |
| 2 | Kitchen |
| 3 | Grocery |

## Sub Category Table

| product_sub_cat_id(PK) | product_sub_category | catego |
|---|---|---|
| 55 | Book | 1 |
| 201 | cookware | 2 |
| 202 | Pot | 2 |
| 144 | cereal | 3 |
| 145 | Processed | 3 |

## Product Table

| product_sub_cat_id(FK) | product_id(PK) | product_name |
|---|---|---|
| 55 | 101 | Manohar Pothi |
| 55 | 102 | Lucent GK |
| 201 | 103 | Stove |
| 201 | 104 | Microwave |
| 202 | 105 | Pressure Cooker |
| 202 | 106 | Plate |
| 144 | 107 | Besan |
| 144 | 108 | Maida |
| 144 | 109 | Atta |
| 145 | 110 | Bread |
| 145 | 111 | cheese |

11. Now we can use joins..and get all the data required
12. Consider another table

## Product Table

| product_cat_id | product_category | product_sub_cat_id | product_sub_category | product_id | product_name |
|---|---|---|---|---|---|
| 1 | Education | 55 | Book | 101 | Manohar Pothi(10) |
| 1 | Education | 55 | Book | 101 | Manohar Pothi(20) |
| 2 | Kitchen | 202 | Pot | 105 | Pressure Cooker(2L) |
| 2 | Kitchen | 202 | Pot | 105 | Pressure Cooker(5L) |
| 2 | Kitchen | 202 | Pot | 105 | Pressure Cooker(3L) |
| 2 | Kitchen | 202 | Pot | 105 | Pressure Cooker(4L) |
| 2 | Kitchen | 202 | Pot | 106 | Plate |
| 3 | Grocery | 144 | cereal | 107 | Besan |
| 3 | Grocery | 144 | cereal | 108 | Maida |
| 3 | Grocery | 144 | cereal | 109 | Atta |
| 3 | Grocery | 145 | Processed | 110 | Bread |
| 3 | Grocery | 145 | Processed | 111 | cheese |

13. From the table above..if we clearly observe..we dont have any single column thats a primary key..
14. Here product_id is repeated multiple times…so its not a primary key
15. So here we will use both product_id and product_name as a key which identifies each row uniquely …so this type of key is called composite key

## Features of Primary key

(I) It can not be null. ⇒ Null

(II) It can not be duplicate. ⇒

(III) It can't be empty. ⇒

(IV) In one table we can have (just one PK)

(V) PK can be numeric, or Alpha numeric.

16.

## Features of Foreign key

(I) FK can have null values. ⇒ ✓

(II) It can have duplicate values. ⇒

(III) It is a (PK) of some other table. ⇒

(IV) It is used for establishing the relationship between tables.

17.

Surrogate key and Natural key

1. Lets consider a sales transaction table
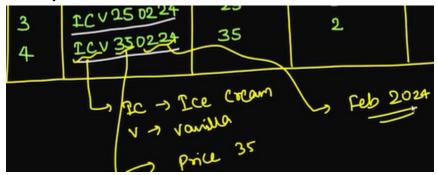
Sales-trxn-tbl

| Id | Product | Price | quantity | customer-id |
|----|---------|-------|----------|-------------|
| 1 | SUG 4025 | 40 | 2 | 1 |
| 2 | SUG 4524 | 45 | 1 | 7 |
| 3 | ICV 25 0224 | 25 | 1 | 72 |
| 4 | ICV 35 0224 | 35 | 2 | 89 |

2. Product name SUG4025 - (SUGAR's price is 40/kg and its expiry is 25(2025))..so this kind of key is called business keys



3. So likewise we can derive for every product  its price and expiry date
4. Similarly for ICV250224



5. Similarly for pharmacy table..we can design product name as



6. So here natural key can be uniquely identified or not
7. Now lets consider these table

| Id | Name | Contact | Address |
|---|---|---|---|
| MAN2123456 | Manish Kumar | 123456 | Bangalore |
| ROHA987654 | Rohan Kumar | 987654 | Delhi |
| MAN2123456 | Manish Kumar | 123456 | Bangalore |

Customer-dim-tbl

8. For surrogate key ..refer vide0

A surrogate key is an artificial identifier assigned to each record in a database table. Unlike a natural key, which is derived from the data itself (e.g., customer email address), a surrogate key has no inherent meaning and serves the sole purpose of uniquely identifying a record.

Here's why surrogate keys are useful:

- **Reliable Uniqueness:** Natural keys might not always be unique. For instance, customer names can be duplicated. Surrogate keys guarantee uniqueness, essential for database operations like record retrieval and relationship management.
- **Data Changes:** If a natural key component changes (e.g., customer email update), all linked records need modification. Surrogate keys remain unaffected by data changes within the table, simplifying maintenance.
- **Performance:** Surrogate keys are often simple numbers (e.g., auto-incrementing integer) and can be efficiently indexed for faster data retrieval compared to potentially complex natural keys.

9.

**Example: E-commerce Store**

Consider an e-commerce store with a "Customers" table.

- **Natural Key Approach:**

  - Use "Customer Email" as the primary key.
  - Issue: Emails might not be unique (duplicate accounts, typos). Joins with other tables based on email could be problematic.

- **Surrogate Key Approach:**

  - Introduce a new column "Customer ID" as the primary key (auto-generated integer).
  - "Customer Email" remains for user identification but is not the key for database operations.
  - Benefit: Guaranteed unique identifier simplifies data manipulation and relationship management with other tables.

In essence, surrogate keys provide a reliable, efficient, and independent way to identify and manage data within a database table.

## Features of Natural key

(I) Natural key may or may not be primary key.

(II) It has a business meaning associated with it.

(IX) Larger in size so takes more memory to store.

10.

## Features of surrogate key

(I) Used in fact and dimension table in Data warehousing.

(II) It is guaranteed to be unique key.

(III) Sequential numeric digit. Indexing is manged in a better way.