

we are checking data. [Transformation]

→ 1. Table.

cust-info.

checking for duplicates.

Select
cust_id,
count(*)

From bronze.cust-info

Group By cust_id.

Having count(*) > 1 or cust_id IS NULL

Gives the info of
duplicates in the
record.

Rank the duplicates.

Select

*

From (

Select

*,

ROW_NUMBER() OVER (Partition By cust_id ORDER BY cust-create_date

DESC) as flag-last

From bronze.com-cust-info

~~where cust-~~

) + where flag-last = 1

we ensure the
record without
any duplicates

The Table of strings columns.

checking for the unwanted spaces.

Remove the leading & trailing spaces.

TRIM (cust-firstname) as cust-firstname
cust.name

we will know it there only

Select

cust_firstname

From bronze.com-cust-info where cust_firstname != TRIM(cust-f

2. In gender column we see there are abbreviation for Female as (F) and male as (M)., NULL as UNKNOWN.(n/a)

So create a rule to give it's full meaning. Female & male.

```
Case when UPPER(TRIM(
  cst_gndr)) = 'F' THEN 'Female'
        WHEN UPPER(TRIM(
  cst_gndr)) = 'M' THEN 'Male'
        ELSE 'n/a'
```

3. Do the same for marital status.

S - Single

M - married

NULL - n/a.

Now we are inserting all the transformed data into.
silver schema (table) silver.com_cust_info.

Insert INTO silver.com_cust_info.

Remove Duplicates.

Ensure only one record per entity by identifying & retaining most relevant row.

Remove unwanted spaces.

Removes unnecessary spaces to ensure data consistency and uniformity across all records.

Data Normalization & Standardization

maps coded values to meaningful, user friendly descriptions

Handling Missing data

Fill in the blanks by adding a default values.

Product table.

1. Check for duplicates.
2. We need sub categories from prod-key column.

Prod-key \Rightarrow Replace (SUBSTRING (Prod-key , 1 , 5) , ' - ' , ' - ') as Cat-id,

(Position)
(How many)

Select prod-id. Check with other table where it's a foreign key
where (substring) NOT IN.
(Select distinct id from bronze.exp-pr-cat-gl-v2)

2nd Substring.

Replace (substring (Prod-key , 7 , Length ()) , ' - ' , ' - ') as
prod-key.

(dynamic until last)

where
(Select \$Sls-prod-key from bronze.com.sales-details)

3. In Prod-cost there are some null values so

use ISNULL (prod-cost) as prod-cost

4. Check for Invalid Date Orders.

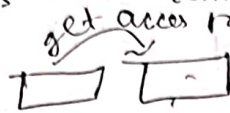
Select *

from bronze.com-prod-info

where prod-end-dt < prod-start-dt

Lead & Lag

to access the next record.



LEAD (prd-start-dt) OVER (PARTITION BY prd-key ORDER BY prd-start-dt) - 1 AS prd-end-dt-test.

Quality checks.

1. Check for Nulls or Duplicates in prd-key

```
SELECT prd-id,  
count (*)  
FROM _____  
GROUP BY prd-id.  
Having count(*) > 1 or prd-id is null.
```

2. Check for unwanted data.

- Unwanted spaces.

Sales table.

where $sis_ord_num \neq Trim(sis_ord_num)$

→ we cannot this prod-key & sis_cust_id .

check Integrity of sis_prod_key .

where sis_cust_id NOT IN (select $cust_id$ from
Silver.crm-cust-paf)

Source -

Challenging here in this table is.

we have integers as date.

we have to convert them to date.

1. check for 0s and Null

Select

NULL IF (sis_order_dt , 0) sis_order_dt

From bronze.crm-sales-details

where $sis_order_dt \leq 0$ or len

(sis_order_dt) $\neq 8$ or $sis_order_dt > 2050101$ or $-$

Check

Outliers by
validating the
boundaries of the
data range.

1 2 3 4 5 6 7 8.
20101229
if len(8) < issue.

fix it with

Case when $sis_order_dt = 0$ or $LEN(sis_order_dt) \neq 8$

Then NULL

Else cast(cast(sis_order_dt as varchar) as date)

End sis_order_dt

check for Invalid Date of Orders.

Select

from bronze.com_sales_details

where $\text{sis_order_dt} \leq \text{sales_ship_dt}$ or $\text{sis_order_dt} > \text{sis_due_dt}$

Business Rules

$\Sigma \text{ sales} = \text{Quantity} \times \text{price}$

Neg, zero, nulls, are not allowed.

Select Distinct

sis_sales,

sis_quantity,

sis_price

from bronze.com_sales_details

where $\text{sis_sales} \neq \text{sis_qunt} \times \text{sis_price}$

or sis_sales IS Null or sis_qunt IS Null or sis_price IS Null

or $\text{sis_sales} \leq 0$ or $\text{sis_qunt} \leq 0$ or $\text{sis_price} \leq 0$

Case when sis_sales IS Null or $\text{sis_sales} \leq 0$ or $\text{sis_sales} \neq$

$\text{sis_qunt} \times \text{ABS}(\text{sis_price})$

Then $\text{sis_qunt} \times \text{ABS}(\text{sis_price})$

End sis_sales

Case when $\text{sis_price} \leq 0$ or sis_price IS Null

Then $\text{sis_sales} / \text{ISNULL}(\text{sis_qunt})$

End sis_price

Rules

If sales is negative, zero or Null, derive it using Quantity & price.

If price is zero or null
calc using sales & Qunt

If price is -ve convert to +ve value.

Data Modeling

Raw data / logical Data model

3 different ways to create DM

Conceptual DM.
what are entities

✓ logical.
diff. col & relⁿ
which are PK
connection

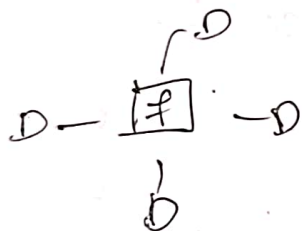
Physical DM.
dT, lines of d

Big Picture

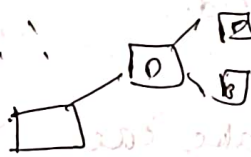
Blue Point

Implementation

Star or snowflake Schema



simple & Easy
Big Dimension.



more complex
Large Dimension.

Dimension

Descriptive info
gives context to
your data

Who what
where

Fact: up
Quantitative
if events
flow much.

How many?

Gold Layer

Views

Data integratin.

Agg.
Business logic & Rules

Star Schema
Agg. Object
Flat Tables

We have 2 gender coln when we joined the tables.
across customers.

So.

Data Integration

```
Select  
  ci.cust_gndr,  
  ca.gen  
Case when ci.cust_gndr != 'n/a' Then ci.cust_gndr  
      Else ca.gen COALESCE (ca.gen, 'n/a')  
End as gender.  
From silver.com-cust-info ci,  
Left join silver.exp-cust-az12 ca.  
on ci.cust_key = ca.cid.
```

now change name as snake case.

Surrogate key

System-generated unique identifiers assigned to each record in a table.

DOL-based generation

Query-based using (Row_no)

Row_number () Over (order by cust-id) as customer-rc

Create obj:

Create view gold.dim_customers as

Select distinct genders from gold.dim_custo

fact table

Sales transaction table.

Data look up \rightarrow joining tables with surrogate key

Building facts

inc. dim's surrogate key insted of IDs so easily connect facts with dimensions

Data schema model

May.

Many

1-many

1-cust who haven't placed any order yet

2. none. one —
3. none multiple.

Data catalog

To understand the data better.