# INFO 7370 Designing Data Architecture Business Intelligence

# Assignment 4

Srishti Ashok Mishra

NU ID : 001305178
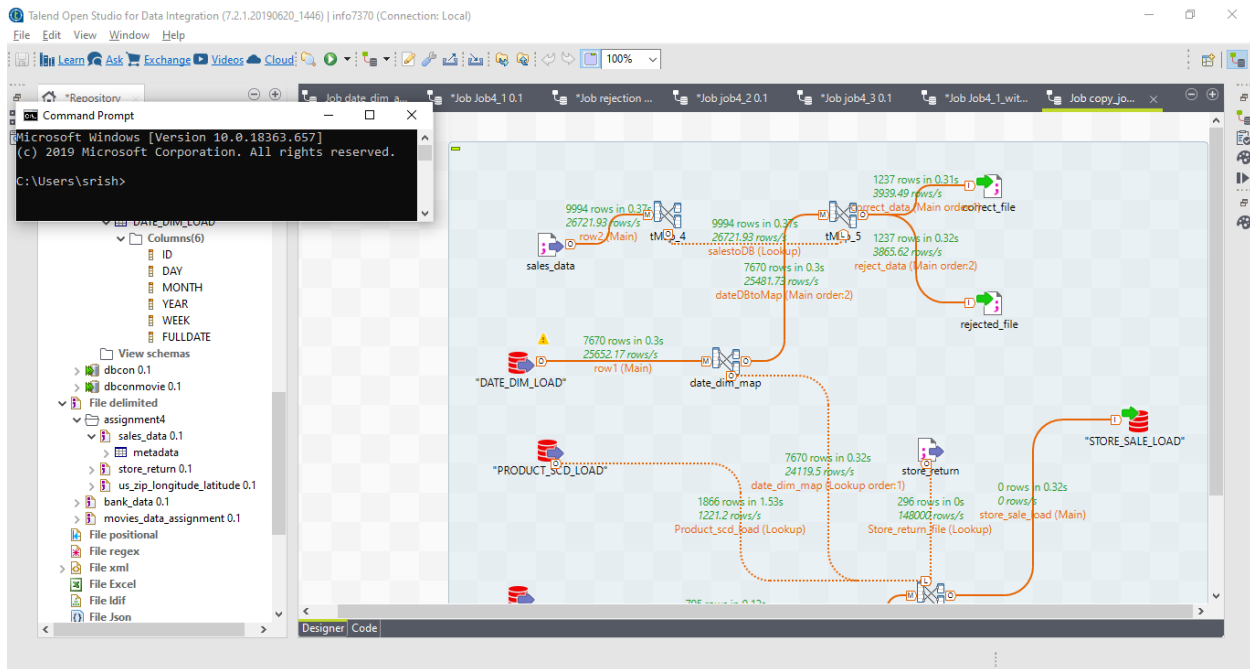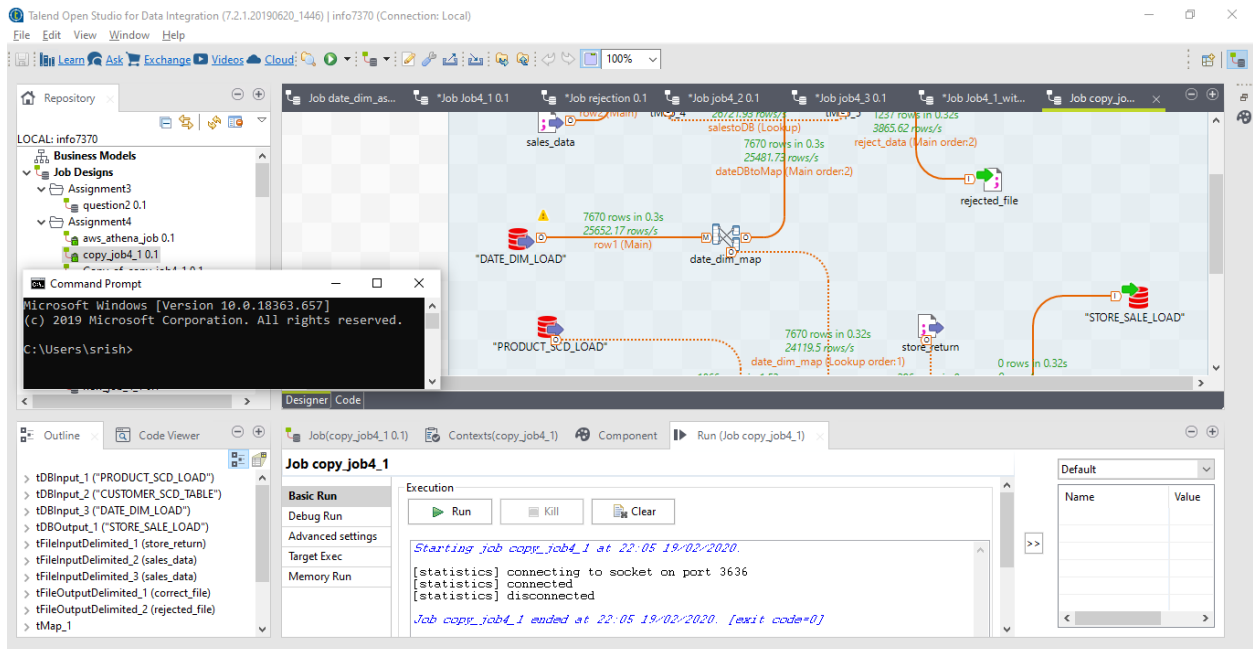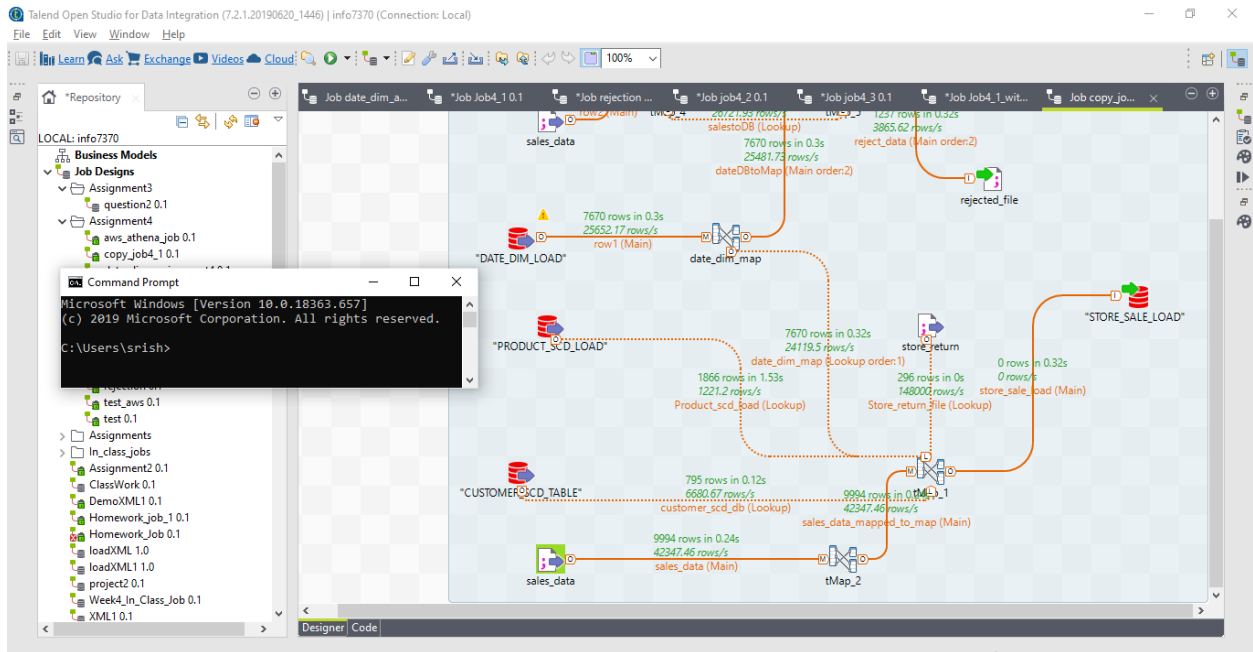
# JOB:

## SALES LOAD:

For this Job:

*Customer_dim_id, product_dim_id, order_dt_id* are based on the other jobs and are linked by surrogate key.

*Ship_date, ship_mode, qty_sold, unit_price, discount_percentage, order_id* where obtain and/ or calculated using store_sales file in tmap

**Data in DB:**

## MAPPING:



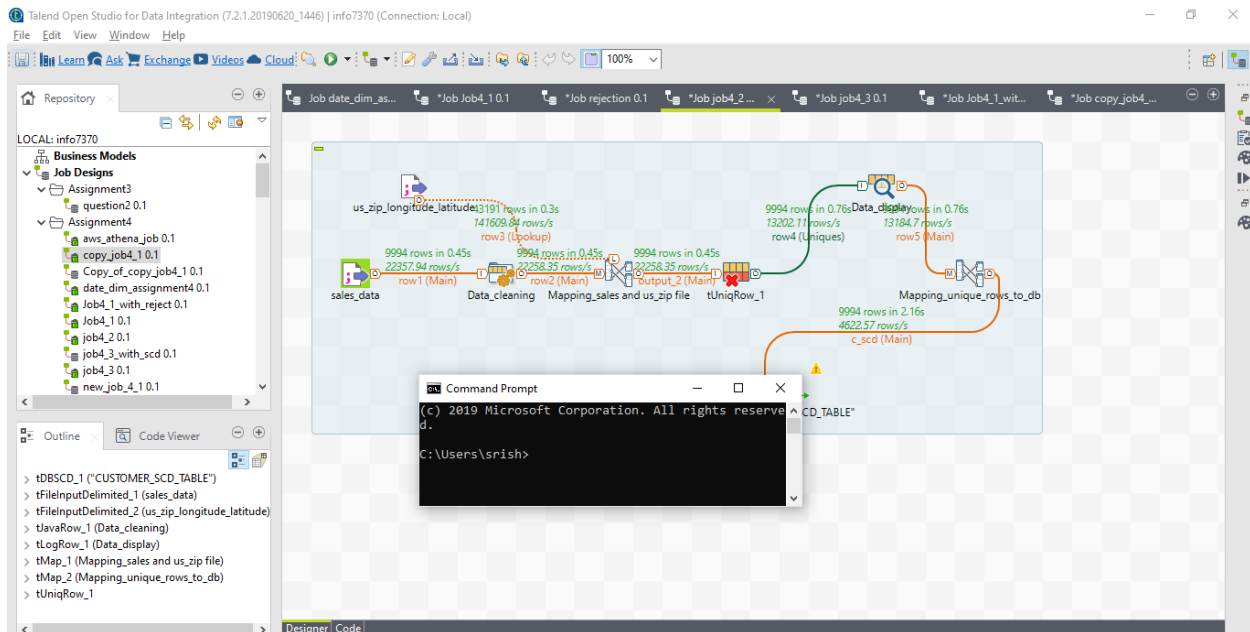## CUSTOMER_LOAD:

For this Job:

Data for *latitude, Longitude, city, state* was obtained by from file us-zip-code-latitude-and-longitude.csv lookup based on zipcode

Data for *First name* and *last name* was obtained from splitting up full name in *tmap* component

Data for *segment* and *region* was directly from *store_sales* file by using it as a delimited component.

## Data in DB:



## MAPPING:

## PRODUCT_LOAD:

For this Job:

All the data was directly from *store_sales* file by using it as a delimited component.



## Data in DB:

# MAPPING:



# SCD EDITOR:

Here the implementation is based on SCD-2.

Product_load

The source delimited file contains various details including *product_id*, *category*, *sub_category*, *product_name*. An *id* column helps ensuring the uniqueness of the data.

Here source Id is  *product_id*

Scd type 2 field is *category*

There is a new record inserted in the dimensional table with a separate key each time the category changes.

The Table will retain both old and new value and the latest value will be shown active in the *scd_is_active* column.

Here the SCD type 2 columns are:

*Scd_start_date* to store the start date

*SCD_end_date* to store the end date

*SCD_is_active* stores which of the column is active

*SCD_version* to store the version of the scd

Similarly for Customer_load

## REJECTED DATA MAPPING:



Here the file is rejected based on the order_id and date_dimension table.

*dateDBtoMap.Date_to_String.equals(  salestoDB.Order_Date )*

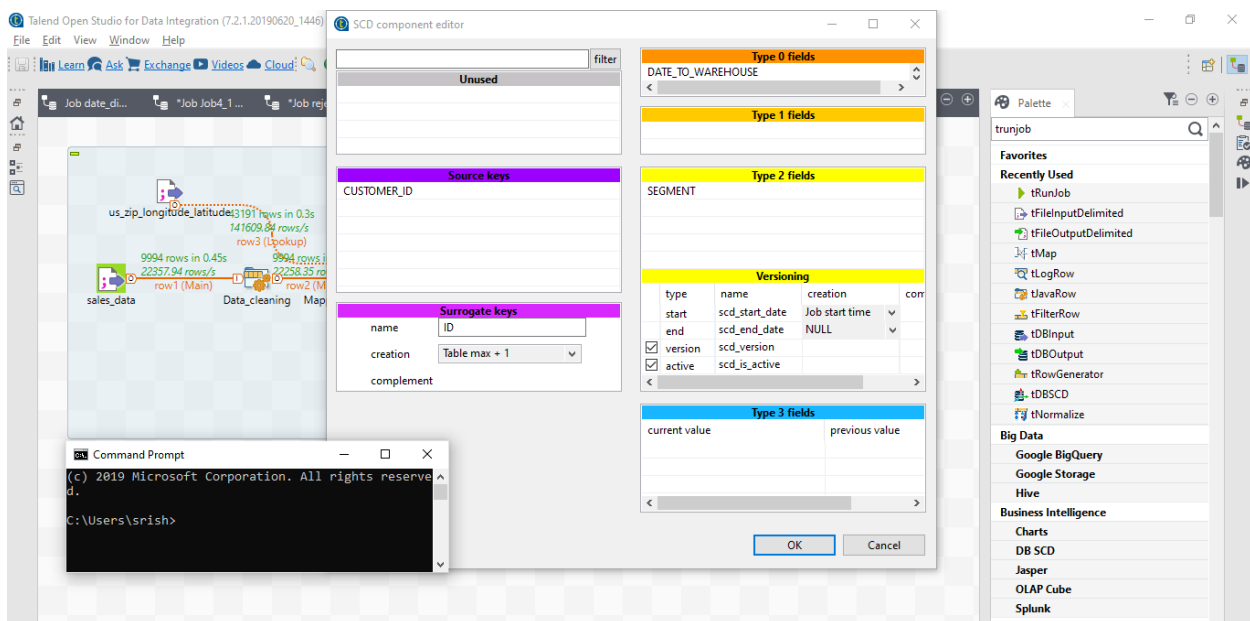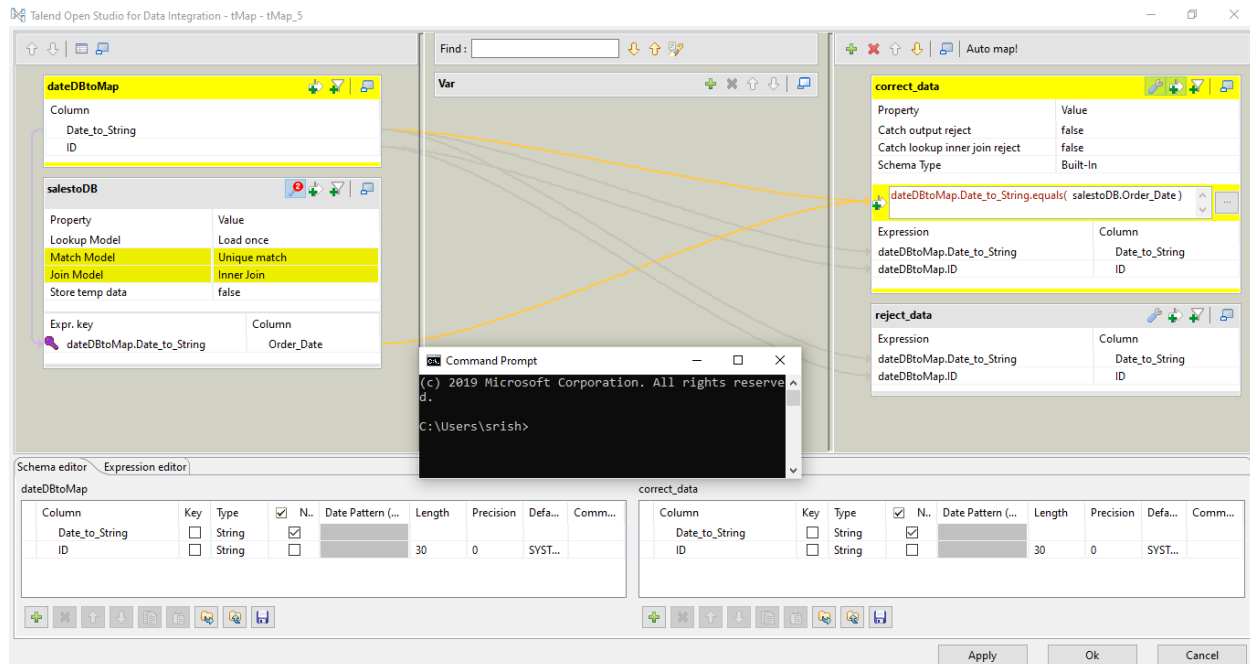if the data falls into the formula, the value will be a part of the correct data, else goes into the reject data csv.

## Master Job:

Master Job consist of Customers_load followed by products_load job. Once these both are successfully completed then it invoke sales_load

It uses tRunJob component to merge all the jobs into a master job