

# Documentation for Capstone Project on Data Warehouse Design and ETL Implementation Using Informatica

## US CHRONIC DISEASE INDICATOR DATASET

~ KRUTHIK BALLARI ~

### PROBLEM STATEMENT

This is a scalable data warehouse design project aiming at the solution to the problem of complex and voluminous dataset analysis and the derivation of actionable insights from those datasets. Due to an increased demand in businesses for efficient mechanisms in data storage and processing, this capstone project is about building a reliable ETL pipeline using Informatica on transforming raw data into structured information housed within a well-designed data warehouse. The solution will be identifying relevant dimensions and facts, modeling the data warehouse schema, and implementing ETL workflows that ensure data accuracy, consistency, and scalability. Finally, SQL-based queries will be used to extract meaningful insights to support business decision-making processes.

### DATASET DESCRIPTION

The US Chronic Disease Indicator dataset is a comprehensive data source that provides information on the prevalence and incidence of chronic diseases within the United States. This dataset typically includes key health indicators related to conditions such as diabetes, cardiovascular diseases, obesity, cancer, and respiratory diseases. The data is collected and reported by health agencies and public health departments to monitor trends, support policy-making, and improve public health interventions.

**Total Columns: 34**

**Total Data Rows: 11,85,676**

**Source Data Type: CSV (Comma Delimited)**

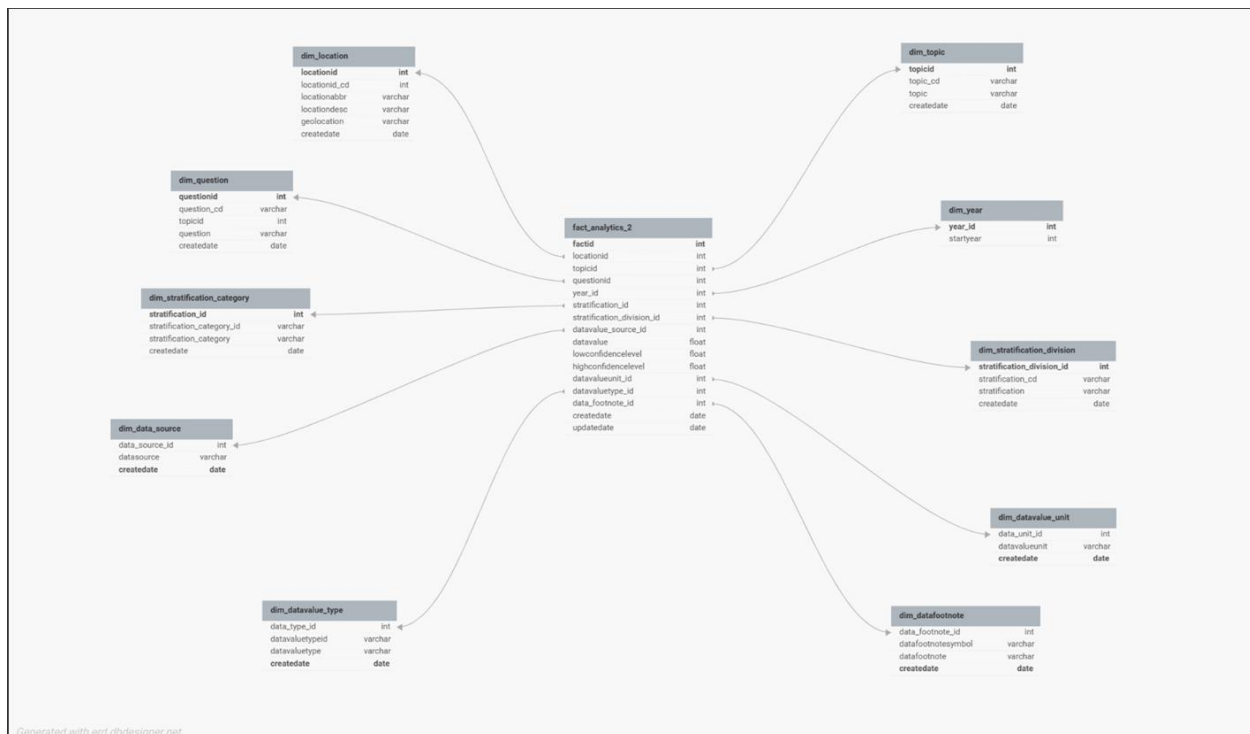
**Dataset Name: US Chronic Disease Indicator Dataset**

1	Column Names	18	StratificationCategory1
2	YearStart	19	Stratification1
3	YearEnd	20	StratificationCategory2
4	LocationAbbr	21	Stratification2
5	LocationDesc	22	StratificationCategory3
6	LocationDesc	23	Stratification3
7	Topic	24	GeoLocation
8	Question	25	ResponseID
9	Response	26	LocationID
10	DataValueUnit	27	TopicID
11	DataValueType	28	QuestionID
12	DataValue	29	DataValueTypeID
13	DataValueAlt	30	StratificationCategoryID1
14	DataValueFootnoteSymbol	31	StratificationID1
15	DataValueFootnote	32	StratificationCategoryID2
16	LowConfidenceLimit	33	StratificationID2
17	HighConfidenceLimit	34	StratificationCategoryID3
		35	StratificationID3

## Staging Table Design:

stg_us_chronic	
yearstart	int
yearend	int
locationabbr	varchar
locationdesc	varchar
datasource	varchar
topic	varchar
question	varchar
response	float
datavalueunit	varchar
datavaluetype	varchar
datavalue	varchar
datavaluealt	float
datavaluefootnotesymbol	varchar
datavaluefootnote	varchar
lowconfidencelimit	float
highconfidencelimit	float
stratificationcategory1	varchar
stratification1	varchar
stratificationcategory2	varchar
stratification2	varchar
stratificationcategory3	varchar
stratification3	varchar
geolocation	varchar
responseid	float
locationid	int
topicid	varchar
questionid	varchar
datavaluetypeid	varchar
stratificationcategoryid1	varchar
stratificationid1	varchar
stratificationcategoryid2	varchar
stratificationid2	varchar
stratificationcategoryid3	varchar
stratificationid3	varchar

## DIMENSION TABLE AND FACT TABLE DB DESIGN



Total Number of Dimension tables are 10 and the Fact tables is 1

## **Steps for Informatica Project on US Chronic Disease Indicator Dataset**

### **1. Data Understanding and Modeling**

#### **1. Dataset Analysis**

- Analyzed the dataset to identify dimensions (e.g., Location, Time, Disease) and facts (e.g., Counts, Prevalence Rate, Age-Adjusted Rate).
- Examined data attributes for primary keys, foreign keys, and granularity levels.

#### **2. Schema Design**

- **Designed a staging table in MS SQL to hold raw data after extraction.**
- **Developed a data model comprising:**
  - Dimension Tables: Location, Time, Disease Type, Demographics.
  - Fact Table: Chronic Disease Metrics (with keys to dimensions and measurable facts).

### **2. ETL Implementation**

#### **1. Staging Table Population**

- Loaded the CSV dataset into the staging table using Informatica mappings.
- Validated data integrity and handled missing values during extraction and loading.

#### **2. Mapping for Dimension Tables**

- **Created mappings for dimension tables using:**
  - Insert Only: For static reference data like locations.
  - SCD Type 1: For attributes where changes overwrite old data, e.g., updated disease descriptions.
  - SCD Type 2: For attributes where historical changes need tracking, e.g., demographic attributes with time validity.

#### **3. Fact Table Loading**

- Defined and implemented mappings to load aggregated metrics into the fact table.
- Ensured correct linkage with dimension tables using surrogate keys.

#### **4. Workflow Implementation**

- **Designed workflows to automate the ETL process for all mappings:**
  - Used Event Wait Task to synchronize data loading steps.

- Applied Command Task for external commands like generating success logs.
- Parameterized database connections and table names for flexibility.

### 3. Data Analysis

#### 1. SQL Query Development

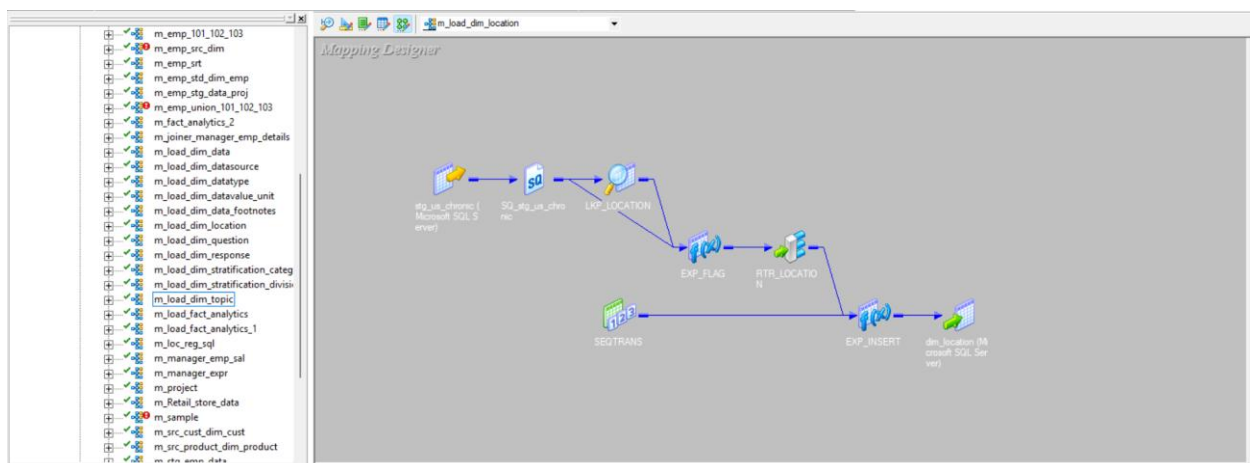
- **Wrote SQL queries on the dimension and fact tables to extract actionable insights.**
  - Example 1: Identify the top locations with the highest prevalence of a specific chronic disease.
  - Example 2: Analyze trends in age-adjusted rates over time for specific demographics.

#### 2. Data Validation

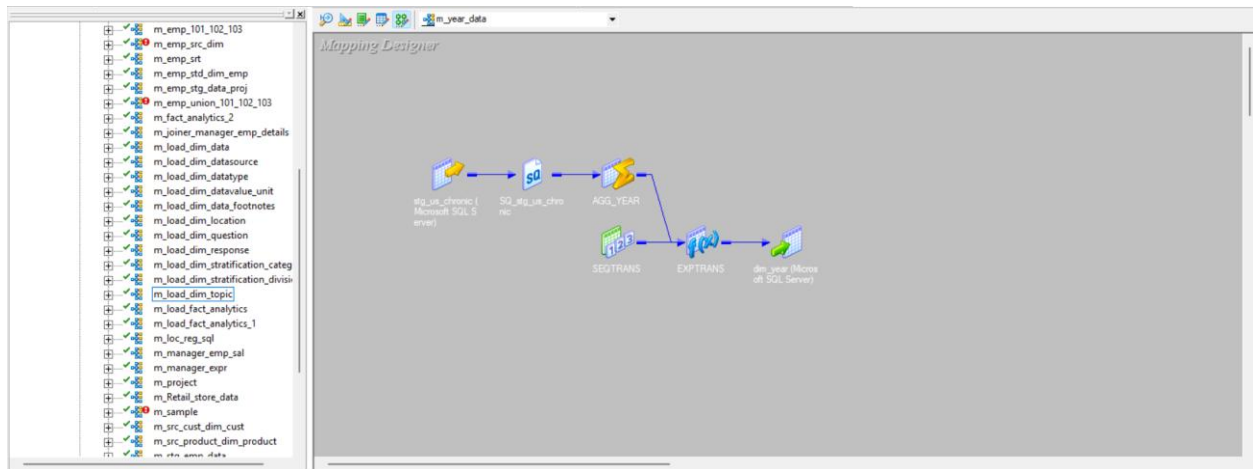
- Validated data accuracy by cross-checking derived metrics against the staging table and input dataset.

### DIMENSION TABLE MAPPINGS:

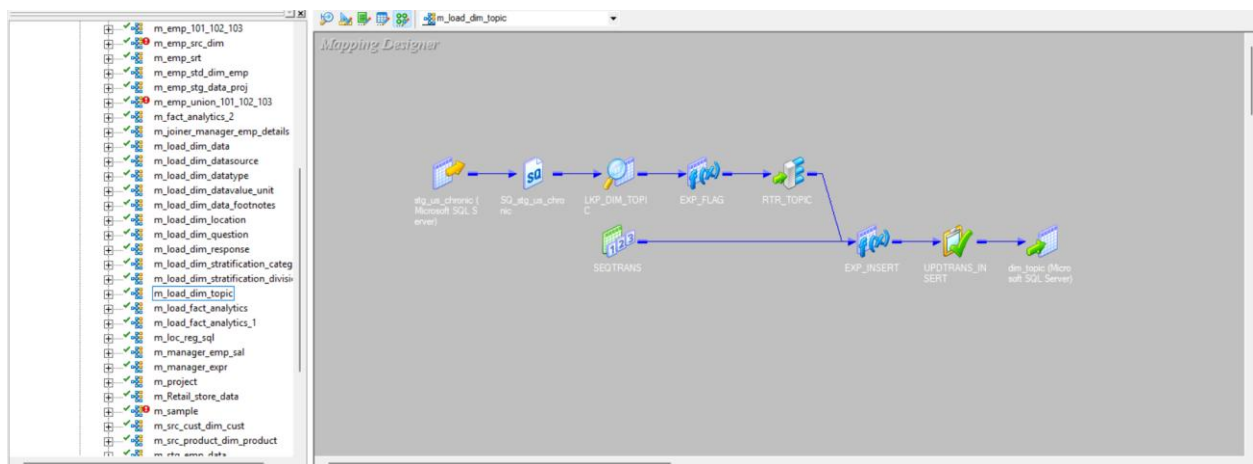
#### DIM\_LOCATION:



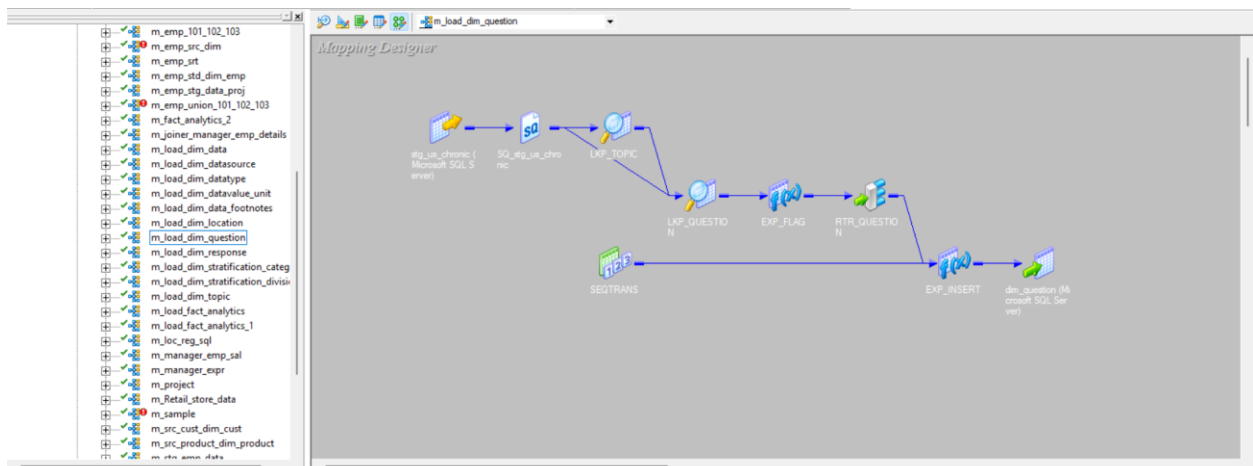
## DIM\_YEAR



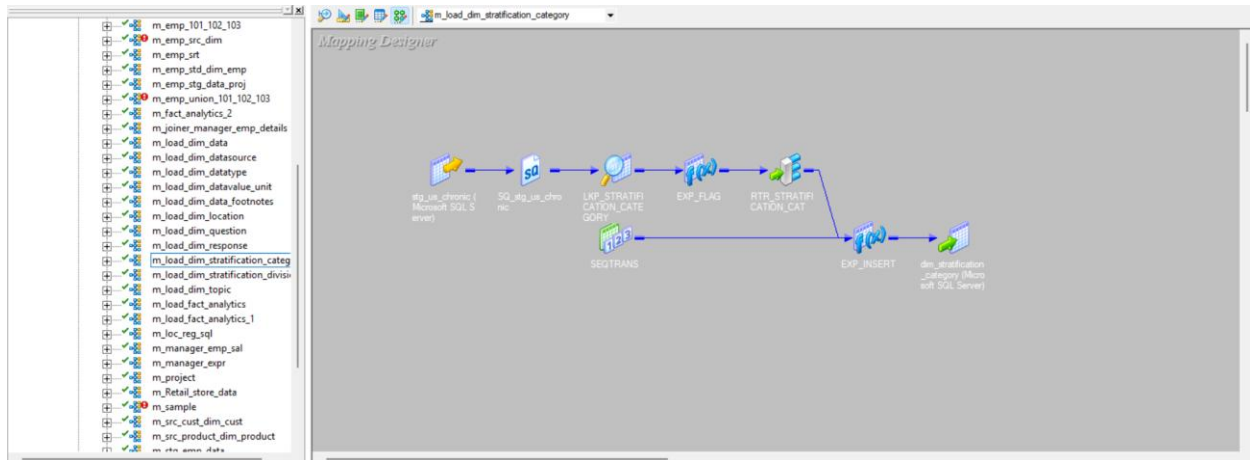
## DIM\_TOPIC



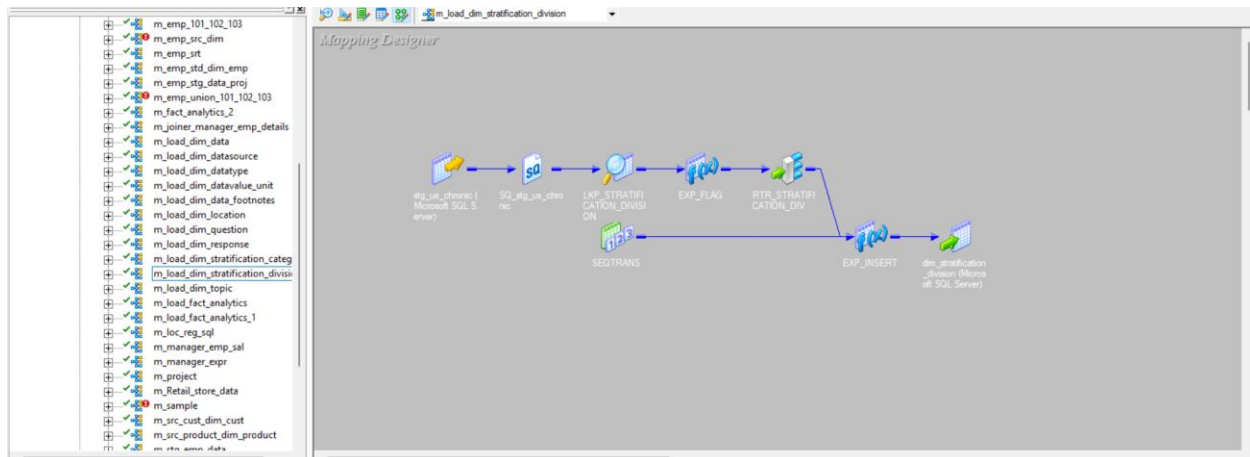
## DIM\_QUESTION



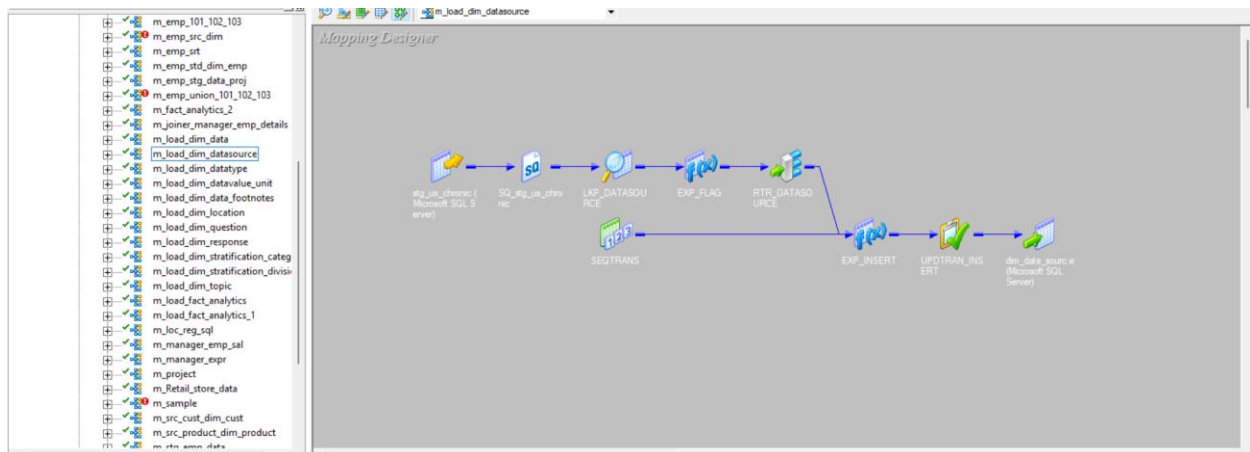
## DIM\_STATIFICATION\_CATEGORY



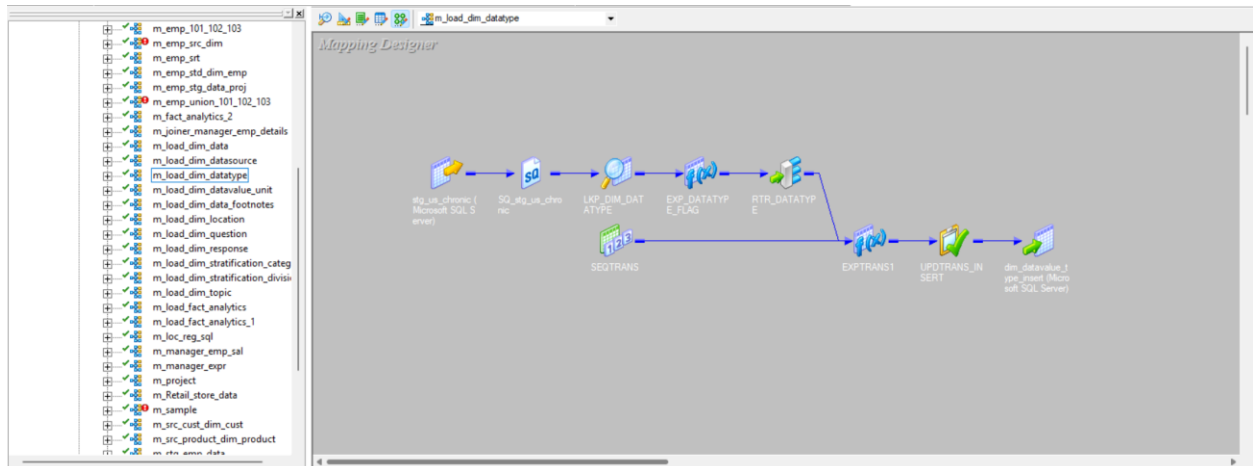
## DIM\_STRATIFICATION\_DIVISION



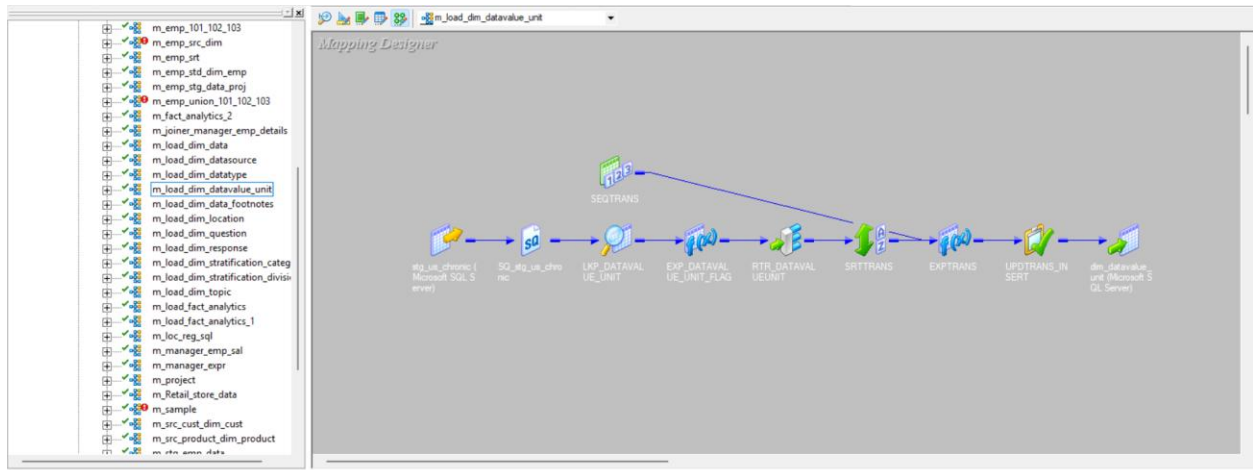
## DIM\_DATSOURCE



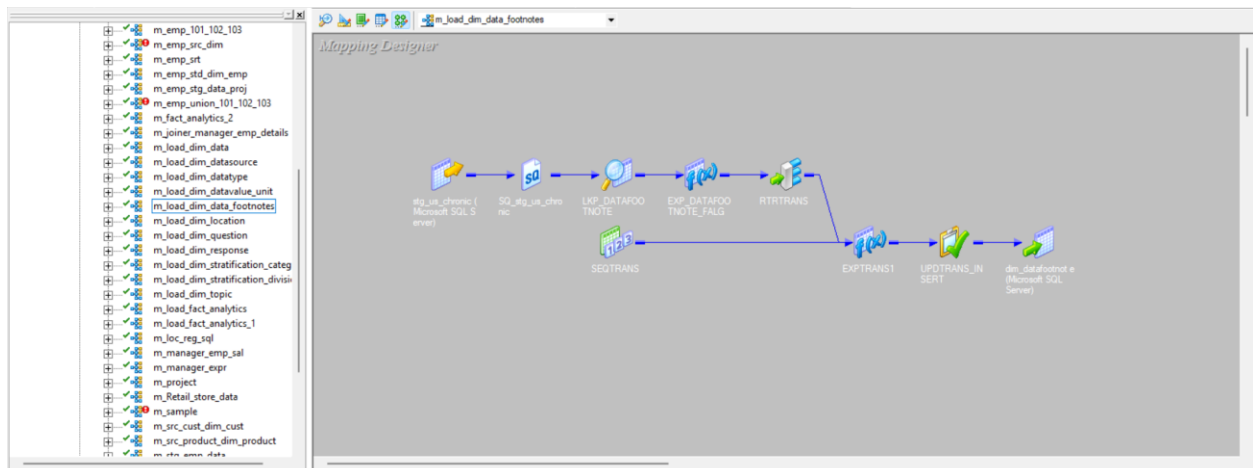
## DIM\_DATATYPE



## DIM\_DATAVALUE\_UNIT



## DIM\_DATA\_FOOTNOTES



## FACT TABLE





**Project Center**

Repository Edit View Tools Task Filters Help

1 Hour

Name	Duration	Status	Dec 8, 2024 12:00am	12:00am	1:00am	2:00am	3:00am	4:00am	5:00am	6:00am	7:00am	8:00am	9:00am	10:00am	11:00am	12:00pm	1:00pm	2:00pm	3:00pm	4:00pm	5:00pm	6:00pm	7:00pm
NFA_DEV_REP	03:13:22	Connected																					
NFA_DEV_INT																							
DEV																							
us_CHRONIC_SRC_ST	00:07:16	Succeeded																					
s_csv_stg_data	00:01:12	Succeeded																					
s_csv_file_check	00:00:00	Succeeded																					
c_stg_success	00:00:03	Succeeded																					
e_stg_success	00:00:02	Succeeded																					
s_load_dim_loc	00:00:04	Succeeded																					
s_load_dim_top	00:00:03	Succeeded																					
s_load_dim_year	00:00:03	Succeeded																					
s_load_dim_question	00:00:05	Succeeded																					
s_load_dim_strat_cat	00:00:04	Succeeded																					
s_load_dim_strat_div	00:00:03	Succeeded																					
s_load_dim_data_source	00:00:03	Succeeded																					
s_load_dim_data_value	00:00:04	Succeeded																					
s_load_dim_datatype	00:00:03	Succeeded																					
s_load_dim_data_isochr	00:00:02	Succeeded																					
c_dim_tables_success	00:00:02	Succeeded																					
e_dim_tables_success	00:00:02	Succeeded																					
B Load fact table	00:01:43	Succeeded																					

Gantt Chart Task View

---

### Task Details

**Instance Name**  
Task Type  
Integration Service Name  
Node(s)  
Start Time  
End Time  
Recovery Time(s)

**Attribute Value**  
**s\_load\_fact\_table**  
**Session**  
NFA\_DEV\_INT  
node01  
12/8/2024 4:05:15 PM  
12/8/2024 4:06:50 PM

---

### Source/Target Statistics

Transformation Name	Node	Applied Rows	Affected Rows	Rejected Rows	Throughput (Rows/Sec)	Throughput (Bytes/Sec)	Bytes	Last Error Code
fact_analytics_2_in	node01	1185676	0	0	17437	1412397	96939756	0
fact_analytics_2_uo	node01	0	0	0	0	0	0	0
SO_stg_us_chronic	node01	1185676	1185676	0	39523	102285524	3068529488	0

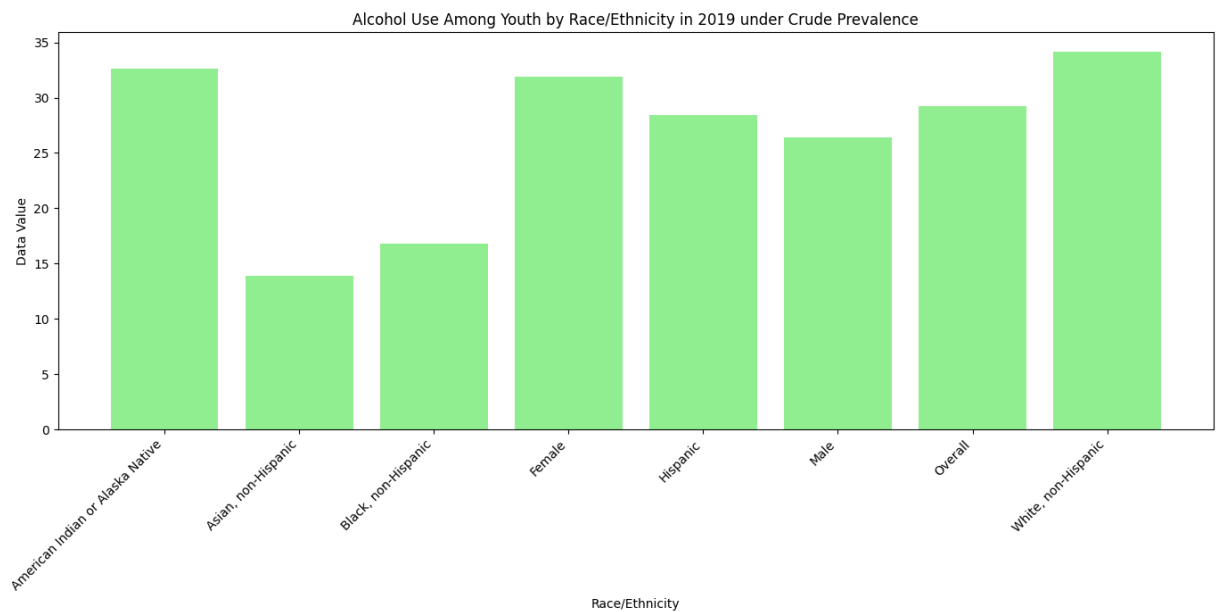
### Partition Details

Performance

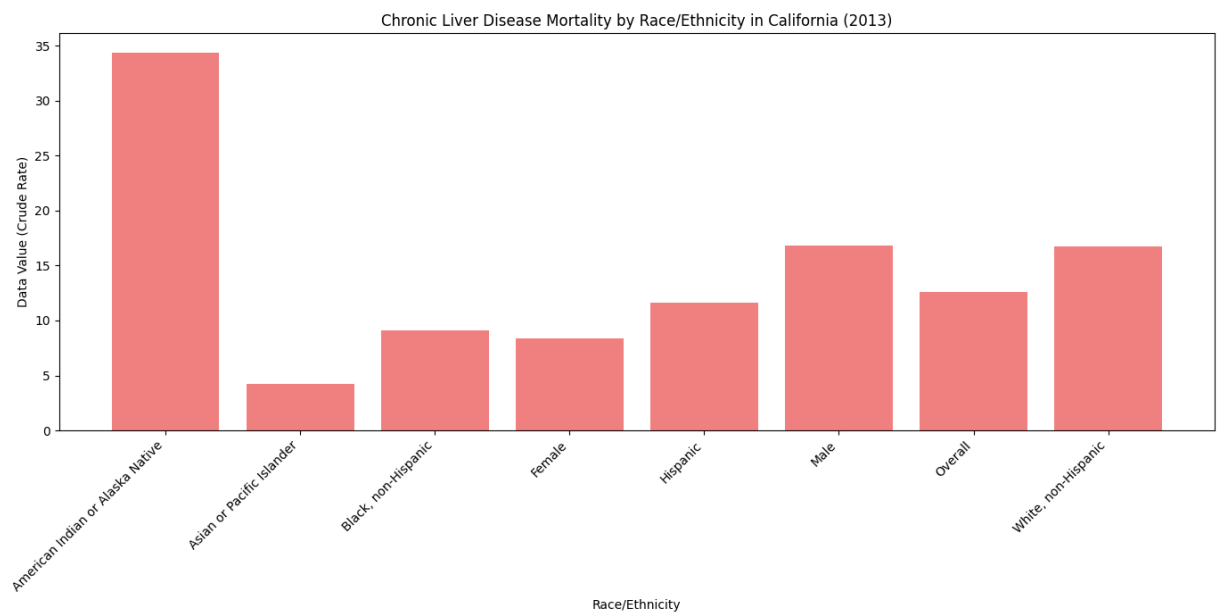
Output Window | Notifications

```
(NF_DEV_INT 12/8/2024 4:01:39 PM) Task Update s_load_dim_top (Running) Start time: 12/8/2024 4:01:39 PM End time: 12/8/2024 4:01:42 PM Task Update s_load_dim_top (Succeeded) Start time: 12/8/2024 4:01:42 PM End time: 12/8/2024 4:02:02 PM Task Update s_load_dim_year (Waiting) Start time: 12/8/2024 4:02:02 PM End time: 12/8/2024 4:02:02 PM Task Update s_load_dim_year (Running) Start time: 12/8/2024 4:02:02 PM End time: 12/8/2024 4:02:05 PM Task Update s_load_dim_year (Succeeded) Start time: 12/8/2024 4:02:05 PM End time: 12/8/2024 4:02:25 PM Task Update s_load_dim_question (Waiting) Start time: 12/8/2024 4:02:25 PM End time: 12/8/2024 4:02:25 PM Task Update s_load_dim_question (Running) Start time: 12/8/2024 4:02:25 PM End time: 12/8/2024 4:02:25 PM Task Update s_load_dim_strat_cat (Waiting) Start time: 12/8/2024 4:02:25 PM End time: 12/8/2024 4:02:42 PM Task Update s_load_dim_strat_cat (Running) Start time: 12/8/2024 4:02:42 PM End time: 12/8/2024 4:02:42 PM Task Update s_load_dim_strat_div (Waiting) Start time: 12/8/2024 4:02:42 PM End time: 12/8/20
```

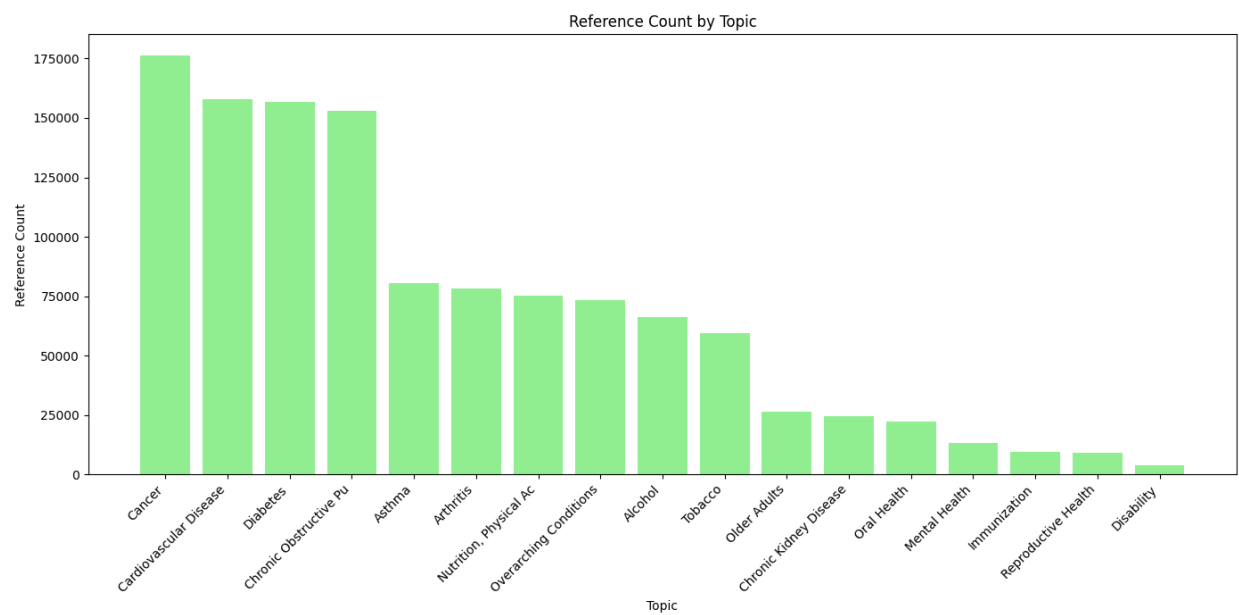
### Alcohol Use Among Youth by Race Ethnicity in 2019 under Crude Prevalence



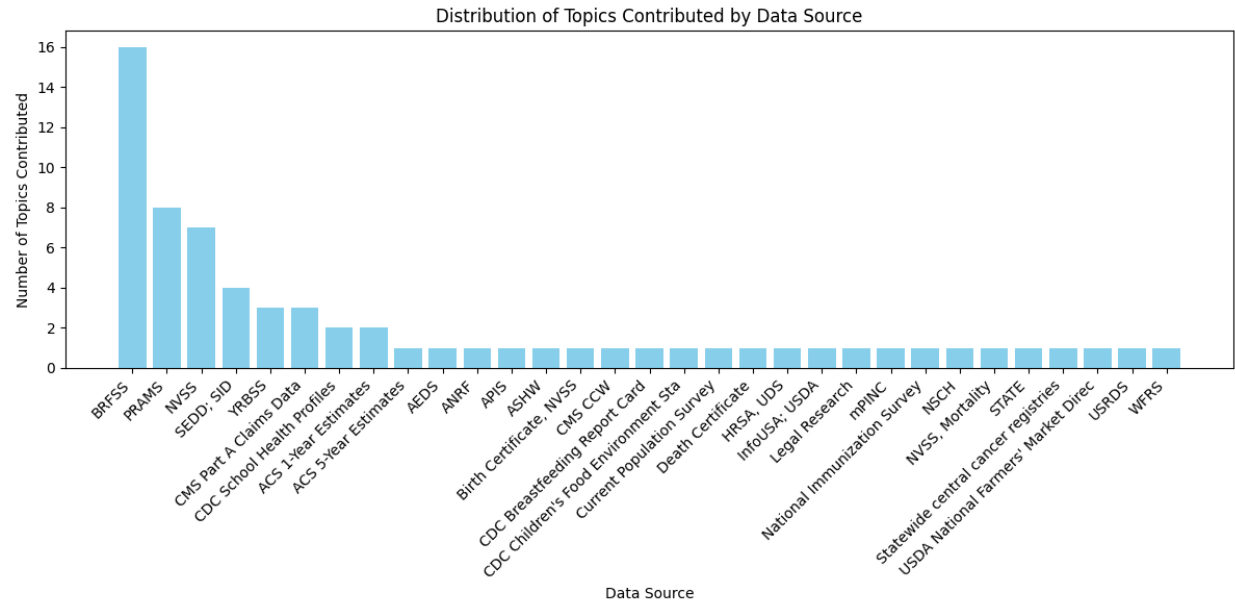
## Chronic liver disease mortality in California in 2013



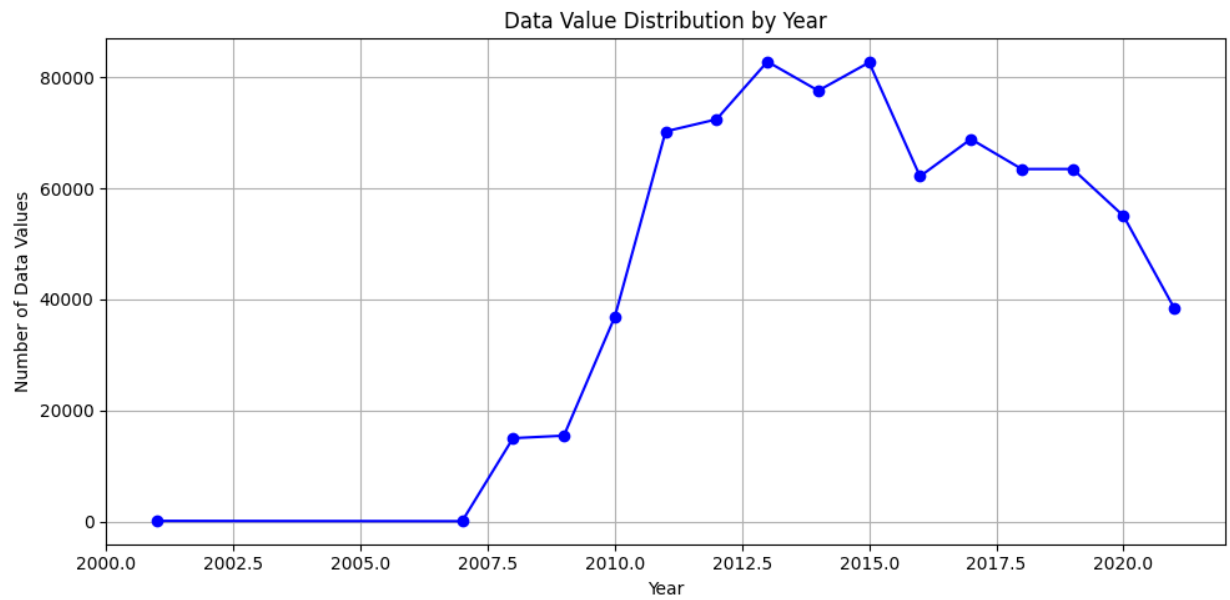
## Reference Count



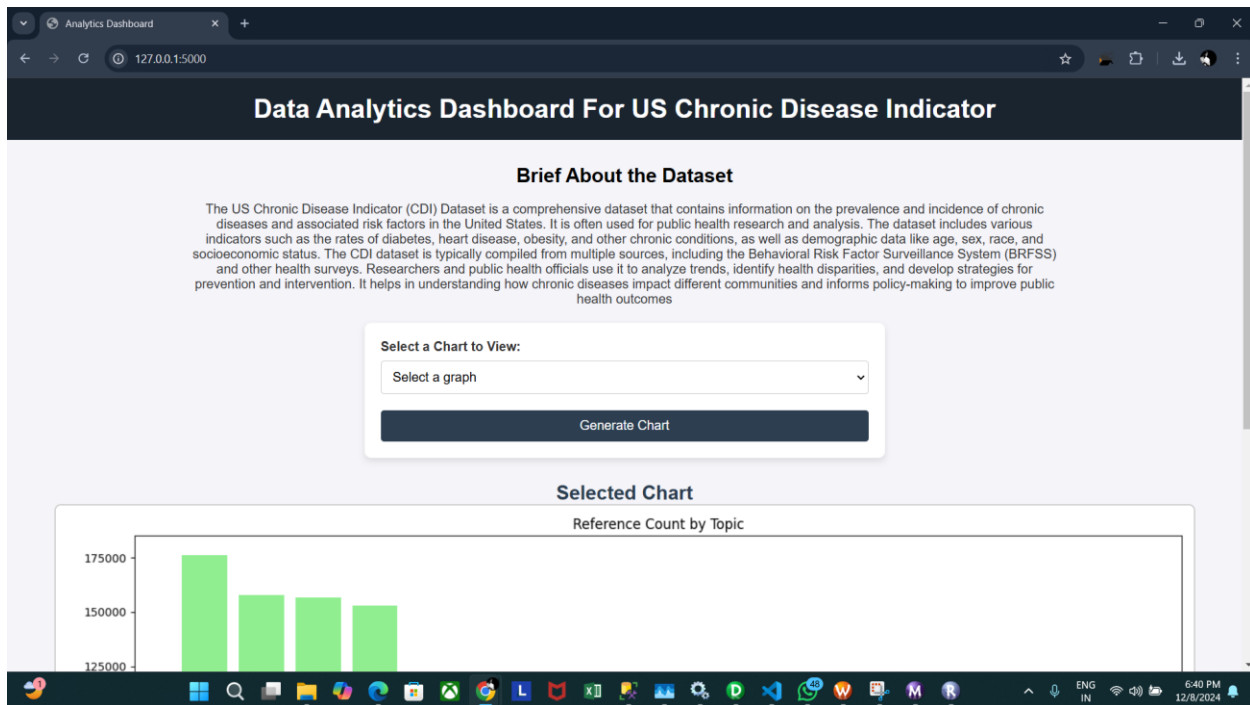
## Topics Distribution



## Year Distribution



## WEBSITE INTERFACE DEVELOPED FOR DATA CHARTS VISUALIZATION



## CONCLUSION

The successful completion of the project on Data Warehouse Design and ETL Implementation Using Informatica demonstrates the practical application of data warehousing concepts and ETL processes in addressing real-world data challenges. By transforming the raw dataset of US Chronic Disease Indicators into a structured data warehouse, actionable insights were derived to support decision-making and analysis.

The project achieved its objectives through:

1. **Comprehensive Data Modeling:** Designing dimension and fact tables with a clear understanding of the dataset's structure and requirements.
2. **Efficient ETL Implementation:** Leveraging Informatica to perform seamless data extraction, transformation, and loading, while adhering to best practices like parameterization and workflow dependencies.
3. **Accurate Data Analysis:** Writing and executing SQL queries to validate the data and extract meaningful insights, highlighting the significance of the data warehouse in real-time analytics.