

Midterm Team Project

Designing Advanced Data Architectures for Business Intelligence

Group 13
Christina Daniel-
Nitish Belagali- 002645514
Shivpriya Mane- 002850626
Shruti Dhamdhere- 002812094

PART 1:

Ydata Profiling:

Food Inspection Report for Chicago:

The screenshot shows a Jupyter Notebook interface with the title "jupyter Untitled". The code cell In [4] contains the command `from ydata_profiling import ProfileReport` and imports pandas. The code cell In [9] reads a CSV file "Chicago.tsv" into a DataFrame df and prints it. The output shows the first 10 rows of the Chicago food inspection dataset.

```
In [4]: from ydata_profiling import ProfileReport
import pandas as pd

In [9]: df = pd.read_csv(r'C:\Users\aadit\Downloads\Chicago.tsv", sep='\t')

# printing data
print(df)
```

	Inspection ID	DBA Name	AKA Name	License #	Facility Type
0	2589772	Thomas, Velma ECC	Thomas, Velma ECC	26891.0	School
1	2589775	TAQUERIA ATOTONILCO #1	TAQUERIA ATOTONILCO #1	33733.0	Restaurant
2	2589762	MANOLOS TAMALES 4	MANOLOS TAMALES 4	2745155.0	Restaurant
3	2589743	LA SERRE/BAR LA RUE	LA SERRE/BAR LA RUE	2939551.0	Restaurant
4	2589753	PLAMONDON ELEMENTARY	PLAMONDON ELEMENTARY	24981.0	School
...
267861	67757	DUNKIN DONUTS/BASKIN-ROBBINS	DUNKIN DONUTS/BASKIN-ROBBINS	1380279.0	Restaurant
267862	70269	mr.daniel's	mr.daniel's	1899292.0	Restaurant
267863	104236	TEMPO CAFE	TEMPO CAFE	80916.0	Restaurant
267864	52234	Cafe 608	Cafe 608	2013328.0	Restaurant
267865	67738	MICHAEL'S ON MAIN CAFE	MICHAEL'S ON MAIN CAFE	2008948.0	Restaurant

The screenshot shows a Jupyter Notebook interface with the title "jupyter Untitled". The code cell In [10] generates a ProfileReport for the DataFrame df with the title "Food Inspection Report". The code cell In [11] attempts to display the report in an iframe, but due to kernel issues, it shows error messages about widget state not being found. Below the notebook, a preview of the "Food Inspection Report" is shown with tabs for Overview, Variables, Interactions, Missing values, and Sample.

```
In [10]: profile = ProfileReport(df, title = "Food Inspection Report")

In [11]: profile.to_notebook_iframe()
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Food Inspection Report

Overview Variables Interactions Missing values Sample

jupyter Untitled Last Checkpoint: Last Thursday at 7:44 PM (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 Logout

Overview

Overview Alerts 5 Reproduction

Dataset statistics

Number of variables	17
Number of observations	267866
Missing cells	84184
Missing cells (%)	1.8%
Duplicate rows	0
Duplicate rows (%)	0.0%
Total size in memory	34.7 MiB
Average record size in memory	136.0 B

Variable types

Numeric	5
Text	8
Categorical	3
DateTime	1

Variables

jupyter Untitled Last Checkpoint: Last Thursday at 7:44 PM (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 Logout

Variables

Inspection ID ▾

Inspection ID
Real number (\mathbb{R})

UNIQUE

Distinct	267866
Distinct (%)	100.0%
Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%
Mean	1730790

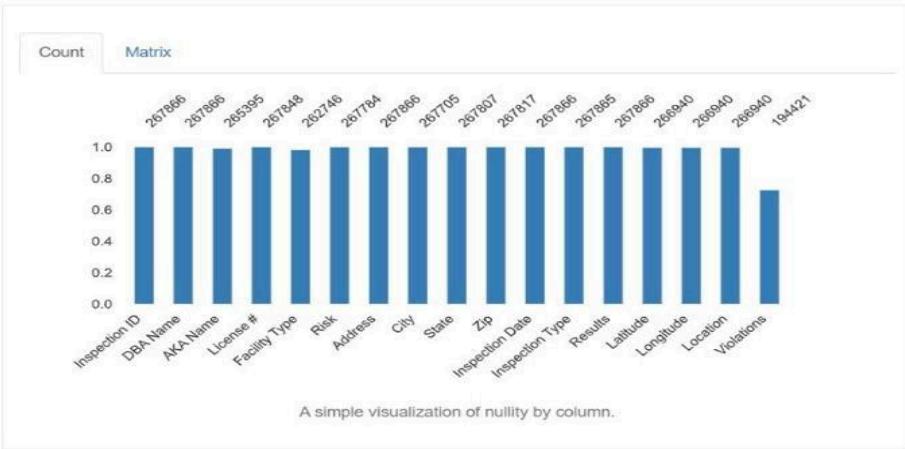
Minimum	44247
Maximum	2589775
Zeros	0
Zeros (%)	0.0%
Negative	0
Negative (%)	0.0%
Memory size	2.0 MiB



More details



Missing values



Sample

In this report, we undertook a thorough data profiling process, analyzing missing values, distinct values, duplicate rows, and inappropriate data formats, alongside a meticulous review of column details. This comprehensive examination not only ensures data quality by identifying and rectifying issues like missing values, duplicate rows, and improper data formats but also enhances the reliability and accuracy of our analysis. Moreover, through this profiling, we pinpointed relevant attributes and measures crucial for our dimensional modeling efforts. These insights will serve as the foundation for the design and implementation phases, facilitating a robust and effective dimensional model.

YData Profiling to get Restaurant report:

jupyter y-prof (1) 1 Last Checkpoint: 21 minutes ago (autosaved)

In [11]: profile.to_notebook_iframe()

Summarize dataset: 100% 128/128 [01:07<00:00, 2.85s/it, Completed]

Generate report structure: 100% 1/1 [01:54<00:00, 114.13s/it]

Render HTML: 100% 1/1 [00:11<00:00, 11.15s/it]

Restaurant report Overview Variables Interactions Correlations Missing values Sample Duplicate rows

Overview

Overview Alerts 115 Reproduction

Dataset statistics	
Number of variables	114
Number of observations	78400
Missing cells	6454357
Missing cells (%)	72.2%
Duplicate rows	42

Variable types	
Text	81
Categorical	29
DateTime	2
Numeric	2

jupyter y-prof (1) 1 Last Checkpoint: 23 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3

Overview

Overview Alerts 115 Reproduction

Dataset statistics	
Number of variables	114
Number of observations	78400
Missing cells	6454357
Missing cells (%)	72.2%
Duplicate rows	42
Duplicate rows (%)	0.1%
Total size in memory	68.2 MiB
Average record size in memory	912.0 B

Variables

Select Columns

jupyter y-prof (1) 1 Last Checkpoint: 21 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3

Restaurant report Overview Variables Interactions Correlations Missing values Sample Duplicate rows

Variables

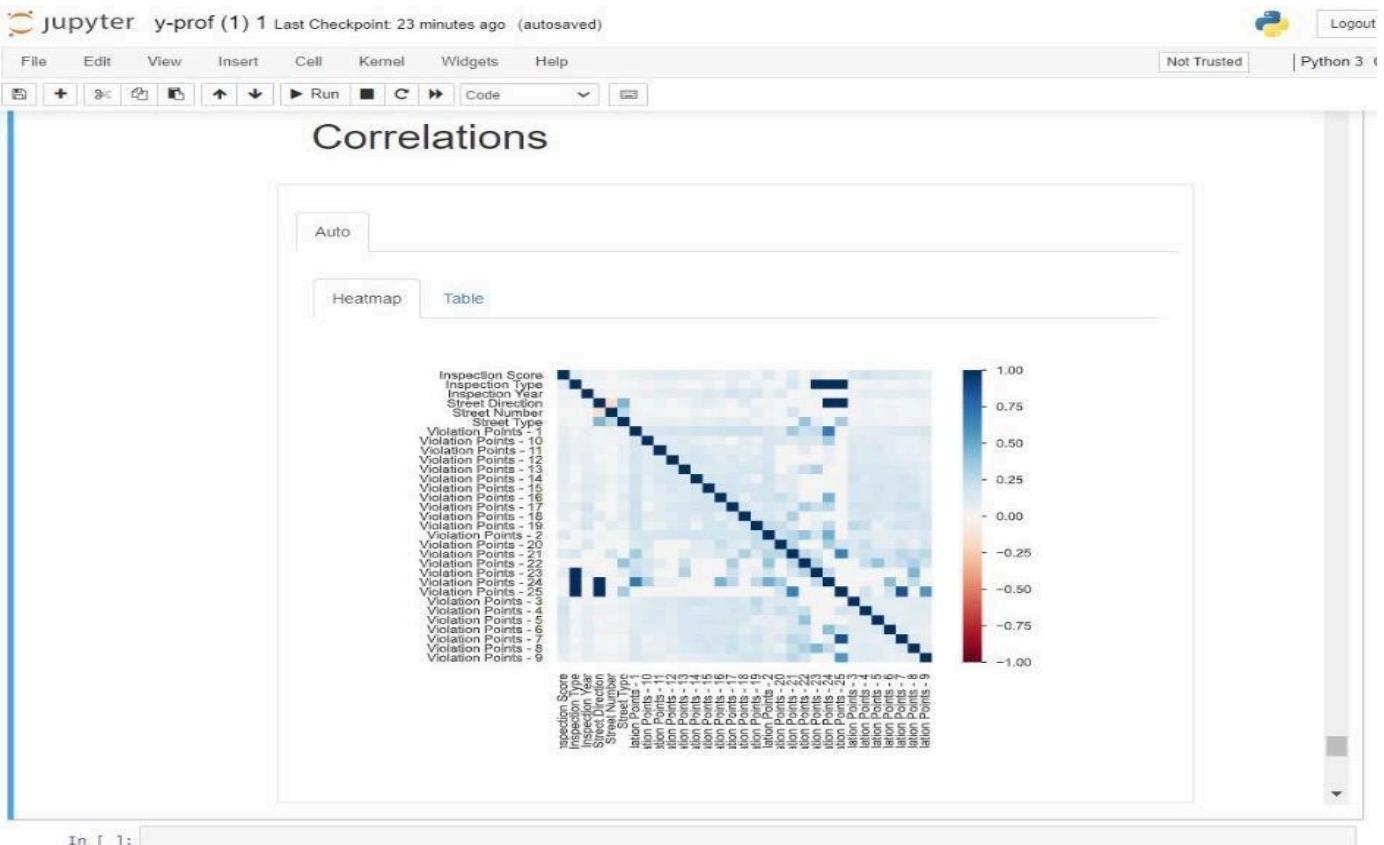
Select Columns

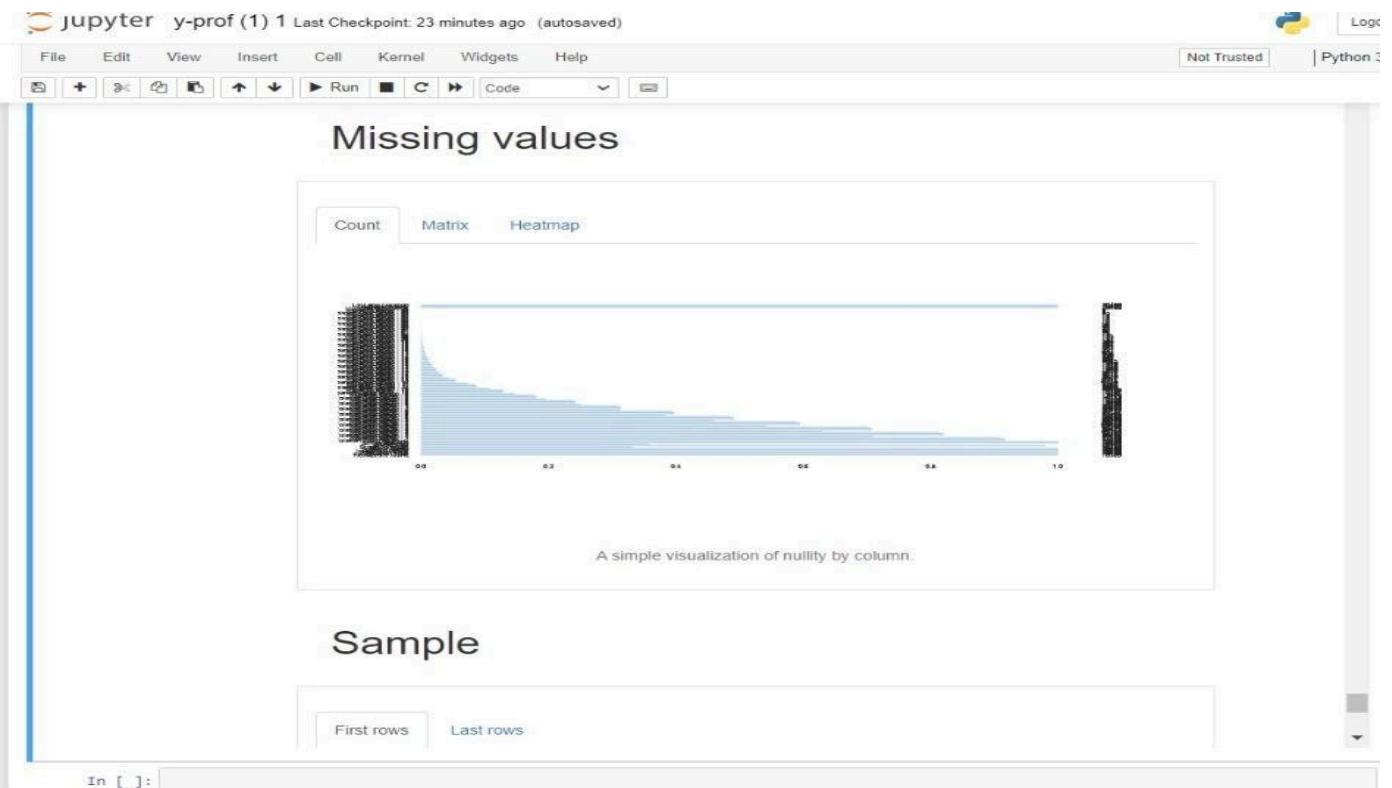
Restaurant Name
Text

Distinct	9136
Distinct (%)	11.7%
Missing	11
Missing (%)	< 0.1%
Memory size	612.6 KiB

Inspection Type

More details





In []:

jupyter y-prof (1) 1 Last Checkpoint: 23 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

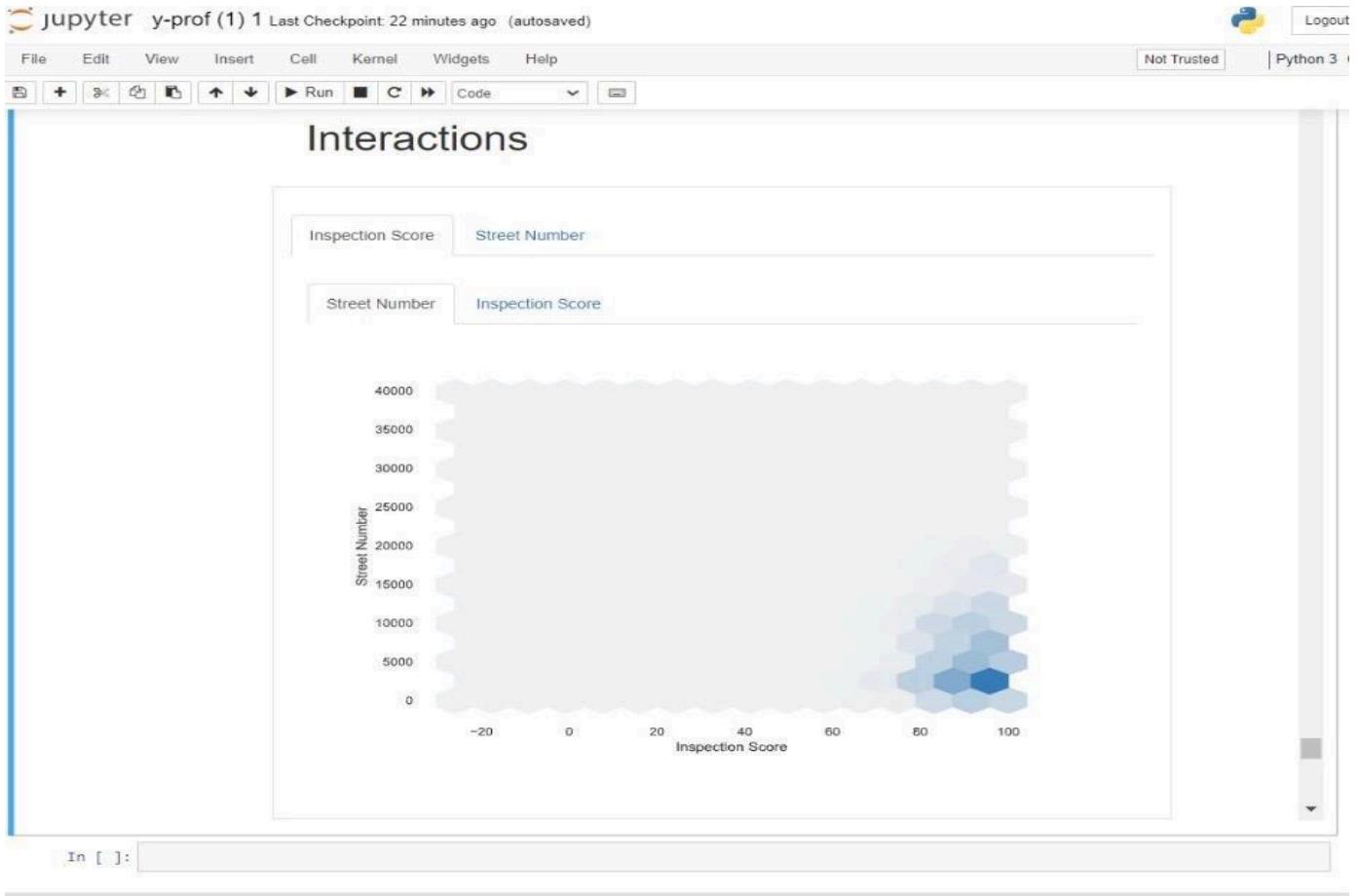
Not Trusted | Python 3

Sample

First rows Last rows

	Restaurant Name	Inspection Type	Inspection Date	Inspection Score
0	MICKLE CHICKEN	Routine	10/30/2019	100
1	TOM THUMB - JUICE BAR	Routine	04/08/2020	100
2	BROOKDALE WHITE ROCK	Routine	07/29/2020	97
3	CHURCH'S CHICKEN #201	Routine	09/14/2020	100
4	PEAK PREPARATORY-PRIMARY SCHOOL	Routine	04/25/2017	98
5	SAMS CLUB #6482 - DEMO ROOM	Routine	04/05/2021	100
6	AMC THEATRES NORTHPARK 15 (BAR)	Routine	04/09/2021	100
7	SAM'S CLUB #8282 - DEMO ROOM	Routine	11/14/2021	100
8	UCHI DALLAS	Routine	03/11/2020	98
9	UNREFINED BAKERY	Routine	03/23/2020	100

In []:



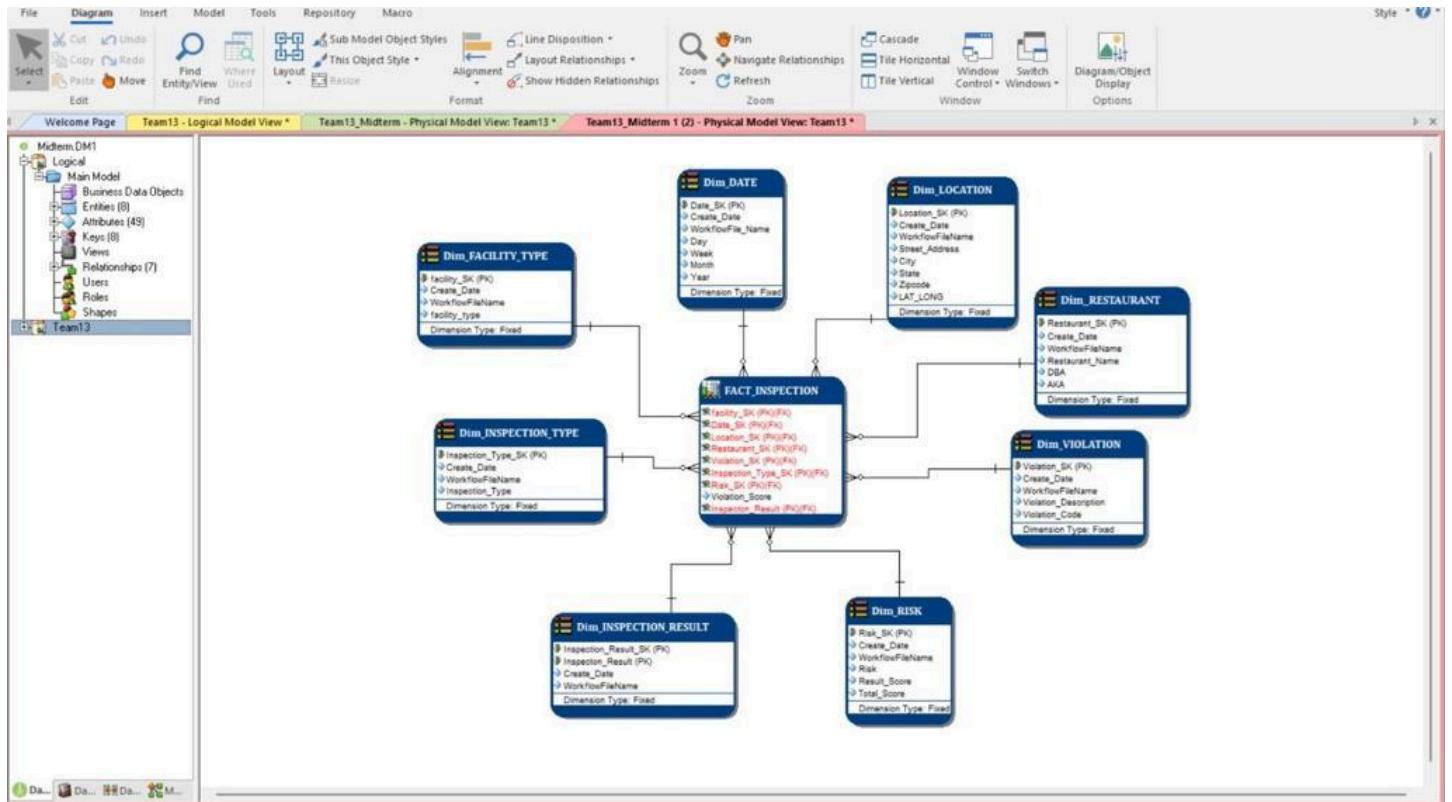
The report presents a comprehensive analysis encompassing various facets of data profiling. We systematically examined missing values, distinct values, duplicate rows, and inappropriate data formats, in addition to conducting a thorough review of column details. This meticulous examination ensures the quality and integrity of our data, laying a solid foundation for subsequent analysis.

Our thorough review of missing values, duplicate rows, and inappropriate data formats is pivotal in ensuring data quality. Identifying and addressing these issues at the outset enhances the reliability and accuracy of our analytical endeavors.

Moreover, our data profiling efforts were instrumental in delineating Dimensions and Facts essential for dimensional modeling. By pinpointing relevant attributes and measures early on, we streamline the subsequent design and implementation phases of our dimensional model. This proactive approach not only facilitates efficient data analysis and reporting but also underscores our commitment to delivering robust insights derived from high-quality data.

Part 2:

Designing dimensional model:



DDL SQL script:

```
/*
```

```
* TABLE: Dim_DATE
```

```
*/
```

```
CREATE TABLE Dim_DATE (
```

```
    Date_SK           int          NOT NULL,
```

```
    Create_Date      date        NOT NULL,
```

```
    WorkflowFileName varchar(100) NOT NULL,
```

```
        Day           varchar(10)      NOT NULL,  
  
        Week          varchar(10)      NOT NULL,  
  
        Month         varchar(10)      NOT NULL,  
  
        Year          varchar(10)      NOT NULL,  
  
CONSTRAINT PK7 PRIMARY KEY NONCLUSTERED (Date_SK)  
)
```

go

```
IF OBJECT_ID('Dim_DATE') IS NOT NULL  
  
PRINT '<<< CREATED TABLE Dim_DATE >>>'  
  
ELSE  
  
PRINT '<<< FAILED CREATING TABLE Dim_DATE >>>'  
  
go
```

```
/*  
* TABLE: Dim_FACILITY_TYPE  
*/  
  
CREATE TABLE Dim_FACILITY_TYPE (
```

```
facility_SK          int           NOT NULL,  
  
Create_Date         date          NOT NULL,  
  
WorkflowFileName    varchar(100)  NOT NULL,  
  
facility_type       varchar(100)  NULL,  
  
CONSTRAINT PK5 PRIMARY KEY NONCLUSTERED (facility_SK)  
)
```

```
go
```

```
IF OBJECT_ID('Dim_FACILITY_TYPE') IS NOT NULL  
  
PRINT '<<< CREATED TABLE Dim_FACILITY_TYPE >>>'  
  
ELSE  
  
PRINT '<<< FAILED CREATING TABLE Dim_FACILITY_TYPE >>>'  
  
go
```

```
/*  
* TABLE: Dim_INSPECTION_RESULT  
*/  
  
CREATE TABLE Dim_INSPECTION_RESULT (
```

```
Inspection_Result_SK      int            NOT NULL,  
  
Inspecton_Result          char(10)       NOT NULL,  
  
Create_Date                date           NOT NULL,  
  
WorkflowFileName           varchar(100)    NOT NULL,  
  
CONSTRAINT PK11_1 PRIMARY KEY NONCLUSTERED (Inspection_Result_SK, Inspecton_Result)  
)
```

```
go
```

```
IF OBJECT_ID('Dim_INSPECTION_RESULT') IS NOT NULL  
  
PRINT '<<< CREATED TABLE Dim_INSPECTION_RESULT >>>'  
  
ELSE  
  
PRINT '<<< FAILED CREATING TABLE Dim_INSPECTION_RESULT >>>'  
  
go
```

```
/*  
* TABLE: Dim_INSPECTION_TYPE  
*/
```

```
CREATE TABLE Dim_INSPECTION_TYPE (
```

```
Inspection_Type_SK      int           NOT NULL,  
  
Create_Date             date          NOT NULL,  
  
WorkflowFileName        varchar(100)  NOT NULL,  
  
Inspection_Type        varchar(100)  NULL,  
  
CONSTRAINT PK11 PRIMARY KEY NONCLUSTERED (Inspection_Type_SK)  
)
```

```
go
```

```
IF OBJECT_ID('Dim_INSPECTION_TYPE') IS NOT NULL  
  
PRINT '<<< CREATED TABLE Dim_INSPECTION_TYPE >>>'  
  
ELSE  
  
PRINT '<<< FAILED CREATING TABLE Dim_INSPECTION_TYPE >>>'  
  
go
```

```
/*  
* TABLE: Dim_LOCATION  
*/
```

```
CREATE TABLE Dim_LOCATION (
```

```
Location_SK           int            NOT NULL,  
  
Create_Date          date           NOT NULL,  
  
WorkflowFileName     varchar(100)   NOT NULL,  
  
Street_Address       varchar(100)   NOT NULL,  
  
City                varchar(50)    NULL,  
  
State               varchar(50)    NULL,  
  
Zipcode             varchar(50)    NULL,  
  
LAT_LONG            decimal(50, 0)  NULL,  
  
CONSTRAINT PK8 PRIMARY KEY NONCLUSTERED (Location_SK)  
  
)
```

go

```
IF OBJECT_ID('Dim_LOCATION') IS NOT NULL  
  
PRINT '<<< CREATED TABLE Dim_LOCATION >>>'  
  
ELSE  
  
PRINT '<<< FAILED CREATING TABLE Dim_LOCATION >>>'  
  
go  
  
/*
```

```
* TABLE: Dim_RESTAURANT
```

```
*/
```

```
CREATE TABLE Dim_RESTAURANT(
```

Restaurant_SK	int	NOT NULL,
Create_Date	date	NOT NULL,
WorkflowFileName	varchar(100)	NOT NULL,
Restaurant_Name	varchar(100)	NULL,
DBA	varchar(100)	NULL,
AKA	varchar(100)	NULL,

```
CONSTRAINT PK2 PRIMARY KEY NONCLUSTERED (Restaurant_SK)
```

```
)
```

```
go
```

```
IF OBJECT_ID('Dim_RESTAURANT') IS NOT NULL
```

```
PRINT '<<< CREATED TABLE Dim_RESTAURANT >>>'
```

```
ELSE
```

```
PRINT '<<< FAILED CREATING TABLE Dim_RESTAURANT >>>'
```

```
go
```

```
/*
* TABLE: Dim_RISK

CREATE TABLE Dim_RISK(
    Risk_SK           int          NOT NULL,
    Create_Date       date         NOT NULL,
    WorkflowFileName varchar(100) NOT NULL,
    Risk              numeric(100, 0) NULL,
    Result_Score      numeric(10, 0) NULL,
    Total_Score       numeric(10, 0) NULL,
    CONSTRAINT PK20 PRIMARY KEY NONCLUSTERED (Risk_SK)
)

go

IF OBJECT_ID('Dim_RISK') IS NOT NULL
PRINT '<<< CREATED TABLE Dim_RISK >>>'
ELSE
```

```
PRINT '<<< FAILED CREATING TABLE Dim_RISK >>>'  
  
go  
  
/*  
 * TABLE: Dim_VIOLATION  
 */  
  
CREATE TABLE Dim_VIOLATION(  
    Violation_SK           int          NOT NULL,  
    Create_Date            date         NOT NULL,  
    WorkflowFileName       varchar(100) NOT NULL,  
    Violation_Description varchar(1000) NULL,  
    Violation_Code         int          NULL,  
    CONSTRAINT PK10 PRIMARY KEY NONCLUSTERED (Violation_SK)  
)  
  
go  
  
IF OBJECT_ID('Dim_VIOLATION') IS NOT NULL  
PRINT '<<< CREATED TABLE Dim_VIOLATION >>>'
```

```

ELSE

PRINT '<<< FAILED CREATING TABLE Dim_VIOLATION >>>'

go

/*
* TABLE: FACT_INSPECTION
*/
CREATE TABLE FACT_INSPECTION(
    facility_SK           int          NOT NULL,
    Date_SK               int          NOT NULL,
    Location_SK            int          NOT NULL,
    Restaurant_SK          int          NOT NULL,
    Violation_SK           int          NOT NULL,
    Inspection_Type_SK     int          NOT NULL,
    Risk_SK                int          NOT NULL,
    Violation_Score        int          NULL,
    Inspector_Result       char(10)    NOT NULL,
    CONSTRAINT PK12 PRIMARY KEY NONCLUSTERED (facility_SK, Date_SK, Location_SK, Restaurant_SK,
    Violation_SK, Inspection_Type_SK, Risk_SK, Inspector_Result)
)

```

```
go
```

```
IF OBJECT_ID('FACT_INSPECTION') IS NOT NULL  
  
PRINT '<<< CREATED TABLE FACT_INSPECTION >>>'  
  
ELSE
```

```
PRINT '<<< FAILED CREATING TABLE FACT_INSPECTION >>>'
```

```
go
```

```
/*  
  
* TABLE: FACT_INSPECTION  
  
*/
```

```
ALTER TABLE FACT_INSPECTION ADD CONSTRAINT RefDim_FACILITY_TYPE1
```

```
FOREIGN KEY (facility_SK)  
  
REFERENCES Dim_FACILITY_TYPE(facility_SK)
```

```
go
```

```
ALTER TABLE FACT_INSPECTION ADD CONSTRAINT RefDim_DATE2
```

```
FOREIGN KEY (Date_SK)  
  
REFERENCES Dim_DATE(Date_SK)
```

go

```
ALTER TABLE FACT_INSPECTION ADD CONSTRAINT RefDim_LOCATION3  
FOREIGN KEY (Location_SK)  
REFERENCES Dim_LOCATION(Location_SK)
```

go

```
ALTER TABLE FACT_INSPECTION ADD CONSTRAINT RefDim_RESTAURANT4  
FOREIGN KEY (Restaurant_SK)  
REFERENCES Dim_RESTAURANT(Restaurant_SK)
```

go

```
ALTER TABLE FACT_INSPECTION ADD CONSTRAINT RefDim_VIOLATIONS5  
FOREIGN KEY (Violation_SK)  
REFERENCES Dim_VIOLATION(Violation_SK)
```

go

```
ALTER TABLE FACT_INSPECTION ADD CONSTRAINT RefDim_INSPECTION_TYPE6  
FOREIGN KEY (Inspection_Type_SK)  
REFERENCES Dim_INSPECTION_TYPE(Inspection_Type_SK)
```

go

```
ALTER TABLE FACT_INSPECTION ADD CONSTRAINT RefDim_RISK7  
FOREIGN KEY (Risk_SK)  
REFERENCES Dim_RISK(Risk_SK)  
  
go
```

```
ALTER TABLE FACT_INSPECTION ADD CONSTRAINT RefDim_INSPECTION_RESULT8  
FOREIGN KEY (Inspection_Type_SK, Inspector_Result)  
REFERENCES Dim_INSPECTION_RESULT(Inspection_Result_SK, Inspector_Result)  
  
go
```

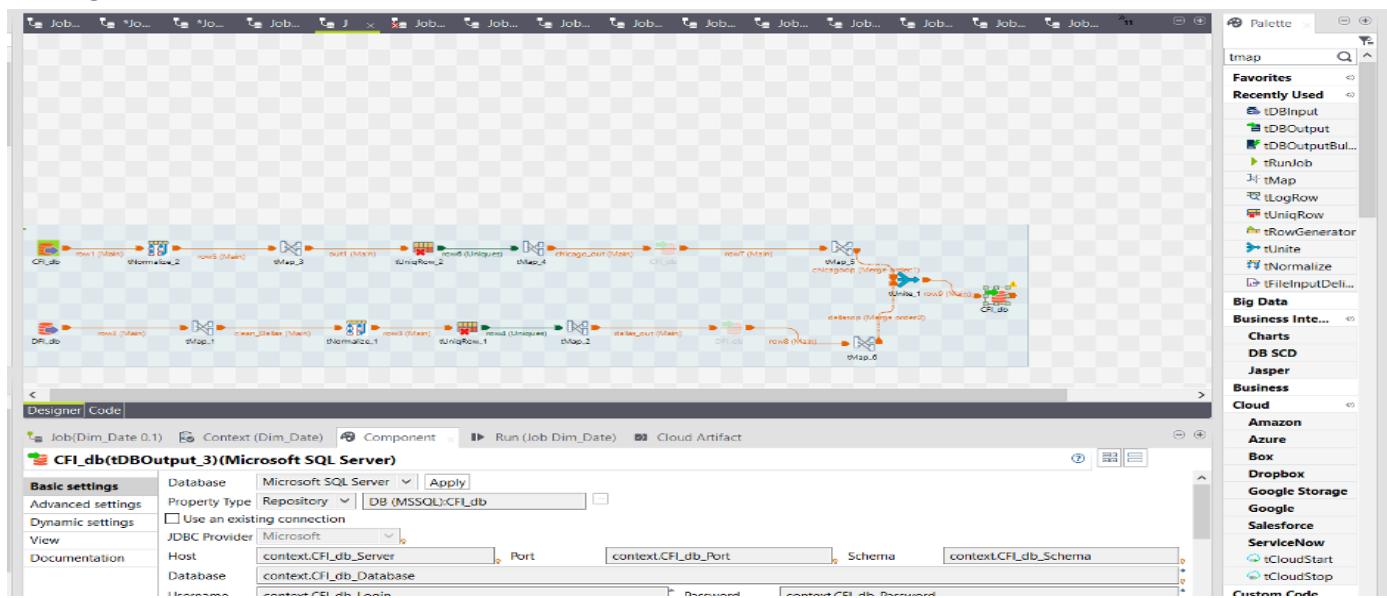
Schema in database:

Please attach ss

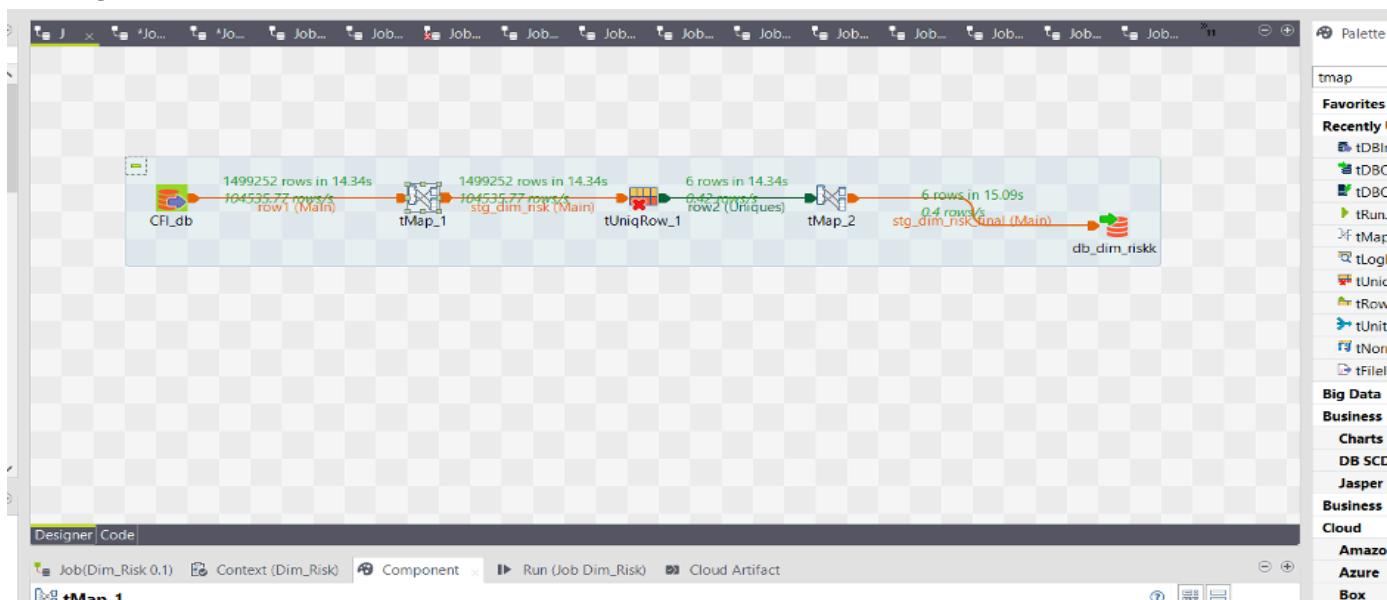
Part 3:

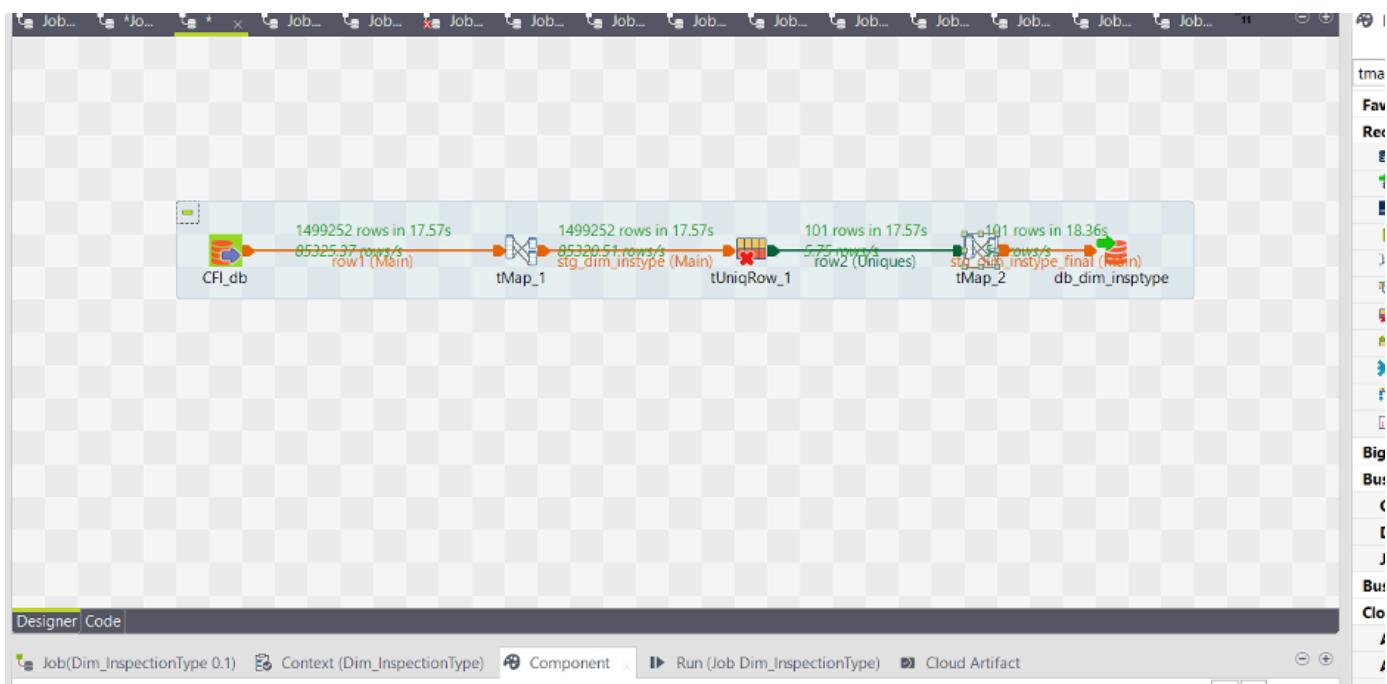
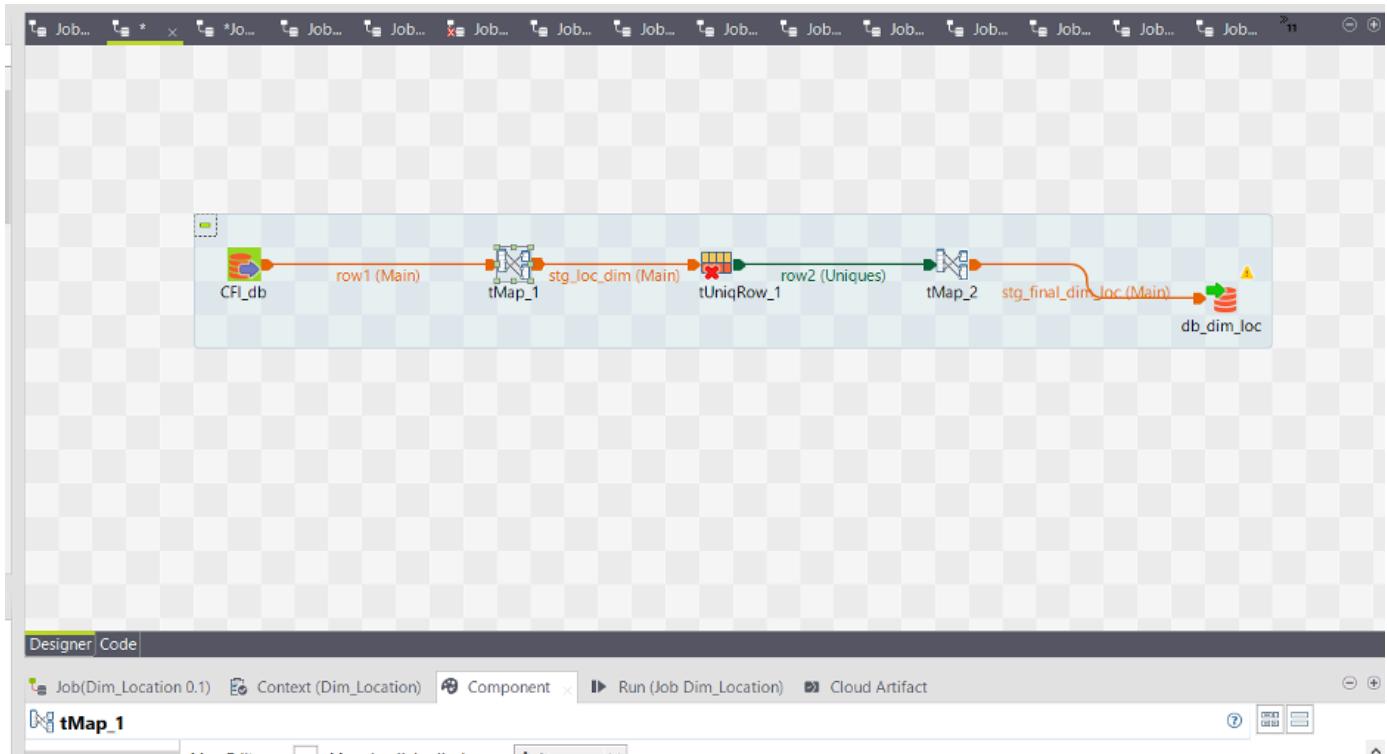
Dimensional Modeling using Talend:

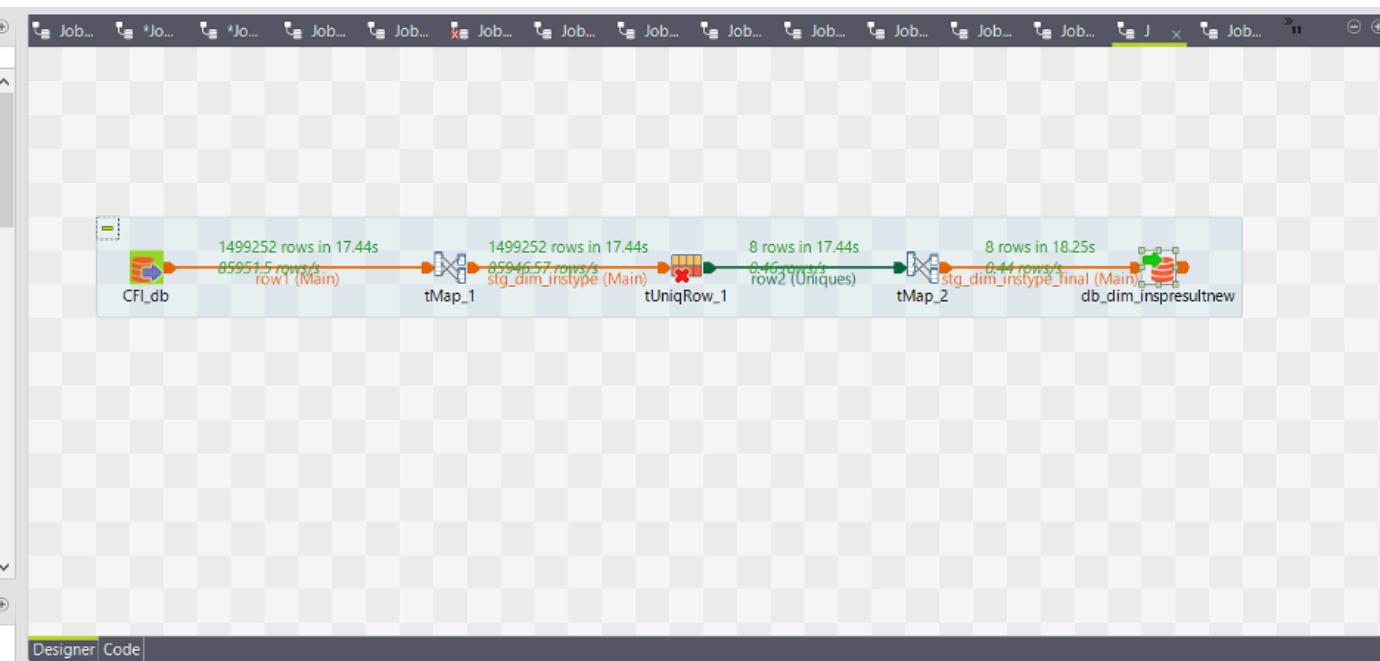
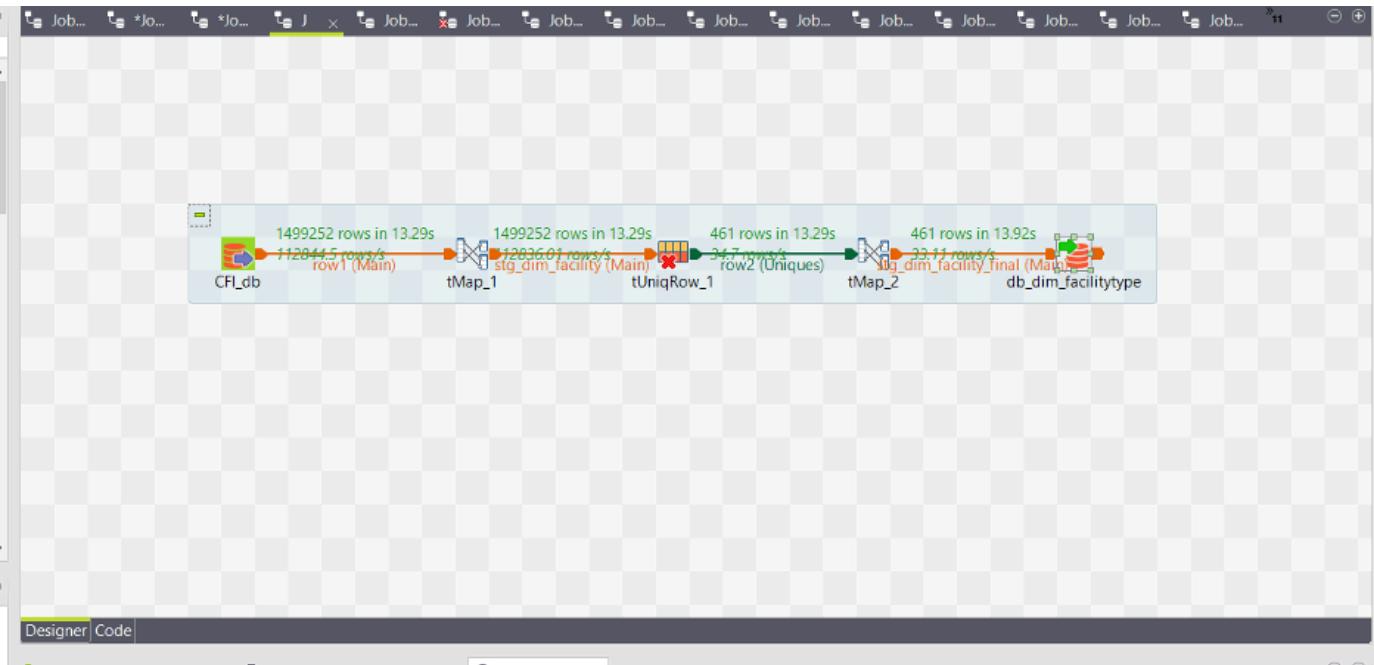
Uniting both datasets:



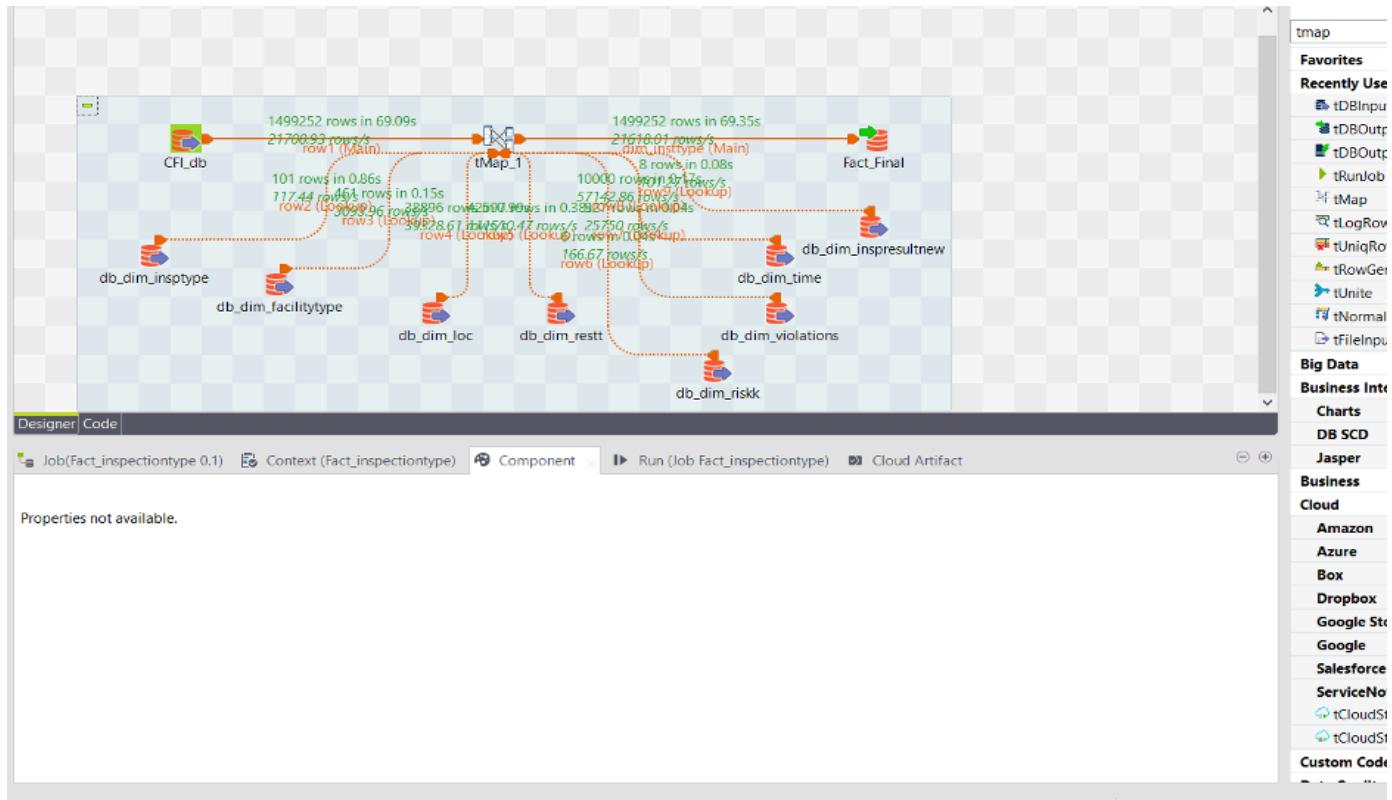
Loading Dimensions:





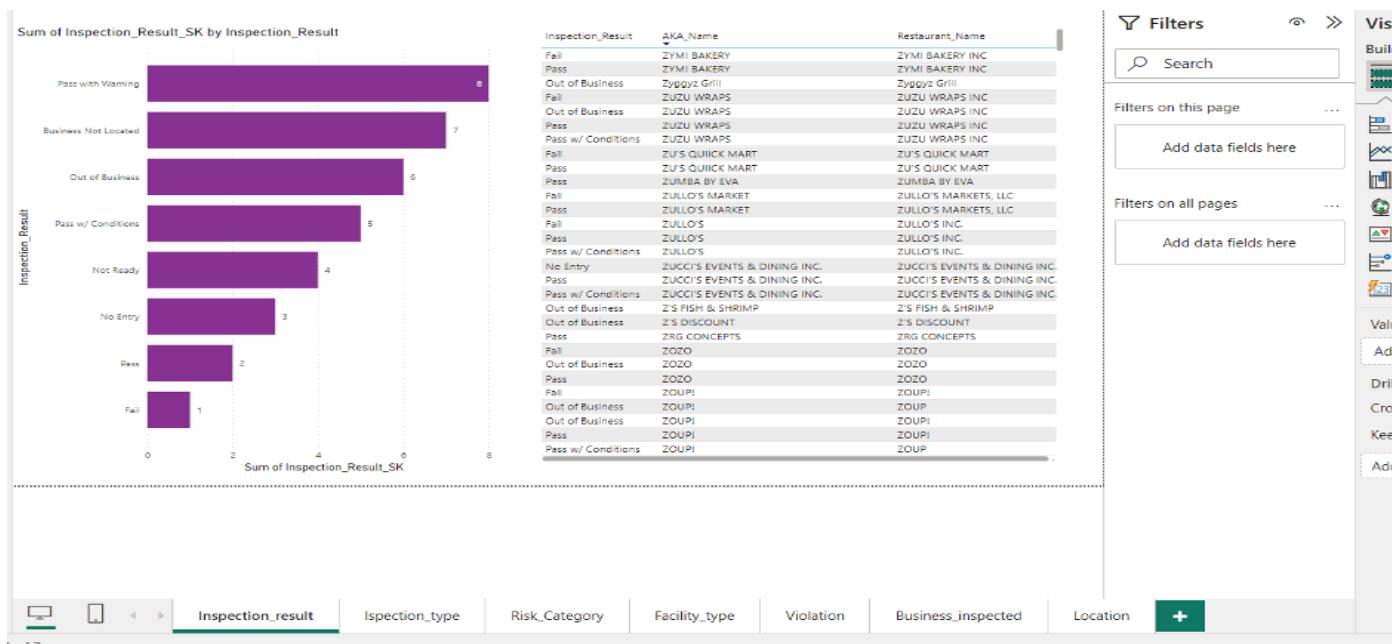


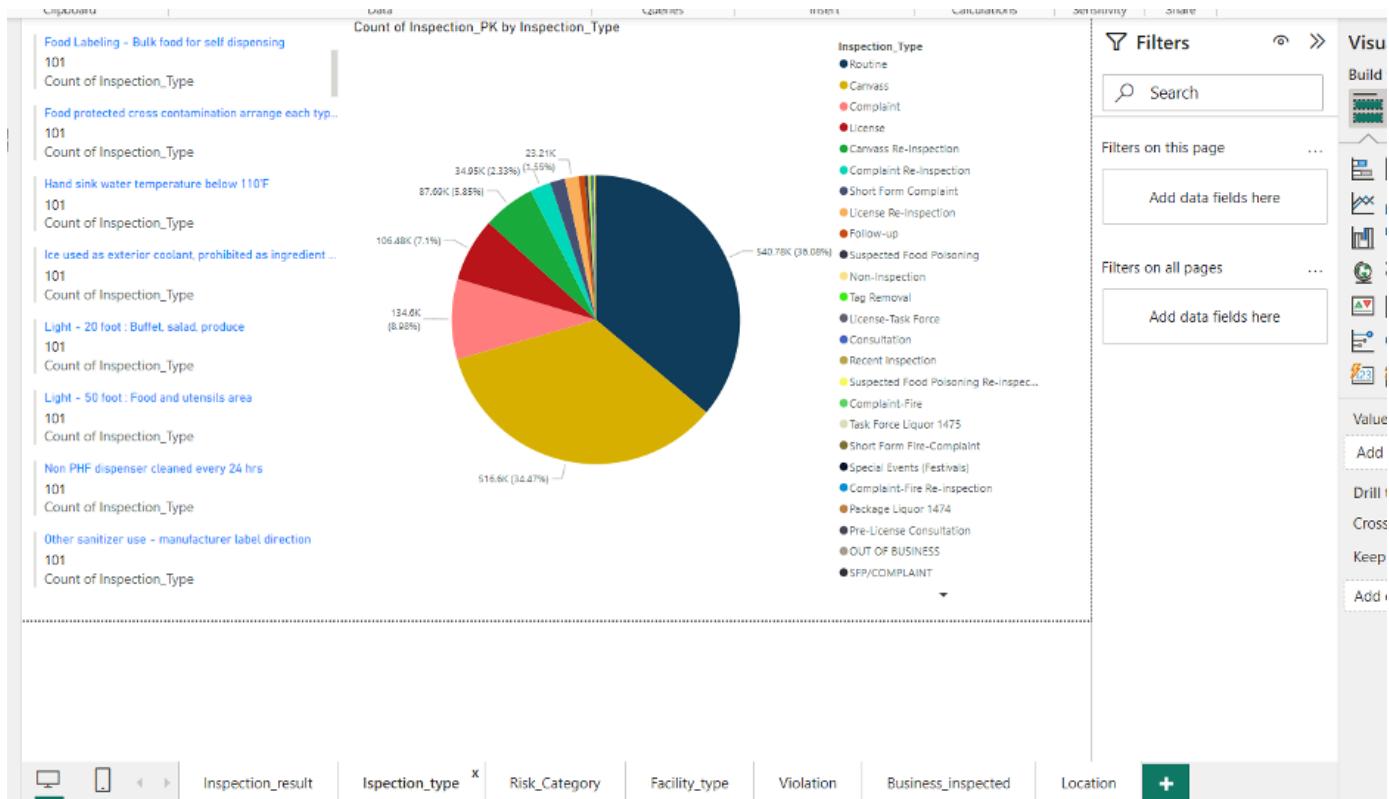
Loading Fact table:



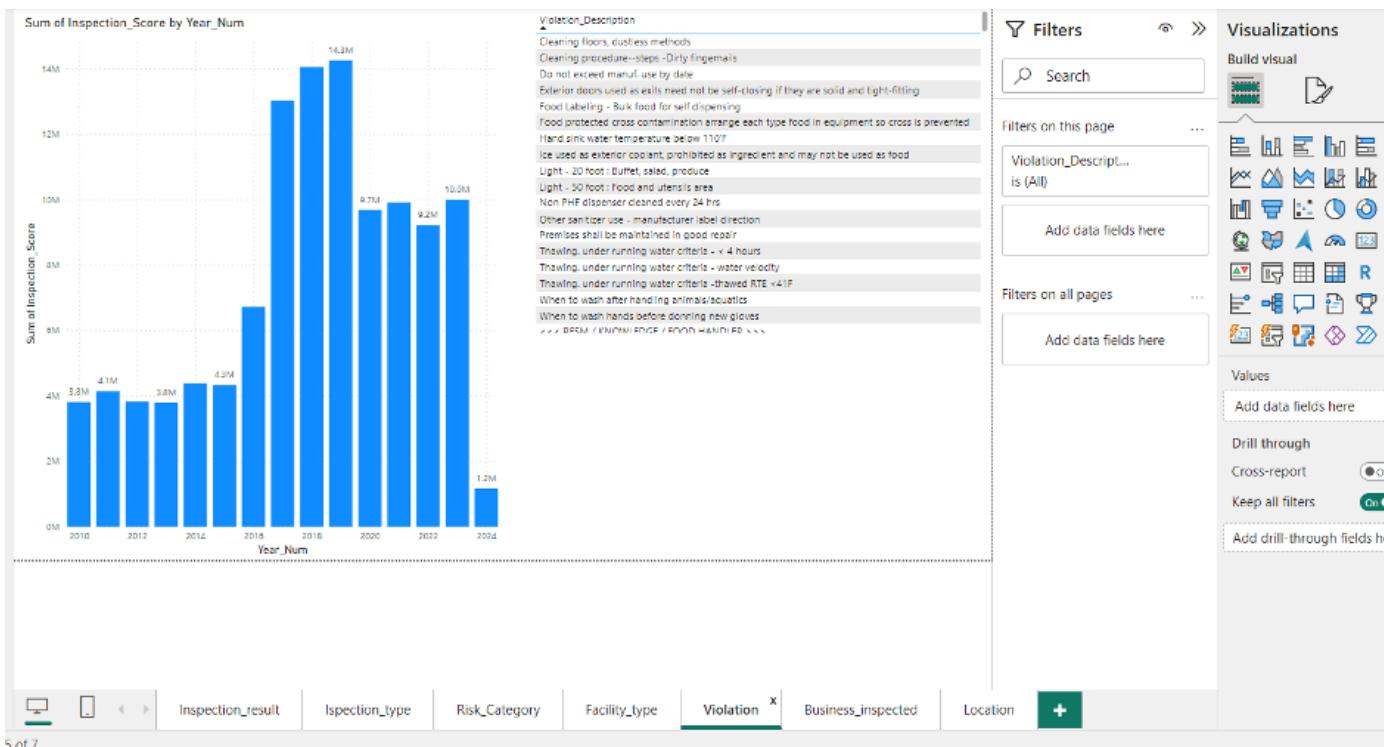
Part 4:

Power BI Dashboards:

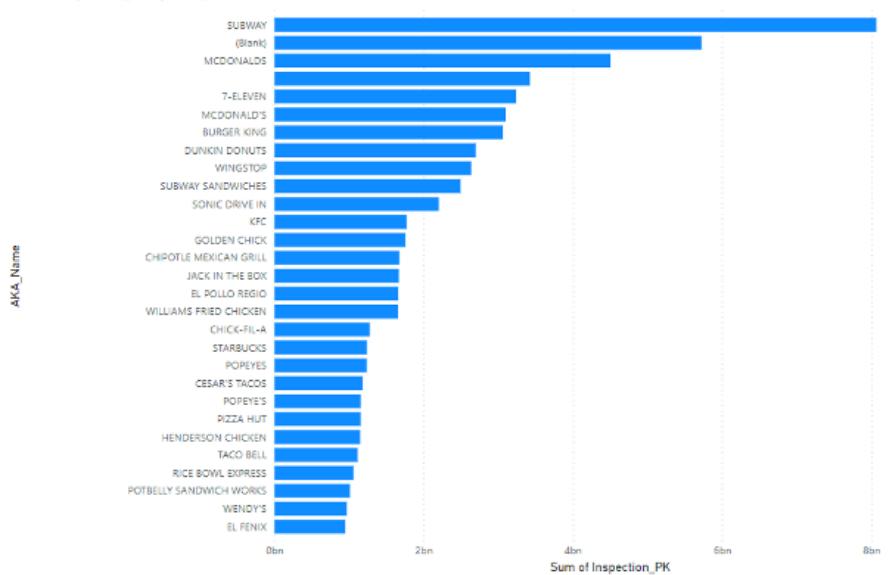




3 of 7



Sum of Inspection_PK by AKA_Name



Filters

Search

Filters on this page

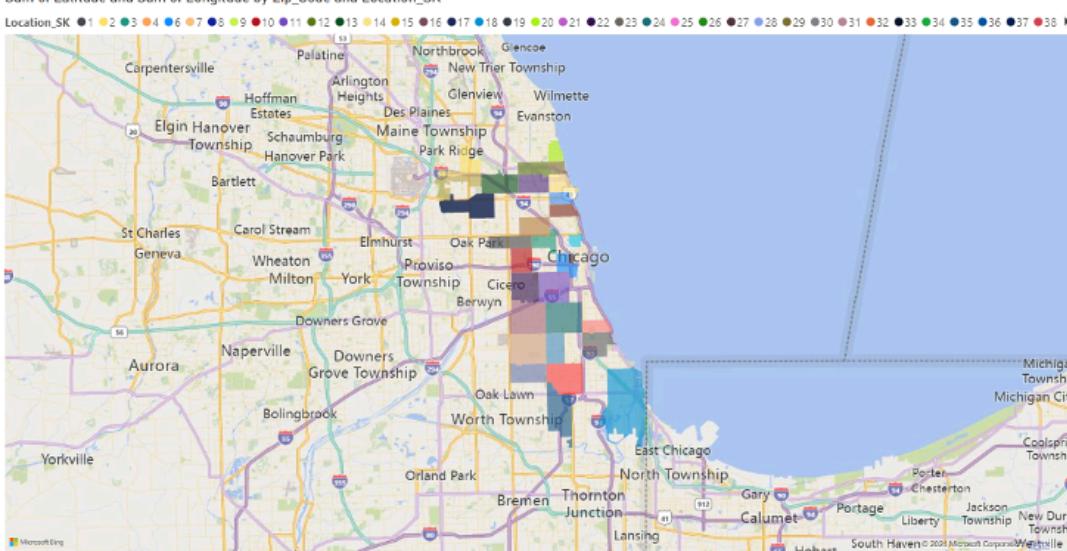
Add data fields here

Filters on all pages

Add data fields here

Inspection_result	Inspection_type	Risk_Category	Facility_type	Violation	Business_inspected	Location	+ Add
-------------------	-----------------	---------------	---------------	-----------	--------------------	----------	-------

Sum of Latitude and Sum of Longitude by Zip_Code and Location_SK



Filters

Search

Filters on this page

Add data fields here

Filters on all pages

Add data fields here

Inspection_result	Inspection_type	Risk_Category	Facility_type	Violation	Business_inspected	Location	+ Add
-------------------	-----------------	---------------	---------------	-----------	--------------------	----------	-------

Tableau Dashboards:

Tableau - Book1

File Data Server Window Help

Connections Add
DESKTOP-JVN42L Microsoft SQL Server

Database master

Table

- Dallas_11
- Dallas_stg
- DAMG1
- DATE
- Dim_DATE
- dim_date_final
- Dim_FACILITY_TYPE
- dim_facilitytype_final
- Dim_INSPECTION_RESULT
- Dim_INSPECTION_TYPE
- dim_inspresult_final
- dim_inspstype_final
- dim_loc_final
- New Custom SQL
- New Union
- New Table Extension

Stored Procedures

- sp_Mscleanupmergepublisher
- sp_Mscreate_subscription

FACT_INSPECTION_new+ (master)

Connections Connection Live Extract Filters 0 | Add

Dim_INSPECTION_RESULT

Risk Category

For horizontal bars try

0 or more Dimensions

1 or more Measures

Horizontal Bar Chart

Sum(Risk SK)

Risk SK

500K 1000K 1500K 2000K 2500K 3000K 3500K 4000K 4500K 5000K

OK

Tables

- dim_date_final
- dim_facilitytype_final
- Dim_INSPECTION_RESULT
- dim_inspstype_final

Marks

- Automatic
- Color
- Size
- Label
- Detail
- Tooltip
- Risk

