

# **CSI-4142**

# **Fundamentals Of Data Science**

Prof. Herna L. Victor

**SUBMISSION DATE**

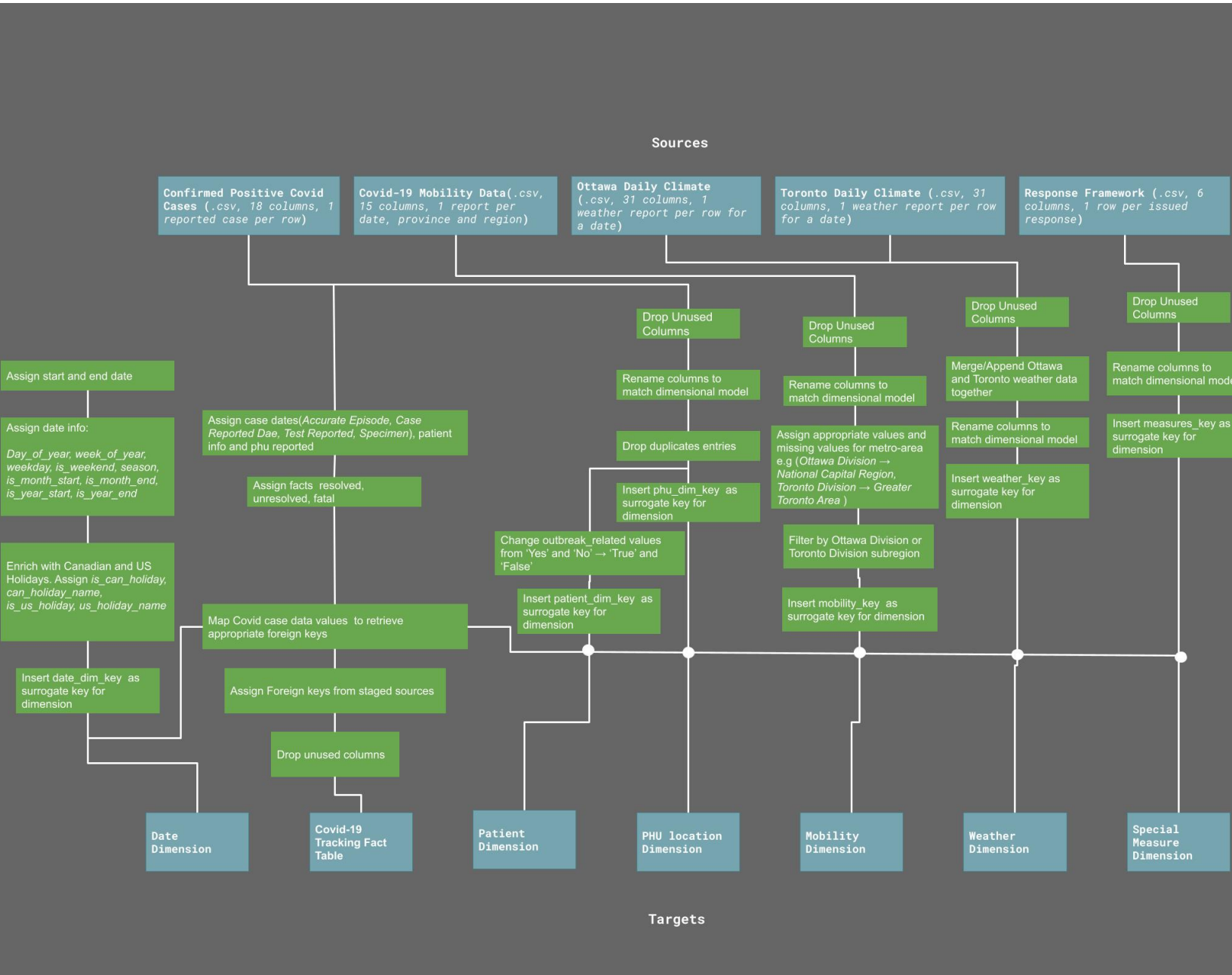
March 10th, 2021

**BILLAL ZAZAI, 8572975**

**YOUSSEF MAKBOUL, 8609614**

**GRAYDON HOPE, 300045044**

# Schema Design



**Figure 1. Covid-19 Data Mart Schema-** Illustrates a high level overview of the data staging level of Data Mart

## Data Quality

Fortunately, the data sets we were provided and found had very little missing data and data quality issues. However, we still encountered poor data quality such as ...

- To minimize type conversion we changed outbreak\_related values from 'Yes' and 'No' → 'True' and 'False' in the Patient Dimension
- The Covid-19 Mobility data set had missing values for the metro\_area column. Therefore, utilizing the sub-region2 column we were able to assign the appropriate values and missing values for metro-area e.g (Ottawa Division → National Capital Region, Toronto Division → Greater Toronto Area )
- When merging Ottawa Daily Climate and Toronto Daily climate checking Referential Integrity(RI) was necessary to ensure the correct data is being merged
- To eliminate repetitive data, removal of duplicate data was required for certain dimensions such as Patient Dimension, PHU Location Dimension etc
- Column renaming was necessary to maintain consistency in the naming standards throughout all the dimensions

Tasks were evenly distributed among team members

