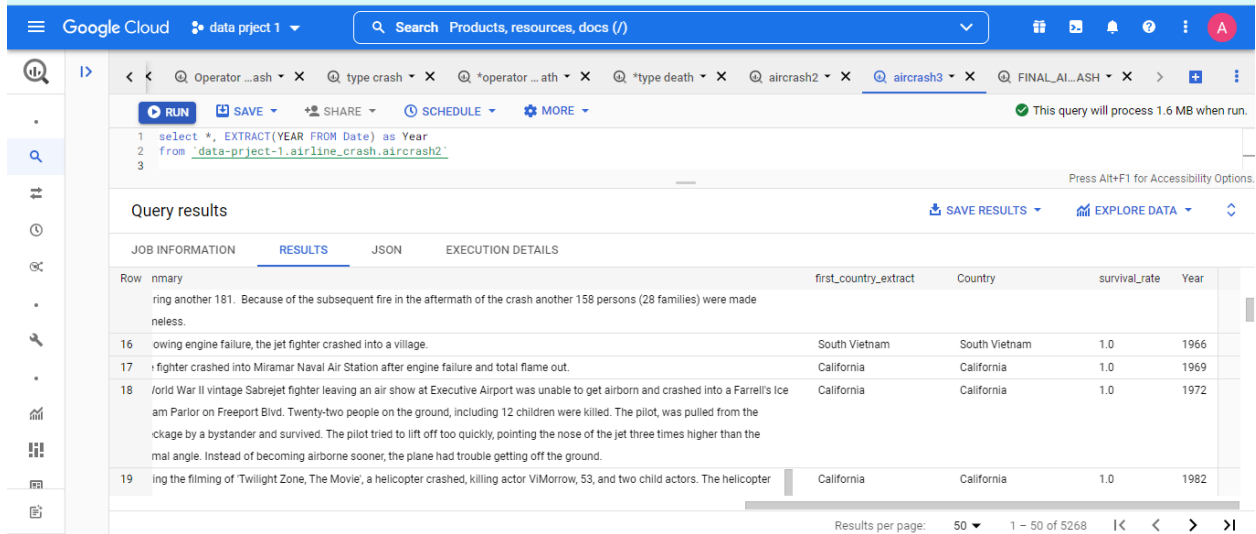# PRE-ANALYSIS

To derive more insights from the data, I created some extra columns from the existing ones using SQL. This also helped in my exploratory analysis to get initial findings from the data.

The added columns and reasons for creating them are discussed below.

## Year Column

Here, I extracted each crash year from the date column to be able to visualize yearly trends.



## Month Column

Extracted each rows month name from the date column. This was done during data transformation in power query.

## Total Deaths

This column shows the total number of deaths for each crash. This column was derived by adding the fatalities and ground columns. This was needed to get the death trend over the years.



## Survived count

This column was created in power query. The survivors count was created by subtracting number of fatalities from number of aboard passengers and crew.

## Survived or not

This column tells us whether or not the crash had survivors. The resulting column returns 0 for crash with no survivors and 1 for crash with survivors. Taking the count of these categories will help in analysis on rate of survivors.

## Exact Location

This column contains the drilled down exact location of the crash. The values in this column is taken from the location column using the string position function to get the last string in the locations and the trim function to remove whitespaces.



## NOTE

Initially, the dataset contained 5,268 rows. During data transformation in power query, 89 rows were dropped or filtered. These rows were dropped because they had null values in either of the following columns: fatalities, aboard, ground, country, operator and type. Columns not needed for the analysis were also removed. The dataset used for the analysis therefore contained 5,179 rows after cleaning.