Milestone 6 | Traffic Collisions in California

INTRODUCTION: Data is often stored across multiple tables to keep the storage requirements compact, and to organize different types of data. Knowing how to use a join is a vital skill when working with data, since bringing tables together can open the door to additional insights that are cumbersome or impossible looking at just one table at a time.

In this Milestone, you'll use your proficiency with joins to help a reporter in California use data to support an article they're writing on the causes of motor vehicle accidents. In particular, they want some information about how many accidents are caused by the influence of alcohol, or due to inattention (such as using a cell phone to text or talk to others), and when these types of accidents tend to occur.

HOW IT WORKS: Follow the prompts in the questions below to investigate your data. Post your answers in the provided boxes: the **yellow boxes** for the queries you write, **purple boxes** for visualizations and **blue boxes** for text-based answers. When you're done, export your document as a pdf file and submit it on the Milestone page – see instructions for creating a PDF at the end of the Milestone.

RESOURCES: If you need hints on the Milestone or are feeling stuck, there are multiple ways of getting help. Attend Drop-In Hours to work on these problems with your peers, or reach out to the HelpHub if you have questions. Good luck!

PROMPT: To help the reporters out, you will be making use of data regarding traffic accidents in the state of California released by the California Highway Patrol. Certain insights can be found by looking at data on the incident level, while other insights are possible by looking deeper at the parties involved in an incident. But to make insights across those two levels, we need a join to be able to relate the unique information contained in each table.

SQL App: <u>Here's that link</u> to our specialized SQL app, where you'll write your SQL queries and interact with the data.

- Data Set **Description**

Data for this Milestone comes from the California Highway Patrol's Statewide Integrated Traffic Records System (SWITRS). The SWITRS data we've provided (switrs.*) consists of two tables from the 2019 data collection: collisions and parties. The tables are related hierarchically. At the top level, there is a unique row and identifier for each incident in the collisions table. Then, in the lower level, each collision is between one or more parties, which include vehicles, pedestrians, etc.

The original collisions table has 469 664 rows and 76 columns, but we'll be focusing on only the following four columns in this Milestone:

- case_id unique identifier for each collision
- collision_time time of day when collision occurred, in 24 hour format
- **day_of_week** day of week when collision occurred. Note that numbering starts at 1 = Monday and ends at 7 = Sunday (instead of 0 = Sunday)
- party_count number of parties involved in the collision

The original parties table has 940 216 rows and 33 columns, with the following five columns of interest:

- **case_id** associated with a collision with matching case_id, may not be unique
- party_number numbering of parties involved, always starts from 1 for each collision
- at_fault Y/N indicating whether party was at fault for collision
- party_sobriety encodings for whether or not the party had been drinking
- oaf_1, oaf_2 encodings for other associated factors

Most of the features in the dataset are coded in some way for efficient data storage, which can make working with highly detailed data like this tricky. This includes the party_sobriety, oaf_1, and oaf_2 columns you'll be investigating in the Milestone. Don't sweat that point, though: the instructions will explain the encoding values relevant to the tasks.

– Task 1: How frequently does alcohol use or lack of attention feature in accidents?

To start, we should run some queries on the parties table to understand how fault, alcohol use, and inattention are attributed to accidents.

A. Write a query that answers the following question: According to this dataset, how many people are at fault for a collision?

```
SELECT
COUNT(at_fault)
FROM switrs.parties
```

B. The party_sobriety field takes on a value of 'B' when the party is known to have been drinking, and under the influence of alcohol. Modify your query from part A to answer the following question: How many parties were found at fault while under the influence of alcohol?

```
SELECT
COUNT(at_fault)
FROM switrs.parties
WHERE party_sobriety = 'B'
```

35059 drivers where under the influence when they crashed

C. The **oaf_1** or **oaf_2** feature takes on a value of 'F' if inattention was a factor in the collision. Modify your query to answer the following question: How many parties were found at fault while lack of attention was a factor in the collision?

SELECT
COUNT(at_fault)
FROM switrs.parties
WHERE oaf_1 = 'F'OR oaf_2 = 'F'

19,722 drivers where at fault due to not paying attention

- Task 2: When do accidents occur by day of the week?

Now that we have a way to identify whether or not a collision can be attributed to alcohol or inattention, let's add in the collisions table to answer the journalist's question of whether or not there are differences between the two accident sources.

A. Let's start with the collisions table on its own. Write a query that returns the number of collisions, grouped by day of the week. Which days have the highest number of collisions, and which days have the least number? Note: Day of week is encoded slightly differently than what comes out of the date_part function: Sunday is indicated by a 7 instead of a 0.

SELECT

day_of_week as day_of_accident, COUNT(case_id) as n_accidents FROM switrs.collisions GROUP BY day_of_accident ORDER BY n_accidents desc

Friday (day 5) had the most accidents with 75,654 accidents. People are more likely to be partying on Fridays

B. The collisions table and parties tables share values in the case_id column. Write a new query that inner joins the two tables on that column, returning the number of rows. How many rows are in the combined output table, and why?

```
SELECT
COUNT(*) as total_rows
FROM
switrs.collisions AS a
INNER JOIN switrs.parties as b ON a.case_id = b.case_id
```

There are 940,216 rows representing each accident

C. Combine the queries from parts A and B to return the number of collisions grouped by the day of the week. Add a condition for the involved parties so that we only count accidents where the party was found to be at fault AND under the influence of alcohol. Which days have the highest number of collisions, and which days have the smallest number?

```
SELECT
   a.day_of_week,
   COUNT(b.case_id) as total_rows
FROM
   switrs.collisions AS a
   INNER JOIN switrs.parties as b ON a.case_id = b.case_id
GROUP BY
   a.day_of_week
```

Fridays have the most accidents and Sundays and Saturdays have the least. I'm guessing that's because people are more irresponsible on Friday nights, then stay home during the weekends. Weekdays typically see the most traffic with people commuting to and from work so it would make sense that these days have a higher accident rate

D. Modify your query to look at the number of accidents by the day of the week where the party was found to be at fault AND inattention was a factor. Which days have the highest number of collisions, and which days have the smallest number?

```
SELECT
a.day_of_week,
COUNT(b.case_id) as total_rows
FROM
switrs.collisions AS a
INNER JOIN switrs.parties as b ON a.case_id = b.case_id
WHERE
b.at_fault = 'B' AND b.oaf_1 = 'F' OR b.oaf_2 = 'F'
GROUP BY
a.day_of_week
```

ORDER BY a.day_of_week;

It would seem the weekends have the highest amount of accidents where the at fault party wasnt paying attention. This could be because people are a bit more switched off during the weekend

- Task 3: When do accidents occur by the time of day?

A data analyst colleague of yours has taken interest in your project with the journalist and has pitched in their own contribution by providing you a summary of the dataset with five features:

- **alcohol_involved** TRUE/FALSE whether or not the party at fault was under the influence of alcohol
- **inattention_involved** TRUE/FALSE whether or not inattention was a factor for the party at fault
- **day_of_week** day of week when collision occurred. Note that numbering starts at 1 = Monday and ends at 7 = Sunday (instead of 0 = Sunday)
- hour_of_day -hour of day when collision occurred, in 24 hour format (0-2300). Values of 2500 indicate an unknown time of day.
- **n_collisions** number of collisions matching the conditions of the first four columns

Let's use this new data summary to look at how accident patterns change based on the time of day. Since the data has already been queried, we'll do this visually within Tableau! <u>Click this link</u> to navigate to the workbook you'll use to complete the remainder of this Milestone. Once you've published your Tableau Workbook in the folder named Upload Workbooks Here, paste the Share Link in the box below. https://prod-useast-b.online.tableau.com/#/site/globaltech/w orkbooks/1226470?:origin=card_share_link

Continue to post your answers in the provided boxes: purple boxes for your visualizations, and blue boxes for text-based answers.

A. On Sheet 1, create a bar chart of the number of collisions by the hour of day. Describe the pattern in the data. Are there times of day where more accidents occur? Does this fit in with your expectations?

HINT: Drag the Hour Of Day pill to the Columns and the N Collisions pill to the Rows.



As I expected the most accidents occur when people are commuting too and from work. With the most accidents happening at the peak of rush hour when people just want to get home. **B.** Copy the chart into a new sheet and add a filter so that the bar chart only shows accidents where the party at fault was found to be under the influence of alcohol. How does this distribution of accidents by time of day compare to the overall distribution?



As I expected the hours with the most instances of DUI is in the evening. Likely as people are leaving the bar or club.

C. Copy the chart into one more sheet, but now change the filter to only look at accidents where inattention was a factor from the party-at-fault. How does this distribution compare to the overall distribution?



It seems that alcohol and not paying attention have simular patterns in the number of accidents and time of day. This correlation makes a lot of sense considering the effects of alcohol

– LevelUp

Simply because an accident was such that inattention was a factor does not necessarily mean that a cell phone was the source of the driver's distraction. In the parties table, there is a column called **sp_info_2**. This feature takes on a value of B, 1, or 2 if a cell phone was known to be in use at the time of the accident.

If you're interested in digging deeper, you might want to try seeing what proportion of accidents were caused by cell phone distraction, and if they differ from other 'inattention' accidents. Keep in mind that the **sp_info_2** column is a string data type, so you'll need to treat the '1', and '2' codes appropriately!

```
SELECT
 a.day_of_week,
 COUNT(DISTINCT a.case_id) as total_accidents,
 SUM(
   CASE
      WHEN b.oaf_1 = 'F'
     OR b.oaf_2 = 'F' THEN 1
      ELSE 0
    END
  ) as inattention_accidents,
 SUM(
   CASE
      WHEN b.sp_info_2 IN ('B', '1', '2') THEN 1
      ELSE 0
    END
  ) as cell_phone_distraction_accidents,
  ROUND(
   AVG(
      CASE
        WHEN b.oaf_1 = 'F'
        OR b.oaf_2 = 'F' THEN 1
        ELSE 0
      END
    ),
    4
  ) as avg_inattention_proportion,
  ROUND(
   AVG(
      CASE
        WHEN b.sp_info_2 IN ('B', '1', '2') THEN 1
        ELSE 0
      END
    ),
    4
  ) as avg_cell_phone_distraction_proportion
```

```
FROM
  switrs.collisions AS a
  INNER JOIN switrs.parties as b ON a.case_id = b.case_id
GROUP BY
  a.day_of_week
ORDER BY
  total_accidents DESC;
```

The data shows that cell phones are the culprit for distracted driving accidents roughly half the time in each day

- Submission

Great work completing this Milestone! To submit your completed Milestone, you will need to download / export this document as a PDF and then upload it to the Milestone submission page. You can find the option to download as a PDF from the File menu in the upper-left corner of the Google Doc interface.