</talentlabs>

# Assignment Chapter 4 - Data Wrangling with SQL

## Instructions

1. This assignment is split into 2 parts. For Part 1, no dataset is required. For part 2 you will need to use the boston_crime.csv dataset that was used during the SQL demonstration lessons.
2. Please answer the questions in the boxes provided.
3. Please submit the assignment through the TalentLabs Learning System.

## Part 1: SQL Queries

**Question 1.1:**
Complete the query below to load data without duplicates.

```
SELECT
   DISTINCT *
FROM
   dataset.tableName
```

**Question 1.2:**
Write a query to select all columns from "cars.database", and all rows which have missing values in the "mileage" column.

```
SELECT
  *
FROM
  cars.database
WHERE
  mileage IS NULL;
```

</talentlabs>

**Question 1.3:**

Following on from question 1.2, write a query to replace the missing values in the mileage column with 0 for rows where the column "condition" has values equal to "new".

```
UPDATE
  cars.database
SET
  mileage = 0
WHERE
  mileage IS NULL AND condition = 'new';
```

**Question 1.4:**

Write a query to select 3 columns ("Date", "Purchase_Price", "Purchase_Desc") from the following table: shop.history. Filter the query to only include data for dates (in "Date" column) between Jan 1st 2019 and April 1st 2022. Finally, order the resulting table by the "Purchase_Price" column with the highest value first.

```
SELECT
  Date, Purchase_Price, Purchase_Desc
FROM
  shop.history
WHERE
  Date BETWEEN '2019-01-01' AND '2022-04-01'
ORDER BY
  Purchase_Price DESC;
```

</talentlabs>

# Part 2 – Data Wrangling with SQL

For part 2 of this assignment you will need to use the boston_crime.csv dataset. Make sure your data set id is boston, and the table name is crime (FROM boston.crime).

**Question 2.1:**
How many entries (rows) does this dataset contain?

The dataset contains 319,073 entries.

**Question 2.2:**
How many unique offense codes are present within the data? Use the Group By command to find your answer. In the box below, please provide your answer to the question and the query used.

Unique offense codes:
- There are 222 unique offense codes.

Query:

```
SELECT
  COUNT(DISTINCT OFFENSE_CODE)
FROM
  boston.crime;
```

**Question 2.3:**
Find out how many OFFENSE_DESCRIPTION entries contain the word "ASSAULT" as the first word?
  – e.g. ASSAULT - AGGRAVATED - BATTERY
In the box below, please provide your answer to the question and the query used.

OFFENSE_DESCRIPTION entries containing "Assault" as the first word:
- There are 23,567 entries where the OFFENSE_DESCRIPTION contains "ASSAULT" as the first word.

Query:

```
SELECT
  COUNT(*)
FROM
  boston.crime
WHERE
  OFFENSE_DESCRIPTION LIKE 'ASSAULT%';
```

</talentlabs>

**Question 2.4:**

Make a new column called TIME which contains the time of the offense from the OCCURRED_ON_DATE column. (Hint: you will need to use the CAST and SUBSTR functions together)

In the box below, please provide the query used as well as a screenshot of the query results containing the new TIME column. The column should look like the one in the Sample Screenshot below.

Query:
SELECT *,
    CAST(SUBSTR(CAST(OCCURRED_ON_DATE AS STRING), 12, 8) AS TIME) AS TIME
FROM boston.crime;

Reasoning:
- First change the OCCURRED_ON_DATE to String type, then extract it, starting from the 12th character and 8 characters after that. Then convert it back from String to Time type, and finally name it as a new column named TIME.

Screenshot: