# 6 SQL Tricks Every Data Scientist Should Know

Part 1 of SQL tricks to make your analytics work more efficient



Data scientists/analysts should know SQL, in fact, all professionals working with data and analytics should know SQL. To some extent, SQL is an under-rated skill for data science because it has been taken for granted as a necessary yet uncool way of extracting data out from the database to feed into pandas and {tidyverse} — fancier ways to wrangle your data.



**Photo Source** 

However, with massive data being collected and churned out every day in the industries, as long as the data reside in a SQL compliant database, SQL is still the most proficient tool to help you investigate, filter and aggregate to get a thorough understanding of your data. By slicing and dicing with SQL, analysts are allowed to possibly identify patterns that worth further looking into, which oftentimes leads to redefining the analysis population and variables to be considerably smaller (than the initial scope).

Hence rather than transferring huge datasets into Python or R, the first step of analytics should be using SQL to gain informative insights from our data.

Working in real-world relational databases, SQL is way more than just SELECT, JOIN, ORDER BY statements. In this blog, I will discuss 6 tips (and one Bonus tip) to make your analytics work more efficient with SQL and its integrating with other programming languages like Python and R.

For this exercise, we will work with Oracle SQL on the toy data table below, which consists of multiple types of data elements,

ID_VAR	SEQ_VAR	EMPTY_STR_VAR	NULL_VAR	NA_STR_VAR	NUM_VAR	DATE_VAR1	DATE_VAR2
ID, used as the	sequence	string (missing value	string (missing value	string (missing value	float number	date1	date2
join key	number	coded as empty)	coded as NULL)	coded as 'NA')			
19017	1		(null)	NA	198.2	11/2/2018	11/30/2018
19017	2		(null)	NA	212.35	11/2/2018	11/30/2018
19017	3		(null)	NA	424.99	11/2/2018	11/30/2018
19017	4		(null)	NA	318.53	11/2/2018	11/30/2018
19017	5		(null)	NA	302.44	11/2/2018	11/30/2018
19064	1		(null)	NA	135.13	1/28/2019	1/28/2019
19064	2		(null)	NA	135.13	1/30/2019	1/30/2019
19064	3		(null)	NA	135.13	1/26/2019	1/26/2019
19064	4		(null)	NA	83.05	1/26/2019	1/26/2019
19064	5		(null)	NA	99.66	1/31/2019	1/31/2019
19228	1	S	(null)	NA	62.71	3/25/2019	3/25/2019
19228	2	S	(null)	NA	62.71	3/27/2019	3/27/2019
19228	3	S	(null)	NA	62.74	3/29/2019	3/29/2019
19272	1		(null)	NA	22226.19	2/24/2019	3/22/2019

Toy data table (with variable definitions)

## 1. COALESCE() to recode NULL / missing data

When it comes to re-coding missing values, the COALESCE() function is our secret sauce, which, under this circumstance, re-codes the NULL to whatever value specified in the second argument. For our example, we can re-code the NULL\_VAR to a character value 'MISSING',

```
1 ---- 1) COALESCE() to recode the NULL value to the character string MISSING

2 SELECT

3 ID_VAR,

4 NULL_VAR,

5 COALESCE(NULL_VAR, 'MISSING') AS RECODE_NULL_VAR

6 FROM

7 CURRENT_TABLE

8 ORDER BY ID_VAR

coalesce.sql hosted with ♥ by GitHub view raw
```

ID_VAR	NULL_VAR	COALESCE_NULL_VAR
19017	(null)	MISSING
19064	(null)	MISSING
19228	(null)	MISSING
19228	(null)	MISSING
19228	(null)	MISSING
19272	(null)	MISSING

COALESCE() to recode NULL

One important note, however, is that in databases, *missing values* can be encoded in various ways besides NULL. For instance, they could be empty string/blank space (e.g., EMPTY\_STR\_VAR in our table), or a character string 'NA' (e.g., NA\_STR\_VAR in our table). In these cases, COALESCE() would not work, but they can be handled with the CASE WHEN statement,

```
--- However, COALESCE() NOT WORK for Empty or NA string, instead, use CASE WHEN
    SELECT
 3
      ID_VAR,
     EMPTY_STR_VAR,
       COALESCE(EMPTY_STR_VAR, 'MISSING') AS COALESCE_EMPTY_STR_VAR,
       CASE WHEN EMPTY_STR_VAR = ' ' THEN 'EMPTY_MISSING' END AS CASEWHEN_EMPTY_STR_VAR,
 8
       NA_STR_VAR,
       CASE WHEN NA_STR_VAR = 'NA' THEN 'NA_MISSING' END AS CASEWHEN_NA_STR_VAR
 9
10
    FROM
11
       CURRENT_TABLE
     ORDER BY ID_VAR
coalesce casewhen.sql hosted with \bigcirc by GitHub
                                                                                             view raw
```

CASE WHEN to re-code empty or NA

ID_VAR	EMPTY_STR_VAR	COALESCE_EMPTY_STR_VAR	CASEWHEN_EMPTY_STR_VAR	NA_STR_VAR	CASEWHEN_NA_STR_VAR
19228	S	S		NA	NA_MISSING
19228	S	S		NA	NA_MISSING
19228	S	S		NA	NA_MISSING
19272			EMPTY_MISSING	NA	NA_MISSING

### 2. Compute running total and cumulative frequency

Running total can be useful when we are interested in the total sum (but not individual value) at a given point for potential analysis population segmentation and outlier identification.

The following showcases how to calculate the running total and cumulative frequency for the variable NUM\_VAR,

```
--- 2) Running total/frequency
 2 SELECT
    DAT.NUM_VAR,
     SUM(NUM_VAR) OVER (PARTITION BY JOIN_ID) AS TOTAL_SUM,
     ROUND(CUM_SUM / SUM(NUM_VAR) OVER (PARTITION BY JOIN_ID), 4) AS CUM_FREQ
 6 FROM
 7
     SELECT
 8
             T.*,
             SUM(NUM_VAR) OVER (ORDER BY NUM_VAR ROWS UNBOUNDED PRECEDING) AS CUM_SUM,
10
              CASE WHEN ID_VAR IS NOT NULL THEN '1' END AS JOIN_ID
     FROM CURRENT_TABLE
13 ) DAT
14 ORDER BY CUM_FREQ
running total.sgl hosted with ♥ by GitHub
                                                                                         view raw
```

NUM_VAR	TOTAL_SUM	CUM_FREQ
62.71	24458.96	0.0026
62.71	24458.96	0.0051
62.74	24458.96	0.0077
83.05	24458.96	0.0111
99.66	24458.96	0.0152
135.13	24458.96	0.0207
135.13	24458.96	0.0262
135.13	24458.96	0.0317
198.2	24458.96	0.0398
212.35	24458.96	0.0485
302.44	24458.96	0.0609
318.53	24458.96	0.0739
424.99	24458.96	0.0913
22226.19	24458.96	1

Output for cumulative frequency

Here is our output (on the left).

Two tricks here, (1) SUM over <u>ROWS UNBOUNDED PRECEDING</u> will calculate the sum of all prior values to this point; (2) create a JOIN\_ID to calculate the total sum.

We use the <u>window function</u> for this calculation, and from the cumulative frequency, it is not hard to spot the last record as an outlier.

#### 3. Find the record(s) with extreme values without self joining

So our task is to return the row(s) with the largest NUM\_VAR value for each unique ID. An intuitive query is to first find the max value for each ID using group by, and then self join on ID and the max value. Yet a more concise way would be,

```
--- 3) Find the record having a number calculated by analytic functions (e.g., MAX) without sel
   SELECT *
   FROM
4 (
5
    SELECT
6
      DAT.*,
       CASE WHEN (NUM_VAR = MAX(NUM_VAR) OVER (PARTITION BY ID_VAR)) THEN 'Y' ELSE 'N' END AS MAX
7
8
    FROM
     CURRENT_TABLE
9
                      DAT
10 ) DAT2
11 WHERE MAX_NUM_IND = 'Y'
```

NECOTAS VVICITARIO ITIAN VALUE

this query should give us the following output, showing rows having the max NUM\_VAR grouped by ID,

ID_VAR	SEQ_VAR	EMPTY_STR_VAR	NULL_VAR	NA_STR_VAR	NUM_VAR	DATE_VAR1	DATE_VAR2	MAX_NUM_IND
19017	3		(null)	NA	424.99	11/2/2018	11/30/2018	Υ
19064	1		(null)	NA	135.13	1/28/2019	1/28/2019	Υ
19064	2		(null)	NA	135.13	1/30/2019	1/30/2019	Υ
19064	3		(null)	NA	135.13	1/26/2019	1/26/2019	Υ
19228	3	S	(null)	NA	62.74	3/29/2019	3/29/2019	Υ
19272	1		(null)	NA	22226.19	2/24/2019	3/22/2019	Υ

Output for records with the max NUM\_VAR value

#### 4. Conditional WHERE clause

Everyone knows the WHERE clause in SQL for subsetting. In fact, I find myself using conditional WHERE clause more often. With the toy table, for instance, we want only

to keep the rows satisfying the following logic,

- if SEQ\_VAR in (1, 2, 3) & diff(DATE\_VAR2, DATE\_VAR1)  $\geq 0$
- elif SEQ\_VAR in (4, 5, 6) & diff(DATE\_VAR2, DATE\_VAR1)  $\geq 1$
- else diff(DATE\_VAR2, DATE\_VAR1)  $\geq$ 2

Now the conditional WHERE clause comes in handy,

```
1 -- 4) Conditional where clause
 2 SELECT
 3 DAT.ID_VAR,
 4 DAT.SEQ_VAR,
    DAT.NUM_VAR,
    DATE_VAR1,
     DATE_VAR2,
     TRUNC(DATE_VAR2) - TRUNC(DATE_VAR1) AS LAG_IN_DATES
 9 FROM
10 CURRENT_TABLE DAT
11 WHERE
    (TRUNC(DATE_VAR2) - TRUNC(DATE_VAR1)) >= CASE WHEN SEQ_VAR IN (1,2,3) THEN 0 WHEN SEQ_VAR IN
12
    ORDER BY ID_VAR, SEQ_VAR
13
```

SEQ_VAR	NUM_VAR	DATE_VAR1	DATE_VAR2	LAG_IN_DATES
1	198.2	11/2/2018	11/30/2018	28
2	212.35	11/2/2018	11/30/2018	28
3	424.99	11/2/2018	11/30/2018	28
4	318.53	11/2/2018	11/30/2018	28
5	302.44	11/2/2018	11/30/2018	28
1	135.13	1/28/2019	1/28/2019	0
2	135.13	1/30/2019	1/30/2019	0
3	135.13	1/26/2019	1/26/2019	0
1	62.71	3/25/2019	3/25/2019	0
2	62.71	3/27/2019	3/27/2019	0
3	62.74	3/29/2019	3/29/2019	0
1	22226.19	2/24/2019	3/22/2019	26
	1 2 3 4 5 1 2 3 1 2	2 212.35 3 424.99 4 318.53 5 302.44 1 135.13 2 135.13 3 135.13 1 62.71 2 62.71 3 62.74	1 198.2 11/2/2018 2 212.35 11/2/2018 3 424.99 11/2/2018 4 318.53 11/2/2018 5 302.44 11/2/2018 1 135.13 1/28/2019 2 135.13 1/30/2019 3 135.13 1/26/2019 1 62.71 3/25/2019 2 62.71 3/27/2019 3 62.74 3/29/2019	1       198.2       11/2/2018       11/30/2018         2       212.35       11/2/2018       11/30/2018         3       424.99       11/2/2018       11/30/2018         4       318.53       11/2/2018       11/30/2018         5       302.44       11/2/2018       11/30/2018         1       135.13       1/28/2019       1/28/2019         2       135.13       1/30/2019       1/30/2019         3       135.13       1/26/2019       1/26/2019         1       62.71       3/25/2019       3/25/2019         2       62.71       3/27/2019       3/27/2019         3       62.74       3/29/2019       3/29/2019

Output for conditional where clause

The logic aforementioned should eliminate the sequences 4, 5 of ID = 19064 because the difference between date2 and date1 = 0, and this is exactly what the query returns

above.

#### 5. Lag() and Lead() to work with consecutive rows

Lag (looking at the previous row) and Lead (looking at the next row) probably are two of the most used <u>analytic functions</u> in my day-to-day work. In a nutshell, these two functions allow users to query more than one row at a time without self-joining. More detailed explanations can be found <u>here</u>.

Let's say, we want to compute the difference in NUM\_VAR between two consecutive rows (sorted by sequences),

```
--- 5) LAG() or LEAD() function
     SELECT
     DAT.ID_VAR,
     DAT.SEQ_VAR,
 5
     DAT.NUM_VAR,
      NUM_VAR - PREV_NUM AS NUM_DIFF
 6
     FROM
 8
     (
 9
             SELECT
10
              T.*,
              LAG(NUM_VAR, 1, 0) OVER (PARTITION BY ID_VAR ORDER BY SEQ_VAR) AS PREV_NUM
11
             FROM
12
              CURRENT_TABLE
13
     ) DAT
14
     ORDER BY ID_VAR, SEQ_VAR
lan sal hosted with C by GitHub
```

The LAG() function returns the prior row, and if there is none (i.e., the first row of each ID), the PREV\_NUM is coded as 0 to compute the difference shown as NUM\_DIFF below,

ID_VAR	SEQ_VAR	NUM_VAR	PREV_NUM	NUM_DIFF
19017	1	198.2	0	198.2
19017	2	212.35	198.2	14.15
19017	3	424.99	212.35	212.64
19017	4	318.53	424.99	-106.46
19017	5	302.44	318.53	-16.09
19064	1	135.13	0	135.13
19064	2	135.13	135.13	0
19064	3	135.13	135.13	0
19064	4	83.05	135.13	-52.08
19064	5	99.66	83.05	16.61
19228	1	62.71	n	62.71

	-	02.71		02171
19228	2	62.71	62.71	0
19228	3	62.74	62.71	0.03
19272	1	22226.1	9 0	22226.19

Output from LAG()

#### 6. Integrate SQL query with Python and R

The prerequisite of integrating SQL queries into Python and R is to establish the database connections via ODBC or JDBC. Since this is beyond the scope of this blog, I will not discuss it here, however, more details regarding how to (create ODBC or JDBC connections) can be found <u>here</u>.

Now, assuming that we already connected Python and R to our database, the most straightforward way of using query in, say Python, is to copy and paste it as a string, then call pandas.read\_sql(),

```
my_query = "SELECT * FROM CURRENT_TABLE"
sql_data = pandas.read_sql(my_query, connection)
```

Well, as long as our queries are short and finalized with no further changes, this method works well. However, what if our query has 1000 lines, or we need to constantly update it? For these scenarios, we would want to read .*sql* files directly into Python or R. The following demonstrates how to implement a getSQL function in Python, and the idea is the same in R,

```
import pandas as pd
    def getSQL(sql_query,
                place_holder_str,
                replace place holder with,
               database con):
        ...
        Args:
             sql query: sql query file
             place_holder_str: string in the original sql query that is to be replaced
10
             replace place holder with: real values that should be put in
             database con: connection to the database
11
         sqlFile = open(sql query, 'r')
         sqlQuery = sqlFile.read()
16
```

```
sqlQuery = sqlQuery.replace(place_holder_str, replace_place_holder_with)

df = pd.read_sql_query(sqlQuery, database_con)

database_con.close()

return df
```

Here, the first arg sql\_query takes in a separate standalone .*sql* file that can be easily maintained, like this,

```
1 SELECT
2 *
3 FROM
4 CURRENT_TABLE DAT
5 WHERE
6 ID_VAR IN ('ID_LIST')
7 ORDER BY ID_VAR, SEQ_VAR

SQL_FILE.sql hosted with ♥ by GitHub view raw
```

The "ID\_LIST" is a placeholder string for the values we are about to put in, and the getSQL() can be called using the following code,

```
1 seq12_df = getSQL('SQL_FILE.sql', 'ID_LIST', "','".join(['19228', '19272']), database_con=conn)
```

## Bonus tip, regular expression in SQL

Even though I do not use regular expression in SQL all the time, it sometimes can be convenient for text extraction. For instance, the following code shows a simple example of how to use REGEXP\_INSTR() to find and extract numbers (see <a href="here">here</a> for more details),

```
-- Find and extract numbers between 0 - 9 that consecutively happens 5 times

SELECT

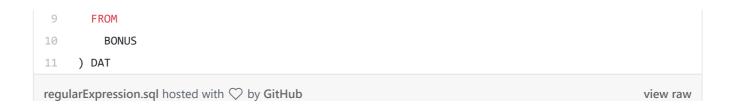
SUBSTRING(LONG_TEXT, REG_IDX, REG_IDX+5) AS NUMBER_LIST_FOUND

FROM

(
SELECT

REGEXP_INSTR(LONG_TEXT, '[0-9]{5}') AS REG_IDX,

LONG_TEXT
```



I hope you find this blog helpful, and the full code along with the toy dataset is available in my github.

P.S. head over to Part 2 of this mini-series for more SQL analytics tips.

