
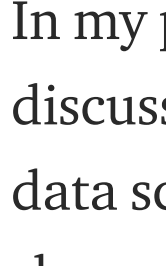


Extra SQL Tricks Every Data Scientist Should Know

Part 2 of getting more out of SQL to step up your analytics

Yi Li [Follow](#)

Apr 20 · 6 min read ★



prev

vious b

used the most us



Today, we will be working with a new toy table shown below containing multiple data types, and for demonstration purpose, this exercise will be implemented in MS SQL server 2017,

ID	Var	Num	Var	Gen	Var	Date	Var
----	-----	-----	-----	-----	-----	------	-----

0155	60.79	Female	6/2/2018
0155	79.77	Female	9/29/2018
0155	122.82	Female	9/30/2018
0180	1810.47	Female	6/30/2016
0188	2732.62	Female	6/30/2016

0792	886.17	Male	8/31/2016
0792	954.71	Male	9/30/2016
0792	1048.69	Male	10/31/2016

Toy data table (with variable definitions)

ROW_NUMBER() to return a subset of rows

The ROW_NUMBER() function in SQL creates a unique incremental integer value to each row of the result. This column of values is considered a **pseudo-column** as it does not inherently exist in our data table, and because of which, the result is returned in the order determined by the analysts in the ORDER BY clause.

With this pseudo-column, we can solve the “top N” questions using

The following demonstrates how

column RowNumber from the output),

RowNumber	ID_Var	Num_Var
1	0155	122.82
2	0155	79.77
3	0155	60.79
1	0180	1810.47

	1	2	3
1	0792	1048.69	
2	0792	954.71	
3	0792	886.17	

Output: top 3 records with the highest Num_Var value

being only one row, ID C
IDs (with more than 3 rows)
ard to whether the raw data
ified. Straightforward!

**h of consecutive days w
vanced)**

of the ROW_NUMBER()

possible application in time series, where

What do I mean? Well, let's consider our toy data being customers logging in to our website. Customer ID 0155 logged in on 06/01, 06/02, and 06/03 (i.e., 3 consecutive days), and then there were another 2 consecutive logins on 09/29 and 09/30 each. Visually it is not that simple for us to do this calculation if the dates are not sorted

An intuitive way would be first using the ROW_NUMBER() function to assign integer row numbers ordered by dates within each customer ID. Then, the integer row numbers happen to be incremental by 1 also. The difference between consecutive days is just the difference between the last row number and the

This brings up an important note on the ROW_NUMBER() function:

09/29 are not consecutive). Thus, simply applying ROW_NUMBER() cannot get us the desired output. We need some tweaks, and here's how,

Step 1: Create the RowNumber (order by date ASC)

- RowNumber to group the consecutive observations (the starting points for counting)
- RowNumber to group the consecutive observations (the starting points for counting)

turn, which is labeled as


h_grouping_var,

Interim output: Creating a grouping variable

Starting_Count_DT because it's meaning is "starting count date" variable to basically tell SQL that these rows are incremental by 1).

Count_DT

our final output!



Output: length of consecutive days

To put everything together,

Phew, problem solved! If you don't get it the first time, no worries.

exact problem in your next data scientist position interview with
big tech companies 🙏

The previous query, as w

... by leveraging the WITH clause, we broke down the n into 2 individual temp views for better code readability. This is a good practice for the relations among the multiple WITH queries.

independent or dependent

referenced multiple times. Without repeatedly writing the same sub-query again and again (i.e., the DRY / Don't Repeat Yourself principle), we can introduce upfront and reuse it later.

4. Concatenating to re-format your data structure

Moving on, our next task is to create an aggregated report for our

observing that the ID_Var and Gen_Var both contain duplicated values across rows, therefore to make it suitable for human reading, we immediately think of the STRING_AGG() function in MS SQL 2017,

this code snippet yields,

ed by simply specifying the ORDER BY statement (e.g., ID ORDER BY).

Despite of being more commonly used for string concatenation, the **STRING_AGG()** function also works for other data types (e.g. date and number variables in our example).

As a couple of *callouts* regarding this aggregation functionality:

- Although it is available in major SQL databases, the specific function name varies. In Oracle/PLSQL (11g), it's the **LISTAGG()** function. In Microsoft SQL Server, it's the **STRING_AGG()** function. In MySQL along with IBM Netezza, it's the **GROUP_CONCAT()**, in PostgreSQL server 2017 and PostgreSQL, it's the **STRING_AGG()**.

analysis (by the way, check out [my previous post for tips of extracting data from SQL database to Python](#) if you haven't done so). Instead of transferring the entire raw data table as is with most columns having

Python. In this concatenation, Null values will be excluded, and thus we don't have to worry about missing value handling at this point.

As always, all the code snippets together with the toy data are available [here in my Github](#) 😊


SQL tricks to make your analytics work more efficient

towardsdatascience.com

	(row)	N/A
	(row)	N/A
	(row)	N/A
	(row)	N/A
	(row)	N/A
	(row)	N/A
	(row)	N/A
5	(row)	N/A
5	(row)	N/A
5	(row)	N/A
5	(row)	N/A

Sign up for The Daily Pick

delivered Monday to Thursday. Make learning your daily ritual. [Take a look](#)

 [Get this newsletter](#)

Emails will be sent to sreeaurovindh@gmail.com.
Not you?

[Sql](#) [Data Science](#) [Analytics](#) [Python](#)



WRITTEN BY

Yi Li

Data Scientist

Follow

Discover Medium
Welcome to a place where

Make Medium yours
Follow all the topics you care

Explore your membership

ads in sight. Watch

access to insightful stories from amazing thinkers and storytellers. [Browse](#)

Medium

[About](#) [Help](#) [Legal](#)