Database Designs

Section 4: Subqueries, Constraints & Advanced Concepts

4.1 Subqueries Introductor

So far we have explored how to query from tables and virtual table to generate a results table. SQL Subqueries is a fairly advanced concept which allows us to query from a SQL query i.e. a query within a query.

So far we have seen queries as standalone commands that fetch data from a database; however, in reality, queries are generally plug-and-play - what do we mean by this? Plug-and-play means the ability to use queries in places where you would not expect them to be used, this is because the results of queries are tables.

Tables can be real tables, tables that are generated by joins or tables that are a result of queries. Therefore, queries are plug-and-play into other pieces of SQL.

A query is a command that returns a table (columns and rows). If you imagine the result of a query as a table by itself, then this mentality will open all the possibilities that you can do with the results of a query. For example:

We could calculate the Union, Intersection or Differences of two queries.

We could use one query inside another (via Subqueries).

We could use a subquery to populate a table via an **INSERT**.

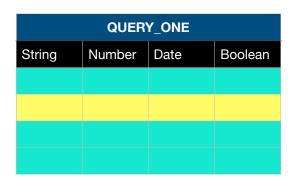
To conclude, a subquery is nothing more than a query of which the results are used in another query, otherwise there is no logical differences between a subquery and a normal standalone query.

4.2 Union, Union All, Intersect and Except

We can calculate the Union, Intersection and Difference between two queries provided they have the same columns i.e. the number, order and type of the columns are identical across the two queries.

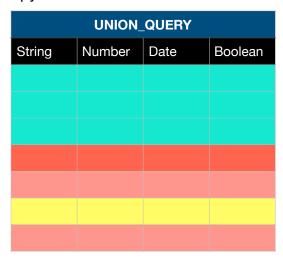
To demonstrate the **UNION**, **UNION ALL**, **INTERSECT** and **EXCEPT** set operators subqueries, below are two tables QUERY_ONE and QUERY_TWO. The QUERY_ONE table row data is coloured turquoise while the QUERY_TWO table row data is coloured peach. The common data rows of both tables are coloured in yellow.

This will help clearly illustrate how SQL creates the subquery tables using each Set operators and how we can use subqueries to perform more advanced queries and open our minds to all the possibilities SQL offers using the plug-and-play mentality.

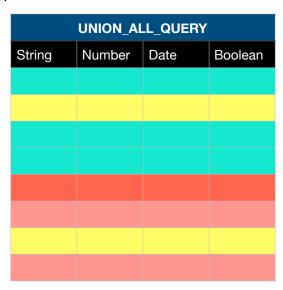


QUERY_TWO			
String	Number	Date	Boolean

The **UNION** set operator creates a new query table which combines the two queries together but creates one copy of the common row i.e. removes duplicate rows.



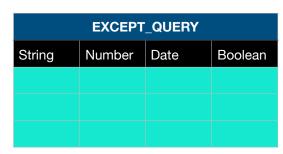
The **UNION ALL** set operator creates a new query table which combines the two queries together containing all duplicate rows.



The **INTERSECT** set operator finds the commonality between the two sets that it is intersecting i.e it returns the common row between the two query tables.

INTERSECT_QUERY			
String	Number	Date	Boolean

Finally, the **EXCEPT** set operator creates a new query table which removes everything in the first table (i.e. the left table) that is present in the second table (i.e. the right table). This returns the difference between the two sets of input queries.





Important Note: The Except query can return either of the above results depending on which table was made as the left table in the subquery. The left table is always the first table mentioned in the query within the **FROM** clause.

The **UNION**, **UNION ALL**, **INTERSECT** and **EXCEPT** are all a kind of Set operations. We know that a query is a command that returns a table and not really a Set. The Entity Relationship theory tells us that a table is a collection of tuples (a tuple relates to one row within the table). A Set cannot contain duplicates but a table can.

By default the **UNION** will eliminate duplicates but SQL has the Union All operator to keep all the duplicate tuples if we seek that operation behaviour. There is another rule which applies to the **UNION** Set operator which is where individual queries that participate in a **UNION** operator cannot have the **ORDER BY** clause. This is because the elements of a Set are not ordered. Below is an example demonstrating this:

Pet_Owners		
AptNumber	Name	
123	John	
345	Tim	
349	Vandana	
567	Bilal	

Flat_Owners		
FlatNumber	Name	
234	Mary	
567	Bilal	
879	Jane	
903	Ellen	

In this example we have two tables of Pet_Owners and Flat_Owners both have the same columns i.e. they have in common the same number, order and type of columns. The names of the columns is irrelevant. In the above example, both tables have two columns of integer and string in that specific order. Below is the syntax to perform the **UNION** Set operator on the two tables.

(SELECT AptNumber AS FlatNumber, Name FROM Pet_Owners)
UNION

(SELECT FlatNumber, Name FROM Flat_Owners);

It is OK for the column names in the physical tables to be different. This is because in a query we can easily alias one column name to another - notice the AptNumber AS FlatNumber which simply glosses/hides the fact that the original column name was

something different. This is the reason for why the actual name from the physical table is irrelevant. The final subquery table will return a results table with two columns called FlatNumber and Name as seen below:

UNION_RESULTS_TABLE		
FlatNumber	Name	
123	John	
345	Tim	
349	Vandana	
234	Mary	
567	Bilal	
879	Jane	
903	Ellen	

As previously mentioned, individual queries that have been ordered by the **ORDER BY** clause cannot be used with the UNION set operator. Therefore, the below example syntax is invalid and will not work returning an error:

(SELECT AptNumber AS FlatNumber, Name FROM Pet_Owners ORDER BY Name)
UNION

(SELECT FlatNumber, Name FROM Flat Owners ORDER BY Name);

However, we can write the syntax so that the results of the **UNION** is ordered by the column Name. In this case the **ORDER BY** clause will work on ordering the **UNION** results table:

```
(SELECT AptNumber AS FlatNumber, Name FROM Pet_Owners)
UNION
(SELECT FlatNumber, Name FROM Flat_Owners)
) ORDER BY Name;
```

Pay particular attention to the curly brackets which wraps around the whole **UNION** set operator subquery. The **ORDER BY** clause is applied to the results of the **UNION** because it is falls outside the wrapping curly brackets. Therefore, it is only the individually queries participating in the **UNION** that cannot have the **ORDER BY** clause and the reason for why the second syntax works without throwing any errors.

4.3 Query-In-A-Query

Subqueries are very useful because they allow us to write SQL queries entirely free of hardcoded values or very large intermediary tables. Subqueries at first can seem very difficult to get use to at the beginning but once you know how to use them they are

extremely powerful. We will now explore how to use a query within a query (also known as subqueries) using the tables below.

Stores_Data		
storeID	storeLocation	city
1	ASDA Wilmslow	Manchester
2	ASDA Stratford	London

Products_Data		
productID	productName	
1	Bread	
2	Milk	
3	Noodles	
4	Nutella	

Sales_Data			
storeID	productID	salesDate	totalRevenue
1	1	November 20, 2020	7,233.32
1	3	November 20, 2020	3,234.84
1	2	November 20, 2020	5,865.55
1	2	November 20, 2020	6,849.99
2	3	November 20, 2020	2,110.95
2	2	November 20, 2020	4,558.24
2	4	November 20, 2020	2,284.75

If a database existed like the above it is most likely that the business would like to pull a report to see what are the annual revenue is like for various products and make business decisions around this information. The query to extract this report would look like the following:

```
SELECT p.productName, YEAR(date), SUM(totalRevenue)
FROM Sales_Data AS s
INNER JOIN Products_Data AS p
ON s.productID = r.productID
WHERE (p.productName = 'Bread' or p.productName = 'Milk')
AND (YEAR(date) = 2020)
GROUP BY p.productName, YEAR(date);
```

In the above syntax we are matching on the productID column and combining the Products_Data table onto the Sales_Data table using an INNER JOIN. The WHERE clause shows that we are interested in products 'Bread' and 'Milk' which are believed to be the top sellers and we are interested in the current year. Finally, we wan to GROUP BY the productName and the year. In conclusion this query will return back using the SELECT statement the productName, year and the SUM of the totalRevenue for each product within the current year groups.

The above query will only return information for the year 2020 and about 'Bread' and 'Milk' because that is the date and products specified in the **WHERE** clause. This query would get the job done to extract the report from the database for the business. However, there are many issues with the above query which relate to hardcoding the values. Our objective is to re-write the code to avoid any hardcoded values.

First, the query is hardcoded for the year 2020 which means the query would need to be updated every year. To avoid this hardcoding of values we can plug one query into another using subqueries as seen below.

```
SELECT p.productName, YEAR(date), SUM(totalRevenue)

FROM Sales_Data AS s

INNER JOIN Products_Data AS p

ON s.productID = r.productID

WHERE (p.productName = 'Bread' or p.productName = 'Milk')

AND (YEAR(date) = (SELECT YEAR(MAX(date)) FROM Sales_Data))

GROUP BY p.productName, YEAR(date);
```

Here we are using one query inside another larger query. The inner query (highlighted in yellow) returns the maximum value of the year from the Sales_Data table. This query will look for the last transaction from the Sales_Data table because this would be the maximum (latest) date and return the year value from the date. Therefore, the inner query returns a single value i.e. the year in which the last transaction occurred. In the **WHERE** clause of the outer query it simply compares the year of the date that is returned from the inner query. The **WHERE** clause functions as it did before.

Note that the the inner queries are always evaluated first.

The next issue relates to the hardcoded products value within the **WHERE** clause. What if we do not want to hardcode these values. Let's imagine that we instead have a TopSellers_Data table which has the productID and productName columns. This table holds the top products we are interested in at any point in time.

TopSellers_Data		
productID	productName	
1	Bread	
2	Milk	

Instead of hardcoding the products in our query we instead want to retrieve the products from the TopSellers Data table. Again we can do this using a subquery as seen below:

```
SELECT p.productName, YEAR(date), SUM(totalRevenue)
FROM Sales_Data AS s
INNER JOIN Products_Data AS p
ON s.productID = r.productID
WHERE (p.productName IN (SELECT p.productName FROM TopSellers_Data)
AND (YEAR(date) = (SELECT YEAR(MAX(date)) FROM Sales_Data))
GROUP BY p.productName, YEAR(date);
```

Once again we have plugged one query into another (the inner query is highlighted in yellow). The inner **SELECT** statement simply returns all the products from the TopSellers_Data table. This TopSellers_Data table can be modified/updated and the query will always return the latest top products. The inner query returns a range of values which the outer query can then use these values using the **IN** keyword for its **WHERE** clause check for the current top products.

The final issue with the original query is to do with the size of our table. Let us imagine as time passes the business is operating extremely well. This means more sales will be recorded in the Sales_Data table and the possibility that the Sales_Data table eventually becomes enormous. We now have an issue that the table is too large to perform the INNER JOIN on the Sales_Data table which could take too long to generate the report. What we wan to do is reduce the number of rows that we select from the Sales_Data table. We know that at any point in time we are only interested in a few products which are present in the TopSellers_Data table. So How can we reduce the number of rows in the Sales_Data table that should be part of the query? Once again we would use subqueries as seen below:

```
SELECT p.productName, YEAR(date), SUM(totalRevenue)

FROM (SELECT * FROM Sales_Data WHERE productID IN (SELECT productID FROM TopSellers_Data)) AS s

INNER JOIN Products_Data AS p
ON s.productID = r.productID

WHERE (p.productName IN (SELECT p.productName FROM TopSellers_Data)
AND (YEAR(date) = (SELECT YEAR(MAX(date)) FROM Sales_Data))

GROUP BY p.productName, YEAR(date);
```

Instead of only using the Sales_Data table in the **INNER JOIN**, we can use a **SELECT** statement which returns a table to be one of the tables in the **INNER JOIN**. Remember the result from a **SELECT** statement is also a virtual table which can therefore be used as part of any of the joins i.e. wherever we use a table, we can use a query. Once again we have used one query plugged into another query to improve the performance of our query. We are now using a subset of the Sales_Data table which helps improve the performance of our query. The outer query uses the subset table (which is also a table) in the **FROM** clause to do the **INNER JOIN**.

We now have a final subquery that looks like the below:

```
SELECT p.productName, YEAR(date), SUM(totalRevenue)

FROM (SELECT * FROM Sales_Data WHERE productID IN (SELECT productID FROM TopSellers_Data)) AS s

INNER JOIN Products_Data AS p

ON s.productID = r.productID

WHERE (p.productName IN (SELECT p.productName FROM TopSellers_Data)

AND (YEAR(date) = (SELECT YEAR(MAX(date)) FROM Sales_Data))

GROUP BY p.productName, YEAR(date);
```

We can see that this query is entirely free from any hardcoded values or very large intermediate tables all thanks to the use of subqueries rather than the tables itself.

To conclude, subqueries allows us to use non-hardcoded values by returning values from subqueries as well as optimising the performance of our queries by using subset tables created by subqueries. As we can see subqueries can be difficult to grasp at the beginning but once we understand how to harness them they become very powerful and useful tools when querying databases.