

Querying Data With SQL

Introduction to Subqueries





Discussion:

Cost of Returns

The VP of Operations at Superstore needs to determine the operational cost of shipping wrong items to customers. To help her make this critical decision, you've been asked to find out the following:

What is the number of orders (with wrong items) that have been returned for each shipping mode?

Writing queries is all about asking the right questions:

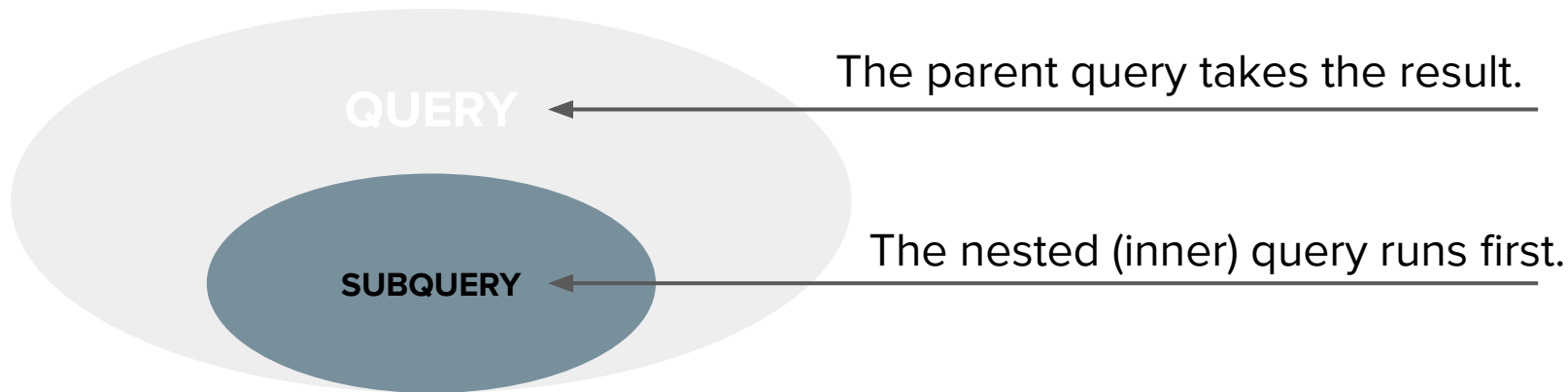


What are two questions we need to answer in order to write this query?

Meet Subqueries

Also known as a **nested query**, a subquery is an SQL statement that combines multiple questions into one query.

How this works: When you run one query and get a result, you can feed it to another query.



Pseudo-Coding Subqueries

Referring back to the request from the VP of Operations, we need to:

- Identify the order numbers that have been returned as “wrong item.”
- Count the number of returned orders by shipping mode.

Table: Returns

order_id [PK] text	return_date timestamp without time	return_quantity integer	reason_returned text	product_id text
AU-2020-1320434	2020-01-09 00:00:00	1	Wrong Item	OFF-IBI-10003541
AU-2020-1104334	2020-02-03 00:00:00	1	Wrong Item	OFF-IBI-10003541
AU-2020-1348776	2020-01-28 00:00:00	1	Wrong Item	OFF-IBI-10003541

Table: Orders (select columns)

order_id [PK] text	product_id [PK] text	order_date timestamp without time z	ship_mode text	sales numeric
AU-2020-1104334	OFF-IBI-10003541	2020-01-08 00:00:00	Second Class	719.88
AU-2020-1320434	OFF-IBI-10003541	2020-01-06 00:00:00	Second Class	719.88
AU-2020-1348776	OFF-IBI-10003541	2020-01-03 00:00:00	Second Class	719.88



Guided Walk-Through: Subquery Syntax

```
SELECT ship_mode,  
COUNT( DISTINCT order_id )  
FROM Orders  
WHERE order_id IN (  
  
    SELECT order_id  
    FROM returns  
    WHERE reason_returned =  
        'Wrong Item'  
  
    )  
GROUP BY ship_mode;
```

A breakdown of the subquery (in red):

- **SELECT** all of the order_ids.
- **FROM** the Returns table, because we want to identify orders that match our criteria for return reason.
- Filter the results for only order_ids with reason_returned equal to “Wrong Item.”



Guided Walk-Through: Subquery Syntax

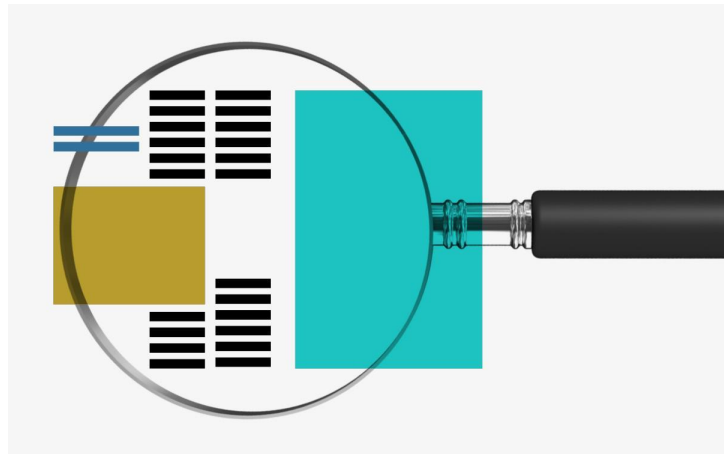
```
SELECT ship_mode,  
COUNT( DISTINCT order_id )  
FROM Orders  
WHERE order_id IN (  
  
    SELECT order_id  
    FROM returns  
    WHERE reason_returned =  
        'Wrong Item'  
  
    )  
GROUP BY ship_mode;
```

A breakdown of the outer query (in red) that reads the results from the inner query:

- **SELECT** takes ship_mode and a count of all rows.
- **FROM** tells us where the information is coming from.
- **WHERE** filters the selected rows from the Orders table using the results of our subquery.
- **GROUP BY** ship_mode because of our aggregation function COUNT.

Subquery Syntax | A Checklist

- ❑ Subqueries are enclosed in parentheses and will execute first.
- ❑ In order to execute, subquery structures **must** have complete query components: **SELECT**, **FROM**, and a specified criteria.
- ❑ Assigning aliases is important for syntax and readability.



Why Do You Need Subqueries...

...when you already have JOINS?

Subqueries and JOINS are both used to combine data from different tables.

Subqueries:

- Can be used to return either a single value or a row set.
- Can run on their own (as queries).
- Results are used immediately.
- Used in SELECT, WHERE, FROM, and HAVING.

JOINS:

- Used to return rows.
- Cannot run on their own.
- The “JOINED” results are available in SELECT statements.
- Used as FROM clauses of WHERE statements.

— Subquery Types

Types of Subqueries

Three common types of subqueries differ by how the results are used by the outer query.

Derived Tables: Creates a “temporary table” to use the results of one query as a table for the *FROM statement* of the outer query.

Single Column: Uses a subquery to filter the results of one query by the results of another query using the *IN operator*.

Nested: Uses comparison operators to filter the results of one query by the results of another query using a *comparison operator*.

Derived Table Subquery Syntax

- Also known as “a query *within* a query”, it allows us to supply an entire query between parentheses as a subquery instead of calling a normal table.
- The outer query then treats result of the inner query as if it were a table. This creates a derived table and requires an alias.

```
SELECT a_column  
FROM (SELECT another_column as a_column  
      FROM a_different_table) AS alias;
```



The Head of Sales asked for insights on whether we've seen an increase in average monthly sales over time. First, let's write a subquery to *calculate* the sales for every month:

```
SELECT DATE_TRUNC('month',order_date) as order_month,  
SUM(sales) AS monthly_sales  
FROM orders  
GROUP BY 1;
```



Then, average the results in the outer query:

```
SELECT DATE_PART('year',order_month) AS year, AVG(monthly_sales) AS  
avg_monthly_sales  
FROM (SELECT DATE_TRUNC('month',order_date) as order_month, SUM(sales)  
AS monthly_sales  
FROM orders  
GROUP BY 1) AS temp  
GROUP BY 1  
ORDER BY 1 DESC;
```



Solo Exercise:

Categorizing Customers

5 minutes



Most customers make several hundred purchases from us as a business supplier. Let's look at how much of our revenue these customers make up.

Inner query: Create a query that categorizes customers by purchase frequency and returns total sales per customer:

- 1000+ = Supplier
- 500+ = Multiple
- All others

Outer query: How much in sales did each segment make?



Solo Exercise:

Categorizing Customers | Solution

```
SELECT customer_type, SUM(total_sales) AS total_sales
FROM
  (SELECT customer_id,
    CASE WHEN COUNT(DISTINCT order_id)>=1000 THEN 'Supplier'
    WHEN COUNT(DISTINCT order_id)>=500 THEN 'Frequent'
    ELSE 'All others' END as customer_type,
    SUM(sales) as total_sales
  FROM orders
  GROUP BY 1) AS temp
GROUP BY 1
ORDER BY 2 DESC;
```

Single Column Subquery Syntax

- The *inner query* retrieves a list of items in the form of **a single column**.
- The *outer query* then tests one of its columns against this list of items.

```
SELECT a_column
FROM a_table
WHERE column_id IN
    (SELECT column_id
     FROM a_different_table
     WHERE a_different_column meets some_condition);
```




Looking at orders in more detail, the VP of Product wants you to focus on the consumer segments and find out:

How many orders had products with a cost to consumer of more than \$500?

First, write a query that first selects the appropriate products:

```
SELECT product_id
FROM products
WHERE CAST(product_cost_to_consumer AS int) > 500;
```

Then, insert it into the outer query's WHERE clause (see next slide).



```
SELECT COUNT(DISTINCT order_id)
FROM orders
WHERE product_id IN (
    SELECT product_id
    FROM products
    WHERE CAST(product_cost_to_consumer AS int) > 500);
```

Here, the inner query must *only* return **one column** for an **IN** filter. Consider that the inner query is *building an array of values* that are then substituted into the **WHERE** clause in the outer loop.



Solo Exercise:

How About **NOT IN**?

5 minutes



Conversely, you can use **NOT IN** in the subquery to exclude data.

Going back to the previous query, we can find the number of **orders that had products with a cost to consumer of more than \$500** in either of two ways with the options below:

Give Option 2 a try, and we'll discuss as a class.

Option 1: Use the **IN** statement and the logical operator **greater than**.

Option 2: Use the **NOT IN** statement and the logical operator **less than**.



Query 1 Use IN

```
SELECT COUNT(DISTINCT order_id)
FROM orders
WHERE product_id IN (
  SELECT product_id
  FROM products
  WHERE
    CAST(product_cost_to_consumer
    AS int) > 500)
```

Query 2 Use NOT IN

```
SELECT COUNT(DISTINCT order_id)
FROM orders
WHERE product_id
NOT IN (
  SELECT product_id
  FROM products
  WHERE
    CAST(product_cost_to_consumer
    AS int) < 500)
```

Nested Subquery Syntax

- The *inner query* retrieves **the average of another column** from a different table.
- The *outer query* then takes the inner query and uses it with a comparison operator to filter a list of values.

```
SELECT *  
FROM a_table  
WHERE a_column <  
      (SELECT AVG(another_column)  
       FROM a_different_table);
```



Now that you have a gauge on the number of orders for products that had a cost of more than \$500 to the consumer, let's shift the focus to profit:

How many orders have more profit than the average product cost to consumer?

Use a subquery to *first* calculate the most expensive item to consumers:

```
SELECT  AVG(product_cost_to_consumer)
FROM    products;
```



Then, apply the subquery in the outer query:

```
SELECT COUNT(DISTINCT order_id)
FROM orders
WHERE profit >
    (SELECT AVG(product_cost_to_consumer)
     FROM products);
```



Single Column: Use SELECT to Find % of Total

Let's take our query from earlier and find the share of sales that each frequency of customer makes up. Add a SELECT subquery in our outer query to find the total sum of all sales.

```
SELECT customer_type, SUM(total_sales) AS total_sales
FROM
    (SELECT customer_id,
     CASE WHEN COUNT(DISTINCT order_id)>=1000 THEN 'Supplier'
     WHEN COUNT(DISTINCT order_id)>=100 THEN 'Frequent'
     ELSE 'All others' END as customer_type,
     SUM(sales) as total_sales
    FROM orders
    GROUP BY 1) AS temp
GROUP BY 1
ORDER BY 2 DESC;
```




```
SELECT customer_type, SUM(total_sales) AS total_sales,  
       SUM(total_sales)::numeric/(SELECT SUM(sales) FROM orders)::numeric AS share_of_sales  
FROM  
    (SELECT customer_id,  
     CASE WHEN COUNT(DISTINCT order_id)>=1000 THEN 'Supplier'  
     WHEN COUNT(DISTINCT order_id)>=100 THEN 'Frequent'  
     ELSE 'All others' END as customer_type,  
          SUM(sales) as total_sales  
     FROM orders  
     GROUP BY 1) AS temp  
GROUP BY 1  
ORDER BY 2 DESC;
```



Discussion:

Key Differences of Subquery Types

Now that we've seen all three types of subqueries, let's review their uses by focusing on the key differences among the three:

- **Derived Table** subqueries can provide a whole table.
- **Single Column** using IN or in the SELECT subqueries must provide a single column.
- **Nested subqueries** can be used to compare single values.



Any other differences that come to mind?

— Additional Subquery Use Cases

Running Multiple Queries at Once

Using the subqueries methodology, you can run multiple SELECT statements at once and get the results on the same screen.

The benefit: If you need fields from multiple data sources or from multiple parts of a single data set (e.g., for comparison).

What this looks like in practice:

Instead of separately running two SELECT statements, we can **wrap each SELECT statement in its own parentheses** to be evaluated separately in a common results pane.

Using CASE Statements in Subqueries

CASE statements can be used in *both* the inner query and outer queries.

What makes it cool? One useful and common construction is to:

- Use the inner (nested) query to create a binary value, sometimes called a “flag.”
- Then, use the outer query to average that value.

Let's build an example using a **CASE** classification and a binary value!



Let's say we classified profit per order as small, medium, large. Now we want to use those categories to count how many orders fit into each group. Let's look at the first classification query:

```
SELECT order_id, profit,  
       CASE  
         WHEN profit >= 1000 THEN 'large'  
         WHEN profit >= 50 THEN 'medium'  
         ELSE 'small'  
       END AS profit_size  
FROM orders
```



CASE in Subqueries (Cont.)

Here is the first query (in red) serving as our inner query, with the outer query counting occurrences:

```
SELECT profit_size, COUNT(DISTINCT order_id) AS number_orders
FROM (SELECT order_id, profit,
        CASE
            WHEN profit >= 1000 THEN 'large'
            WHEN profit >= 50 THEN 'medium'
            ELSE 'small'
        END AS profit_size
    FROM orders) AS temp
GROUP BY profit_size;
```



Imagine you want to get even more granular and find out, ***by product_id***, what **percentage of discounts were more than 25%**. First, create the inner query that classifies each discount as either more or less than 25%:

```
SELECT product_id, profit,  
       CASE  
         WHEN discount > 0.25 THEN 1  
         ELSE 0  
       END AS over_25  
FROM orders;
```




CASE in Subqueries With a Binary Value (Cont.)

Now, wrap query in an outer query that averages the *flag (binary value)* per store:

```
SELECT product_id, ROUND(AVG(over_25), 4) AS pct_over_25
FROM (SELECT product_id, profit,
        CASE WHEN discount > 0.25 THEN 1
        ELSE 0
        END AS over_25
        FROM orders) AS temp
GROUP BY product_id
ORDER BY product_id
LIMIT 100;
```

— Common Table Expressions (CTE)

What to Do When...

Your queries are starting to look like this:

```
SELECT customer_description, product_description,  
estimated_profit, order_year FROM  
(  
    SELECT  
        (t3.segment || ': ' || t3.customer_name) AS customer_description,  
        (t2.sub_category || ', ' || t2.product_name) AS product_description,  
        (CAST(t1.quantity AS INTEGER)  
         * CAST(t2.product_cost_to_consumer AS INTEGER)  
         * (1 - CAST(t1.discount AS NUMERIC))) AS estimated_profit,  
        EXTRACT(YEAR FROM t1.order_date) AS order_year  
    FROM orders AS t1  
    INNER JOIN products AS t2 ON t1.product_id = t2.product_id  
    INNER JOIN customers AS t3 ON t1.customer_id = t3.customer_id ) temp_tbl  
WHERE estimated_profit < 100 ;
```



Meet Common Table Expressions (CTE)

- A technique for creating a temporary result set that can be referenced within the following statements: **SELECT**, **INSERT**, **UPDATE**, or **DELETE**.
- CTEs are defined *outside* of the above statements using the **WITH** operator, making it a convenient way to manage more complicated queries.



After identifying the percentage of discounts by *product_id*, your Head of Sales now wants you to look at another factor that is affecting the overall profit:

Which customers are buying which products across all of our fiscal years with an estimated profit <\$100?

Here is the basic main/outer query using CTE:

```
SELECT *  
FROM t1_cte AS t1  
INNER JOIN t2_cte AS t2 ON ...  
INNER JOIN t3_cte AS t3 ON ...  
WHERE some_criteria
```



To keep your results organized and easy to read, add a CTE.

WITH

```
t1_cte (col_1, col_2, ...) AS (query1 ...),  
t2_cte (col_a, col_b, ...) AS (query2 ...),  
t3_cte (col_x, col_y, ...) AS (query3 ...)
```

Common table expressions (CTEs) consist of two parts:

- Table expression definition
- Query definition

SELECT *

FROM t1_cte AS t1

INNER JOIN t2_cte AS t2 ON ...

INNER JOIN t3_cte AS t3 ON ...

WHERE some_criteria

Main query using CTE

Now Your Query Looks Like This!

```
WITH t1_cte AS
(
    SELECT
        t1.product_id, t1.customer_id,
        (t2.sub_category || ', ' || t2.product_name) AS product_description,
        (CAST(t1.quantity AS INTEGER)
         * CAST(t2.product_cost_to_consumer AS INTEGER)
         * (1 - CAST(t1.discount AS NUMERIC))) AS estimated_profit,
        EXTRACT(YEAR FROM t1.order_date) AS order_year
    FROM orders AS t1
    INNER JOIN products AS t2 ON t1.product_id = t2.product_id
)
SELECT
    (t3.segment || ': ' || t3.customer_name) AS customer_description,
    t1.product_description,
    t1.estimated_profit,
    t1.order_year
FROM t1_cte t1
    INNER JOIN customers t3 ON t1.customer_id = t3.customer_id
WHERE estimated_profit < 100 ;
```



Discussion:

How to Create Multiple CTEs

As it turns out, another analyst already wrote a query with multiple CTEs. Take a look at the code below. **Together, let's describe what this query will execute.**

```
WITH return_cte AS
    (SELECT order_id, reason_returned
     FROM returns
     WHERE reason_returned NOT LIKE 'Not Given'),
  order_cte AS
    (SELECT order_id, sales
     FROM orders)

SELECT order_cte.order_id, order_cte.sales,
       return_cte.reason_returned
FROM return_cte INNER JOIN order_cte
ON return_cte.order_id = order_cte.order_id
LIMIT 100;
```




Convert Derived Table Subquery to a CTE

Let's take the following derived table subquery and convert it to a CTE.

```
SELECT profit_size, COUNT(DISTINCT order_id) AS number_orders
FROM (SELECT order_id, profit,
      CASE
        WHEN profit >= 1000 THEN 'large'
        WHEN profit >= 50 THEN 'medium'
        ELSE 'small'
      END AS profit_size
      FROM orders) AS temp
GROUP BY profit_size;
```



Convert Derived Table Subquery to a CTE (Con't)

The subquery is now a CTE and used in the final query.

```
WITH cte_profit_size AS (  
  SELECT order_id, profit,  
         CASE  
           WHEN profit >= 1000 THEN 'large'  
           WHEN profit >= 50 THEN 'medium'  
           ELSE 'small'  
         END AS profit_size  
  FROM orders)  
SELECT profit_size, COUNT(DISTINCT order_id) AS number_orders  
FROM cte_profit_size  
GROUP BY profit_size;
```

Benefits of Using CTEs

- **Readability:** Option to create chunks of data that would then be combined in a final SELECT statement.
- **Substitute for a View:** Use when you don't have permission to create a view object or the view would only be used in this one query.
- **Limitations:** Overcome SELECT statement limitations, such as referencing itself (recursion) or performing GROUP BY using a scalar subselect or non-deterministic functions.



Solo Exercise:

Optional Homework

Use subqueries to answer the following questions:

1. What is the percent of orders shipped using Standard Class?
2. What percent of all sales in the United States did returns make up in 2019?
3. What is the average monthly profit by region? Which region has the highest monthly profit?
4. Which orders have a return quantity of more than 2 (include the `product_id` and `product_name`)



Solo Exercise:

Optional Homework | Solution 1

```
SELECT
    ROUND(AVG(standard_class) *100 ,2) as avg_sales_by_rep
FROM
    (
        SELECT
            CASE WHEN ship_mode = 'Standard Class' THEN 1 ELSE 0 END as
standard_class
            FROM orders
        ) ship_mode
;
```

Returns 40.49



Solo Exercise:

Optional Homework | Solution 2

```
SELECT
    SUM(sales) as sum_returned_sales,
    (SELECT SUM(sales) FROM orders) as total_sales_all_orders,
    ROUND((SUM(sales) / (SELECT SUM(sales) FROM orders)) * 100 ,2) as
pct_returned_sales
FROM orders ord
JOIN returns rtn ON ord.order_id = rtn.order_id
WHERE order_date BETWEEN '2019-01-01' AND '2019-12-31'
;
```

pct_returned_sales 5.11



Solo Exercise:

Optional Homework | Solution 3

```
SELECT region, ROUND(AVG(monthly_profit),2) AS avg_monthly_profit
FROM
    (
        SELECT DATE_TRUNC('month',ord.order_date) as order_month, reg.region,
        SUM(ord.profit) AS monthly_profit
        FROM orders ord
        LEFT JOIN regions reg ON ord.region_id = reg.region_id
        GROUP BY 1,2
    ) AS temp
GROUP BY 1
ORDER BY 2 DESC;
```

Americas has the highest average monthly profit: \$11,816.12



Solo Exercise:

Optional Homework | Solution 4

```
SELECT
    ord.order_id,
    ord.order_date,
    ord.product_id,
    prd.product_name
FROM orders ord
JOIN products prd ON ord.product_id = prd.product_id
WHERE order_id IN (SELECT order_id FROM returns WHERE return_quantity > 2);
```

Returns 115 rows