

Week1: Organizing data in spreadsheet and SQL

GRIT!



Analysis

- the process used to make sense of the data collected.
- Goal: identify trends and relationships within data, so that you can accurately answer to questions you are asking.
- "* The 4 Phases of Analysis" → make it easier to search, understand, e.g. organize data into a table, hide / show column
 - ① Organize data → observing and organizing data in a way that's easy to reference.
 - ② Format & Adjust data → adjusting, make it more easy to digest like SORT & FILTER
 - finding outliers (data points that are very different from similarly collected data & might not be reliable values.)



- * Sort → arrange data into meaningful order (ASC, DESC), group similar data
- * Filter → seeing data that meets a specific criteria

- ③ Get input from others → checking others' analysis e.g. who experiences on this.
- ④ Transform data → identify patterns, calculations (e.g. correlation relationships between 2 variables),

Spreadsheet Workshop Commands →

- Sort
 - ✓ Sort sheet → all of the data in a spreadsheet is sorted by the ranking of a specific sorted column - data across rows is kept together.
 - ✓ Sort Range → Nothing else on the spreadsheet is rearranged besides the specified in a column.
- ✓ Sort Function → =SORT(A2:D6, sorted_by-column, asc/desc)
- ✓ Customized Sort Order → Sort data in a spreadsheet using multiple conditions:
on Data → Sort range > Advanced Range sorting ops.

Get to know
Sort
&
Order By

SQL Workshop Commands →

- ✓ WHERE clause → to filter only values that meet a specific criteria.
- ✓ SORT BY (ASC by default) → ASC or DESC

Week2: Convert and Format Data Combine multiple datasets

Incorrectly formatted data can:

- Lead to mistakes ✗
- Take time to fix 🕒 ...
- Affect stakeholder's decision-making 😬

Spreadsheet Workshop →

- ✓ Typecast numbers (Unit Conversion)
 - tab above (Format > Number)
 - =CONVERT(E2, "mph", "mls")
 - =CONVERT(E2, "C", "F")
(Convert cell E2 from Celsius to Fahrenheit)
- ✓ VALUE() convert text string to number.
- ✓ LEN, LEFT, RIGHT, FIND
 - Find position of text/string in a cell

v Data Validation (Module 4 Process)

In Spreadsheet:

- * Add dropdown lists with predefined options.
- Create custom checkboxes
- * Protect structured data & formulas (Reject input) will protect from mistyping e.g. is valid
- * Conditional Formatting highlight cells when matches the condition

SQL Workshop →

- * CONCAT
 - CONCAT("google", ".com")
 - CONCAT()
 - CONCAT_ws("www", "google", ".com")
 - CONCAT_WSC()
 - CONCAT with +
- CASTC ~ AS int64)

Get support during analysis

- Encounter error when calculating time (start time is greater than end time)
change start-end to =IF(end>start, end-start, 1-start+end)
- ASK team
- ★ - Ask/find questions online
 - > Best practices for searching online
 - Thinking skills (mental model: your thought process & the way you approach a problem)
 - Data Analytics terms (use the right terms when searching) ⚡
 - Basic knowledge tools (can apply new with old knowledge) 📈

Week 3: Aggregate data for analysis
(VLOOKUP, JOINS, Subqueries)

- Aggregation → collecting or gathering many separate pieces into a whole.
→ organizing pieces of data by Average, min, Max, Sum

Helps:

- Identify trends ↗
- Make comparisons ⚡
- Gain insights 🌟

Tools for doing Aggregation:

* ① VLOOKUP (map one product/value to another table)

=VLOOKUP(A2, A2:A14, 2, FALSE)
 search key ↑ range ↑ index ↑ exact match
 (return column) *

Spreadsheet

3

Troubleshooting Questions:

- How should I prioritize these issues?
- In a single sentence, what's the issue I'm facing?
- What sources can help me solve the problem?
- How can I stop this problem from happening in the future?

* You can protect sheet by clicking on Data > Protected sheets & ranges *
to protect the entire sheet / cell from editing *

I - JOINS

- (INNER) JOIN (by default) 🔪
- LEFT JOIN
- RIGHT JOIN 🔪
- FULL OUTER JOIN 🔪 (fetch all records)



SQL

- COUNT [COUNT not count duplicates] COUNT DISTINCT return distinct values in that range.
- * Inner query executes first *
- Subqueries (or Nested Query / Inner Query / Inner Select)
 - HAVING (WHERE first) to aggregate oh statistics functions.
 - CASE → return record that matches like IF-THEN
 - Aliasing → easier & shorter to reference.



Nested Query

A	B	C	D
	Hours Worked	Rank	Employee #
1	10	1	FT12578
2	20	2	FT12579
3	20	3	FT12580
4	20	4	FT12581

Fix by using INDEX() or MATCH()

```

1 SELECT *
2   FROM (
3     SELECT warehouse.state,
4            COUNT(DISTINCT order_id) AS num_orders
5   FROM warehouse
6   GROUP BY warehouse.state
7   ORDER BY num_orders DESC
8 ) AS T1
9   INNER JOIN (
10    SELECT warehouse.state,
11           MAX(order_id) AS max_order_id
12    FROM warehouse
13   GROUP BY warehouse.state
14 ) AS T2
15   ON T1.state = T2.state
16 WHERE T1.state = T2.state
17   AND T1.order_id < T2.max_order_id
18   AND T1.order_id <= T2.max_order_id
  
```

EXAMPLE:

* array = collection of values
like price column

Week 4: Performing Data Calculations (Pivot, SQL calculations, Data-validation, Temp tables)

Spreadsheet:

- ① - COUNTIF(range, condition) - MAXIFS()
- SUMIF(range, condition, sum-range)

functions - SUMPRODUCT(range, [range]) → multiply each record horizontally, then sum it up
e.g. FORMULAS

- ② Pivot Table

(row, column, values, filter)
e.g. Quantity × Sales × Margin of each row
and plus it together



SQL:

- GROUP BY ()

- EXTRACT(YEAR FROM start-time) → to extract a part from a given date. from date/date time

* Types of Data Validation: → to make data become complete and accurate, secure, consistent

- ① Data Type → check data type
- ② Data Range → check data falls within an acceptable range
- ③ Data Constraints → check data meets certain conditions/criteria.
- ④ Data Consistency → check that data makes sense in the context of other related data.
- ⑤ Data Structure → check that data follows / conforms to a set structure
- ⑥ Code Validation → check that the application code systematically performs e.g. More than 1 data type allowed

 Temporary table → PROS: Don't have to query in large dataset many times when using often.
a database table that is created and exists temporarily on a database server.

Global Local  The WITH clause is a type of temporary table that you can query from multiple times. (like CREATE VIEW on MySQL) that'll be deleted when you end your SQL session.

WITH temp_table_name AS (

SELECT * FROM

—

)

CREATE TABLE statements

CREATE TEMP TABLE statements

not supported on BigQuery  SELECT ~ INTO temp_table_name FROM ~

preferred when finishes using :

- *- Dropping Table: removes info in the table and table variable definitions.
- Deleting Table: removes only info in the table.

Time to
"Share"
your data

