

Week 2 Course 05: Analyze Data to Answer Questions

let's get organised

Page No. _____
Date 06/02/23

Instructor : Ayanna
Global Insights Manager

Best way to know if your creating value is if you have evidence i.e. data

Created a learning log & took brief overview of what's to come

2/5

Analysis:

The process used to make sense of data collected

The goal of analysis is to identify trends and relationship within data so you can accurately answer the questions you're asking.

4 Phases of Analysis :

- ① Organize Data : online registry / etc
- ② Format and adjust data : sort & filter
- ③ Get input from others : viewpoint you don't have
- ④ Transform data : Identifying key patterns

"The most important thing that you need to have throughout this learning journey is grit".

Always a need to Organize

- * Sorting and filtering are affected by the type of data we're working with.
- ↳ Just have your data in the right format.

Sorting ↗

Ranking data Sorting:
 Based on specific metric
 When you arrange data into a meaningful order to make it easier to understand, analyze and visualize

Filtering ↗

Showing only the data that meets a specific criteria while hiding the rest.

```
SELECT * FROM movie_data.movies
WHERE genre = 'Comedy';
```

* Advanced Sorting:

① 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.

① Sort Range :-

Nothing else on the spreadsheet is rearranged beside the specified cells in a column.

Doesn't keep information across rows together.

→ 2 methods for sorting :

① Menu ✓

② Writing Sort function

Sheets :-

Select specific column > Data > Sort by Sheet
sort by range

In the same dataset we want to keep all of the information together, so first we'll sort sheet

* SORT Function :

(index)

= SORT (Range, Sort-by, , TRUE/FALSE)

= SORT (A1:D6, 2, TRUE)

→ sort the range by B col in Asc order

* Customized sort orders

When you sort data in spreadsheets using multiple conditions. [order is important]

Data > sort > Advanced > Add multiple columns

Sorting using SQL

* [ORDER BY] → Sorting

↓
Last clause in a query!

[WHERE] → Filtering

[AND] → Customized Sorting / multi-sort

Ex. ① SELECT * FROM genre
ORDER BY Name;

→ Alternative, Word

(2) SELECT * FROM Tracks
WHERE Composer = 'Chris Cornell'
ORDER BY genreID DESC;

→ You Know My Name

In Video 17
SELECT * FROM movies
WHERE genre = "Comedy"
AND Revenue > 300 000 000
ORDER BY Release-Date DESC;

You're as good at an analyst as your ability to ask eight questions.

In correctly formatted data can:

- lead to mistakes
- Take time to fix
- Affect stakeholder's decision-making

* Sheets > Format > More format / choose

* Changing unit of measurement!

$^{\circ}\text{F}$ → $^{\circ}\text{C}$ Temperature

= CONVERT(C1, "F", "C")

Tip: After adding / performing any function. Copy & paste the ~~data~~ as values

paste special > values only

= CONVERT(D2, "mph", "m/s")

Data Validation (function)

Allows you to control what can and can't be entered in your worksheet

- ① > Add dropdown lists with predetermined options

-
-
-

We wanna add a list of options

Select desired column

> Data > Data Validation

Criteria: [Not yet started, In Progress, Ready]

Done ✓

other things we can do with data validation:

- ② > Custom checkboxes

Criteria: checkbox

Use custom cell values

Checked : Approved

Unchecked : Not Approved

- ③ Protect structured data & formulas.
reject invalid inputs

Conditional formatting

Select desired > Format > |
 ↓

column to Set rule & apply to
 ① ②

Text is Exactly color ⇒ Red

[Not Yet Started]

[In Progress]

[Ready]

color = Yellow

color = Green

[Review date] : Orange, Today

- ① Highlighting cells that contains the word (target)
- ② color-coding cells that contain dates after today

* IMP

CONCAT for only 2 variables and
 CONCATENATE for more than two

CONCAT(A2, B2)

CONCATENATE(A2, " ", B2)

CONVERT(col, "F", "C")

week 2

Page No.

Date

10/10/23

2/3

Merging & multiple sources

Query describing user type their route and number of trips, duration.

```
SELECT user_type,
       CONCAT(start_station, " to ", end_station) AS route,
       COUNT(*) AS num_trips,
       ROUND(AVG(CAST(duration AS int64)/6, 2)) AS avg_time,
       FROM T_name
      GROUP BY start_station, end_station, user_type
      ORDER BY num_trips DESC
      LIMIT 10;
```

→ → →

LEN, LEFT, RIGHT, FIND in spreadsheet

Same as that of in SQL

LEN(C2) → 14

FIND(C2, " ", C2)

↑
substring

finds space in C2.

Right 8 chars RIGHT(C2, 8) → 12345678
Left 11 chars LEFT(C2, 11) → 1234567890123

Left 11 chars LEFT(C2, 11) → 1234567890123

* `CONCAT ('Google', '.com')`; → Google.com

* `CONCAT_WS ('.', 'www', 'google', 'com')`

↓
→ www.google.com

Adds 2 or more strings with a separator

* `CONCAT` with + is allowed,

i.e. 'Google' + '.com'

```
SELECT CASE WHEN A>10 THEN "A"+" is great",
           END AS Col
FROM LMAO;
```

→ → →

* `CASE Revision` ↗

```
SELECT CASE
        WHEN condition THEN Sol1
        WHEN Condition -1
        WHEN Condition 2
        END As NewColumn
```

From Tname;

→ → →

CAST `SELECT CAST (myDate AS STRING) FROM Tname;`

DATETIME

`SAFE_CAST` → returns null instead of error

* What to do when you get stuck?

- > Ask a senior or mentor
- > Consult a team member
- > Search online
- > Utilize problem solving techniques.
- > Use online forums.

Reaching out is the way to go!

Layla: Analytical Lead (Google (obv))

↓
Someone who helps advertisers understand the value of their money.

"The analyze stage is where you become Expert about your dataset"

You get to tell the story once the dataset is cleaned and loaded in Excel/v. tool

Creating a pivot table to get the numbers, ~~predict~~ what is happening & what should they expect.

* Narrow down the error to ① Formula or
double check the formula & if not then data.

② Data

Page No. 1

* Finding answers online to your problems can be
empowering and can give you new knowledge
for future.

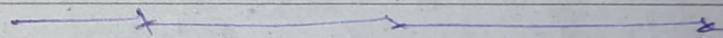
Best practices for searching online is

- > Thinking skills
- > Data Analytics terms
- > Basic knowledge of tools.

Mental model :-

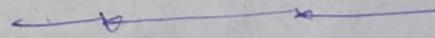
Your thought process and the way you approach
a problem.

* Use the right terms.



* Upcoming we'll learn about R

> R is a programming language popularly
used for statistical analysis, visualization and
other data analysis.



* Being able to modify example code is a
must know skill.



Vlookup for data Aggregation

Data Aggregation ↗

The process of gathering data from multiple sources in order to combine it into a single summarized collected

Puzzle pieces = data

Organisation = aggregation

Pile of pieces = summary

Putting the pieces together = gaining insights

- Identify trends
- make comparisons
- Gain insights.

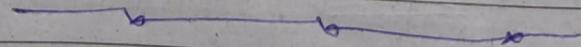
Data can also be aggregated over a given time period to provide statistics such as

- > Averages
- > Minimums
- > Maximums
- > Sums

functions help make data aggregation possible.

* Subquery ↗

A query within a query.



VLOOKUP(value, Range, Return, ~~TFP~~)
Index

$F \Rightarrow \underline{\text{Exact}}$

Data Aggregation tool

1

Vlookup :

A function that searches for a certain value in a column to return a corresponding piece of information.

VALUE :

A function that converts a text string that represents a number to a numerical value.

Before using the lookup()

* Make sure of ↗

> Data is formatted correctly : VALUES()

> There are no extra spaces : TRIM C1

> Duplicates : Remove duplicates

Lookup Action ↗

=VLOOKUP(103, A2:B26, 2, FALSE)

* Most commonly Vlookup is used to populated data from one sheet to another sheet.

From the Example is

=VLOOKUP(A2, 'Emp Sheet'!\$A\$2:\$B\$5, 2, FALSE)

Lock the sheet > Data
> protect sheets and ranges

No.	Date		

VLOOKUP (key, range, index, T/F)

* Limitations of VLOOKUP & common problems

Troubleshooting

- ① How should I prioritize these issues?
- ② In a single sentence, what's the issue I'm facing?
- ③ What resources can help me solve the problem?
- ④ How can I stop this problem from happening in the future.

(*) VLOOKUP only returns the first match it finds

(*) Can only get the data from the right, can't look left
[data we want] [lookup value]

* Absolute Reference:

A ref. that is locked so that rows and cols don't change when copied.

* MATCH

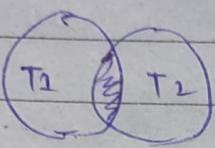
A function used to locate position of specific lookup value

TRUE : Approximate match

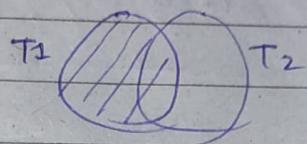
FALSE : Exact match

Default Join \Rightarrow Inner Join~~Start working with ~~Joins~~ Joins~~Page No.:
Date: 15 02 23Joins \rightarrow

A SQL clause that is used to combine rows from two or more tables based on a related column.

Common Joins \rightarrow 

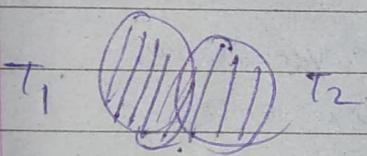
Inner Join



Left Join



Right Join



Full Outer Join

Inner Join \rightarrow A function that returns records with matching values in both tables.

Left Join \rightarrow A function that will return records from the left table and only the matching records from the right table.

The table mentioned first is left and table mentioned second is right.

Right Join \rightarrow A function that will return records from the right table and only the matching records from the left table.

Outer Join

A function that combines RIGHT & LEFT JOIN to return all matching records in both tables.

If there are no matches it'll create Null values in those places.

Example:

SELECT

employee.name AS employee_name,
 employee.role AS employee_role,
 departments.name AS department_name

FROM

employee_data.employees

INNER JOIN

employee_data.departments ON

emp.dept_id = dept.dept_id

→ → →

If we instead use LEFT JOIN

RIGHT JOIN

FULL OUTER JOIN

values will change accordingly!

We'll be using this to answer the question
"How many?"

Page No.			
Date			

- * COUNT in Spreadsheet :-
Can be used to count the total number of numerical values within a specific range in spreadsheet.
- * COUNT in SQL :-
A query that returns the number of rows in a specified range.
- * COUNT DISTINCT :-
A query that only returns the distinct values in a specified range.
- * Aliasing :-
When you temporarily name a table or column in your query to make it easier to read and write.

Ex. :-

= SELECT orders.*,
ware.ware - alias
ware.state
From ware.orders As orders
~~From~~ JOIN
ware-orders .warehouse ware ON orders.warehouse_id

SELECT COUNT (col)

SELECT COUNT (DISTINCT col)

Queries inside query

Subquery:

A SQL query that is nested inside a longer query.

Nesting dolls → matryoshka (Russian Nesting dolls)

Keep in mind inner query executes first.

Usually subqueries are with From or Where Clause

Example →

SELECT AVG(salary)

From: (SELECT name
WHERE name LIKE '%.%'
FROM table);

IN, ANY, ALL

SELECT column_name FROM T_name
WHERE salary > ALL (SELECT ...);

Salary IN

Salary = ANY

normal flows WHERE & then
GROUP BY

Page No.		
Date		

Wherever we want to have GROUP BY first then WHERE clause, we need to use HAVING

* HAVING ↴

Allows you to add a filter to your query instead of the underlying table that can only be used with aggregate functions.

* CASE :

Returns records with your conditions by allowing you to include if/then statements in your query.

Clauses like HAVING & CASE, paired with subqueries will help you build more and more subqueries complex queries.

Justin, Data Analytics lead!

"Your career path is not always straight forward"

Keep changing little by little, figure out what's exciting about your role RN.

Find a job that lets you do more of what you like

Be curious! Ask why's