

* Data Calculations :-

The better you get at SQL and pulling data from data tables, the faster you'll complete your analysis.

This week:-

- > Formula for basic calculations
- > Conditional formulas that use the IF function
- > The SUMPRODUCT function
- > Pivot tables to organize calculations
- > Queries & calculations in SQL
- > Temporary tables in SQL.

* Common calculation formulas:-

=SUM(B2:B12)

= MAX

= MIN

* Functions

COUNTIF (range, "value")

Summary tables:-

A table used to summarise statistical information about data.

SumIF (range, condⁿ, SumRange)

Page No.	
Date	

Quantity	Count	Revenue Total	Average revenue per transaction
1	25		
>1			

G11 \Rightarrow COUNTIF (B3:B50, "=1")

G12 \Rightarrow COUNTIF (B3:B50, ">1")

H11 \Rightarrow SUMIF (B3:B50, "=1", C3:C50)

H12 \Rightarrow SUMIF (B3:B50, ">1", C3:C50)

I11 \Rightarrow H11 / G11

I12 \Rightarrow H12 / G12

* COUNTIFS \Rightarrow (both True)

= COUNTIFS (range1, cond1, range2, cond2, ...)

* SUMIFS \Rightarrow

= SUMIFS (Sum-Range1, ~~Cond1~~, ~~Cond1~~, ^{Cond1}crit1, ~~Cond2~~, ^{Cond2}crit2, ...)

= SUMIFS (B2:B22, E12:E20, "=1", D10:D15, "=2")

* MAXIFS (maxrange, range1, crit1, range2, crit2 ...)

MAXIFS (B2:B32, D2:D10, "=Apple", E2:E10, "=purchased")

Finds maximum quantity of Apples purchased

* Composite Functions &→

SUMPRODUCT:→

A function that multiplies arrays and returns the sum of those products.

$$= \text{sumproduct}(\text{array1}, \text{array2}, \dots)$$

B	C	
Quantity	Price per item	Profit margin
		$\text{Total} = B_1 \times C_1 + B_2 \times C_2 + \dots$

$$\therefore = \text{sumproduct}(B3:B7, C3:C7) = \$655$$

Profit margin:

A percentage that indicates how many cents of profit has been generated for each dollar of sale.

Apple	sold at 5\$	Profit 1\$	Profit margin 20%
	Price	Profit	
quantity	Per item	margin	

$$= \text{sumproduct}(B3:B7, C3:C7, D3:D7)$$

$$= \text{Total profit}$$

individually multiplies

$$B3 \times C3 \times D3 +$$

$$B4 \times C4 \times D4 +$$

$$B5 \times C5 \times D5 +$$

$$B6 \times C6 \times D6$$

$$\text{Total Profit} = B7 \times C7 \times D7$$

Using Pivot table

Movie dataset analysis:-

Analysis steps:-

- Find out how much revenue was generated each year.
- Build a pivot table to show the revenue per year.
- Find the average revenue per movie.
- Checkout findings for some possible trends.

Insert > Create Pivot Table

Rows > Add Release Date

But there are too many data

Right click on any cell & "group by year"

Now only 4 years

Values > Sum of Box office collections

The drop-down is used to change functions applied to the values.

Add more to values > Average of Box office Collections

Add to values > Count of Box office Revenue
i.e. no. of movies Released each year.

Year 2015 has highest no. of releases but lowest all around Box office Collection.

maybe the movies didn't earn much, etc.

To test our hypothesis that "may-be" money didn't earn that much"

↓
first copy the pivot table and paste it somewhere (maybe in same sheet)

↓
Change the column names

Sum < \$10M Average < \$10M Count < \$10M

* next we are going to use filters to find out how many movies earned less than \$10M revenue in 2015

Then (*) Create a calculated field to determine what percentage of total movies from that year they represent.

Applying filter is

Select cell & add filter, it will be applied to the entire table

Filter by condition: less than
10,000,000.00

there are 20 movies that earned < \$10M but not enough to be called an insight



'for_spaces in names' in Sheet

Page No.	
Date	

Create a new column "Calculated Field"

⊛⊛ A calculated field is a new field within a pivot table that carries out certain calculations based on the values of other fields.

Formula:→

$\text{Sum}(\text{Box office Revenue (\$)}) / \text{COUNT}(\text{BOR})$

Summarize by

Custom

Show as

Default

Percent of total movies

= D11/D2

$[\text{<10m\$} / \text{Total movie}]$

2 make it percentage

We came to know, 16.13% movies in 2019 earned <10m\$

while other years have this number around 10% only.

→

We covered 2 function, Formulas & Pivot Tables.

Learn more SQL calculations

Operators

A symbol that names the type of operation or calculation to be performed in a formula.

+ addition

- subtraction

* multiplication

/ division

In Both SQL & Sheets

SELECT COLA, COLB, COLA+COLB AS COLX
FROM T-name;

For multiple calculation use parenthesis

SELECT COLA, COLB, COLC,
(COLA+COLB) * COLC AS COLY
FROM T-name;

* % Modulo → Returns remainder when one num. divided by another

Spreadsheet Function
Sum

AVERAGE

SQL Function

Sum

AVG

\neq NOT Equal

Embedding Simple calculation in SQL

Page No.	
Date	

```
SELECT Date, Region, Small_Bags, Large_Bags, XLarge_Bags,  
       Total_Bags, Small_Bags + Large_Bags + XLarge_Bags  
       AS Total_Bags_Calc  
FROM avocado_data, avocado_prices;
```

Q. what percentage bags were Small_Bags, Large_Bags and XLarge_Bags

```
SELECT Date, Region, Total_Bags, small_Bags,  
       Large_Bags, XLarge_Bags,  
       (Small_Bags / Total_Bags) * 100 AS Small_Percent,  
       (Large_Bags / Total_Bags) * 100 AS Large_Percent,  
       (XLarge_Bags / Total_Bags) * 100 AS XLarge_Percent  
FROM avocado_data, avocado_prices;
```

multiplying by 100 gives us Percentage

Error \rightarrow You can't divide by zero

If somewhere, total_bags = 0
so find it & fix it

Include a WHERE

WHERE ~~total_bags~~ <> 0 ALL (small, large, Xlarge);

You can also use SAFE_DIVIDE

Extract (YEAR FROM Col)

ORDER BY ASC/DESC

Page No.	
Date	

Group By :

A command that groups rows that have the same values from a table into summary rows.

```
SELECT  
FROM  
WHERE  
GROUP BY
```

Extracts:

Let's us pull one part of date

```
SELECT EXTRACT (YEAR FROM STARTTIME) AS Year  
COUNT (*) AS no-of-rides
```

FROM

BigQ - pub-data. new-york - citibik - trips

GROUP BY Year

ORDER BY Year

— x — x — x —

Data Validation 3

Page No.

Date

> ADDING DROP-DOWN to cells.

Checking and rechecking the quality of your data so that it is complete, accurate, secure and consistent.

* Types of Data validation

- | | | |
|----------------------|--|--|
| ① Data Types | | str, num |
| ② Data Range | | 1-12 |
| ③ Data Constraints | | must be whole numbers. |
| ④ Data Consistencies | | final date can't be before initial date |
| ⑤ Data Structure | | music in mp3 & not HTML. |
| ⑥ Code Validation | | checking system perform previously mentioned validations systematically. |

③ Before you upload files to SQL, you can import them into a spreadsheet in sheets to get comfortable with the data before you start.

⑤ This might not be possible with large datasets, but we should explore as much as possible.

Ask
Prepare
Process
Analyze
Share
Act

} 6 Steps of Data Analytics

Description in Briefly
→ makes it less complex & easier to read.
* Temporary Tables in SQL

Temporary tables →

A database table that is created and exists temporarily on a database server.

> Usually used to store subsets of data from standard data tables for certain period & they're automatically deleted, when you end session!

Ex. Running a query and storing result.

⊕ A Analyst has a large number of records in a table. They want to perform calculations on a small subset of the table. Rather than filtering data over and over

They should use a "Temporary Table"

The WITH clause is a type of temporary table that you can query from multiple times.

Ex.

Temp
created →

```
WITH tempname AS ( SELECT * FROM old_Table  
WHERE Cond^n >= 60)
```

Temp
used →

```
SELECT COUNT(*) AS Co FROM temp_name;
```


There's always more than one way!

WITH

③

Page
Date

SELECT INTO
CREATE TABLE

(another syntax) CREATE TEMP TABLE

* SELECT INTO →

copies data from existing table to new table

> Useful if you wanna make copy with a specific condition

like query with a WHERE clause

Ex: →
SELECT *
INTO AfricaSales
FROM GlobalSales
WHERE region = "Africa" ;

⊗ Good practice when you want to keep Database uncluttered & you don't need other people using the table.

⊗ If lots of people use it, then CREATE TABLE is a better option.

Adds table into database:

```
CREATE TABLE Africa_Sales AS (  
  SELECT * FROM GlobalSales  
  WHERE Region = 'Africa')
```

Adding metadata → to describe the data