

# Exploratory Data Analysis (EDA) Project

## What is EDA?

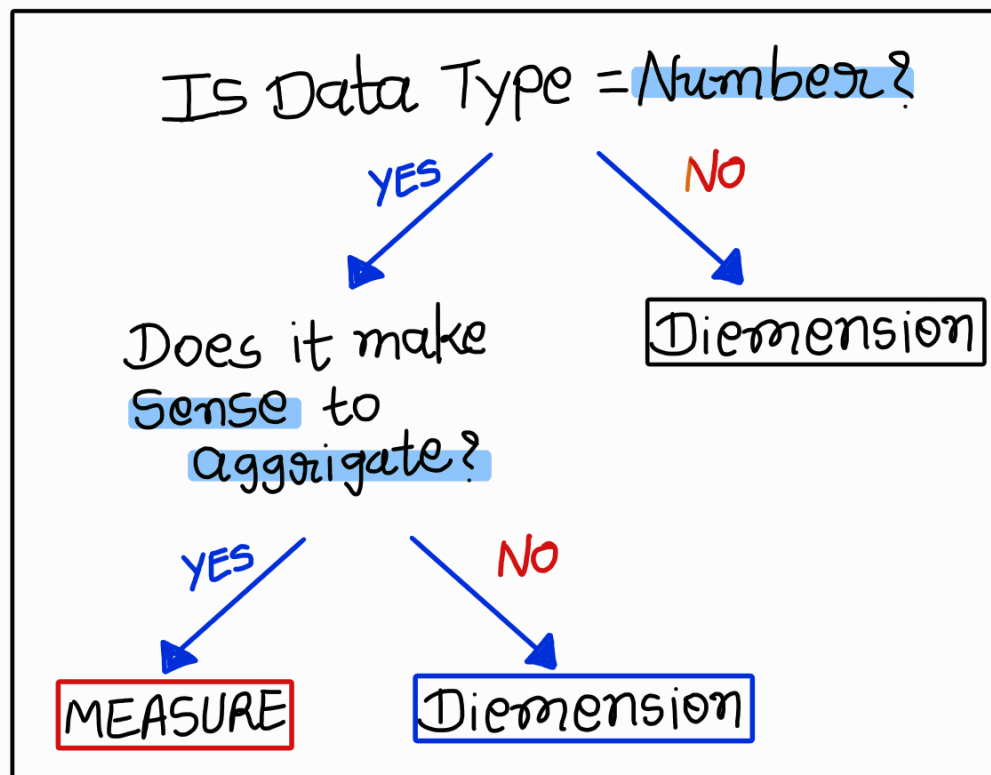
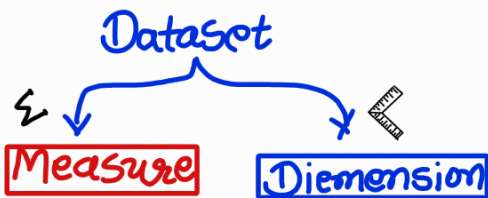
- Understand and cover insights about our datasets.
- This kind of project - How to ask right question and How to Find answer using SQL.

## Exploratory Data Analysis (EDA)

"Understand data"

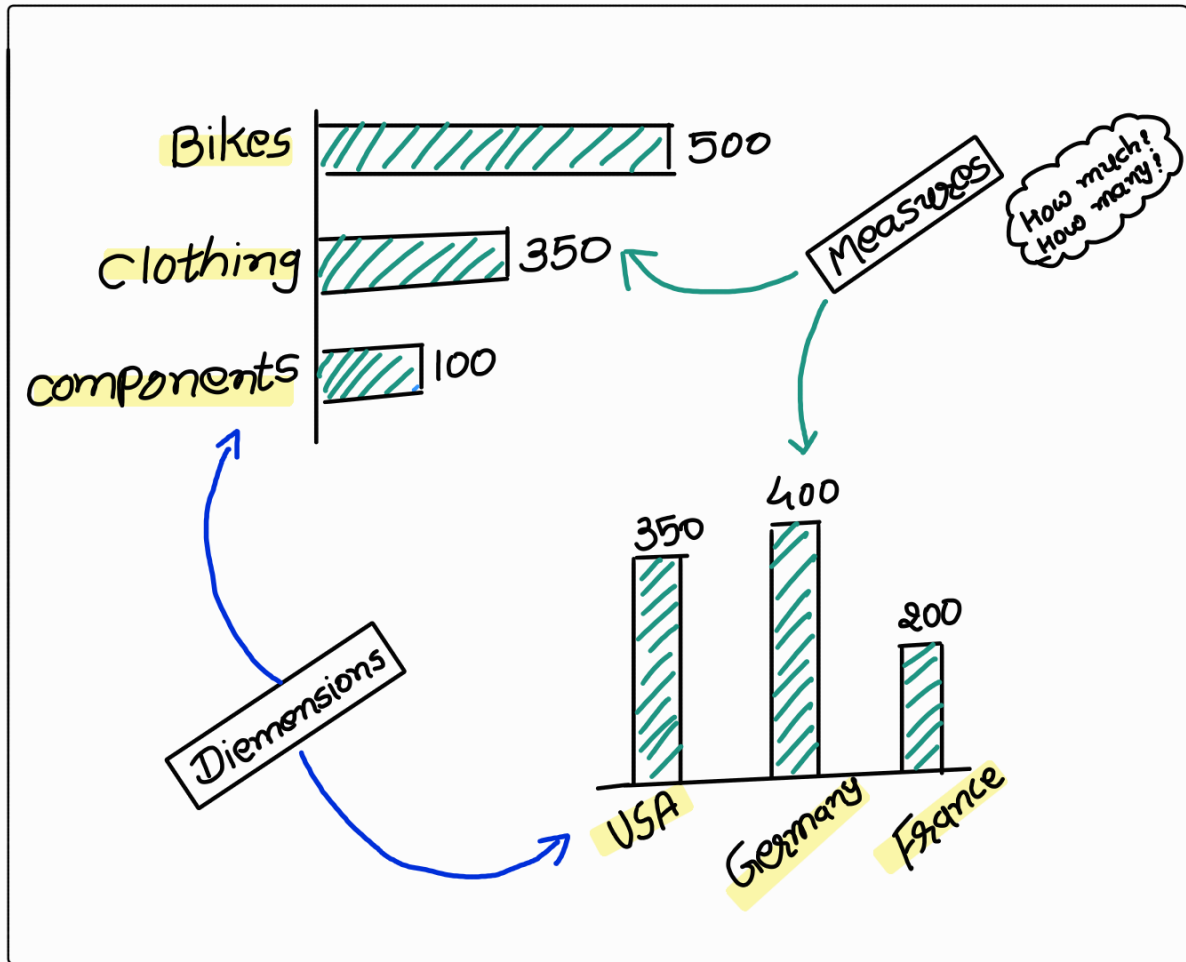
- Basic queries
- Data profiling
- Simple aggregations
- Subquery

## The secret MEASURE & DIMENSION



\* why need measures & Dimensions?

→ For Grouping up data



## 1. Database Exploration

→ Explore the structure of database for basic understanding about the tables, views, columns

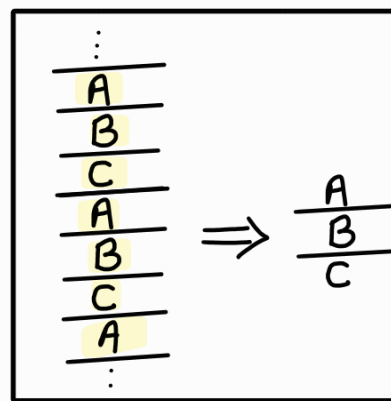
## 2. Dimension Exploration

→ Identifying the unique values (or categories) in each dimension.  
→ Recognizing how data might be grouped or segmented, which is useful for later analysis.

→ Formula:-

**DISTINCT** [Dimension]

Ex. Distinct Country  
Distinct Category  
Distinct Product



### 3. Date Exploration

→ Identify the earliest and latest dates (boundaries).

→ Understand the scope of data and the timespan.

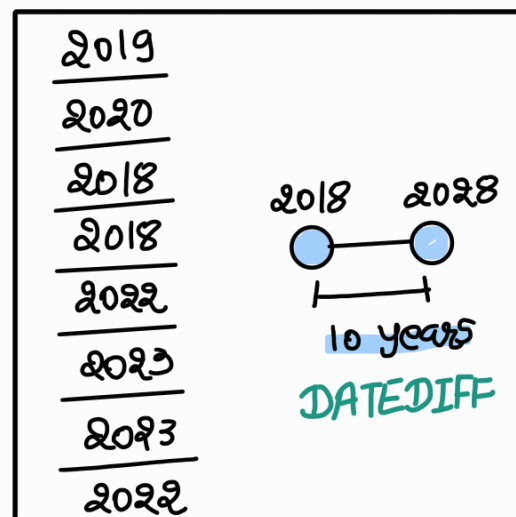
→ Formula:-


**MIN/MAX** [Date dimension]

Ex. **MIN** order\_date

**MAX** create\_date

**MIN** Birthdate



4. 

### Measures Exploration (Big Numbers)

→ Calculate the key metric of the business (Big numbers)

→ Highest level of Aggregation / Lowest level of details

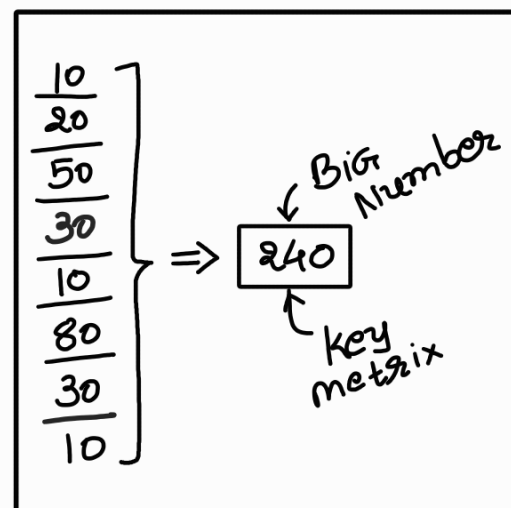
→ Formula:-

**$\Sigma$**  [Measure]

Ex. **SUM** (sales)

**AVG** (Price)

**SUM** (Quantity)



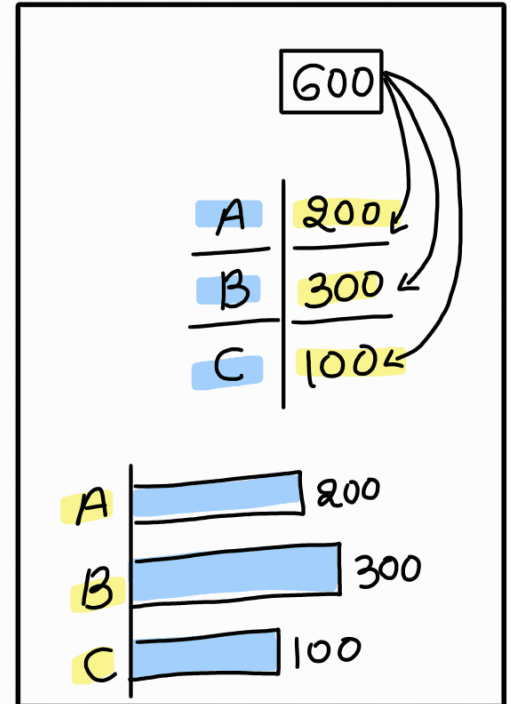
## 5. Magnitude

- compare the measure values by categories.
- It helps us understand the importance of different categories.

→ Formula:-

$\Sigma$  [measure] By [Dimension]

Ex. Total sales by country  
 Total quantity by category  
 Total Price by product  
 Total orders by customer



## 6. Ranking Top N-Bottom N

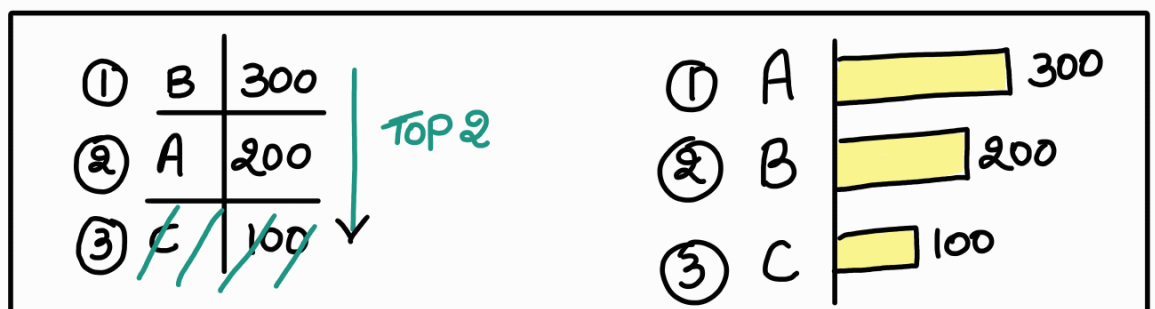
- order the values of dimensions by measures.
- Top N numbers | Bottom N numbers

→ Rank [Dimension] By  $\Sigma$  [Measure]

Ex. Rank category by total sales

Top 5 product by quantity

Bottom 3 Customers By total orders



→ For this use:

- TOP
- RANK()
- DENSE-RANK()
- ROW\_NUMBER()