

# Data Preprocessing-Data Extraction

---

## Data Extraction

Data extraction is the process of obtaining data from a source, such as a database, file, or website, and preparing it for further analysis. It involves identifying, collecting, and retrieving relevant data to support decision-making, research, or other purposes.

### Steps involved in Data Extraction:

1. **Identify the Source:** Determine where the data is located, whether it's a database, spreadsheet, text file, or website.
2. **Define the Requirements:** Specify what data needs to be extracted and why.
3. **Extract the Data:** Use tools or techniques to collect the required data from the source.
4. **Transform the Data:** Convert the extracted data into a format suitable for analysis.

### Example:

Suppose we want to analyze customer purchasing behavior for an e-commerce website. We need to extract relevant data from the database, which contains information about customers, orders, and products.

**Source:** Customer Database

**Requirements:** Extract customer demographic data (name, age, location), order history (date, product name, quantity), and purchase frequency.

### Data Extraction:

Using a SQL query or a data extraction tool like ETL (Extract, Transform, Load) software, we extract the relevant columns from the database:

```
SELECT
  customer_name,
  age,
  location,
  order_date,
  product_name,
  quantity
FROM
  customers
INNER JOIN
  orders ON customers.customer_id = orders.customer_id;
```

### Transformed Data:

The extracted data is then transformed into a suitable format for analysis, such as a CSV file or a pandas DataFrame:

customer_name	age	location	order_date	product_name	quantity
John Doe	32	New York	2022-01-01	iPhone	1
Jane Smith	25	London	2022-01-15	Laptop	2

The extracted and transformed data is now ready for analysis, such as data cleaning, feature engineering, and modeling.

### **Why Data Extraction is important:**

Data extraction is a crucial step in the data science workflow because it ensures that the data used for analysis is accurate, complete, and relevant. Poor or incomplete data can lead to flawed conclusions and decision-making. By extracting the right data from the correct source, we can ensure that our analyses are reliable and actionable.

---

*Curated by Brajesh Kumar*