

# Spark Project-Real World E-commerce data

In this project, we will embark on a detailed exploration of a real-world e-commerce dataset

<https://www.kaggle.com/datasets/olistbr/brazilian-e-commerce>

## Project Modules

### 1. Data ingestion and exploration

→ HDFS

Always NULL

- Set up our Hadoop and spark environment
- Import CSV into HDFS
- Load data into spark DataFrame
- Examining the schema and data types
- Perform EDA on top of our data

### 2. Data cleaning and Transformation

address quality issue and  
transform the data in a  
structured format

- handling of null
- Standardize date formats
- Numerize & scale numerical values
- Create new features

### 3. Data Integration and Aggregation

Combine data from multiple tables to  
Create a unified datasets.

→ Perform join

→ Aggregate data  
to compute metrics

→ Resolve any  
inconsistency  
arising from data  
integration.

### 4. Performance Optimization

enhance the efficiency of data  
processing tasks

→ data partitioning  
to optimize  
spark jobs

→ Vtage Caching  
for iterative  
operations

→ Optimizing the  
spark configuration

### 5. Data Serving

Make the processed data  
available

→ export transformed  
data to an external  
system / databases

→ Create visualization  
to represent data  
trends & patterns

# Exploration and Understanding of Data

Brazilian E-commerce : Olist dataset

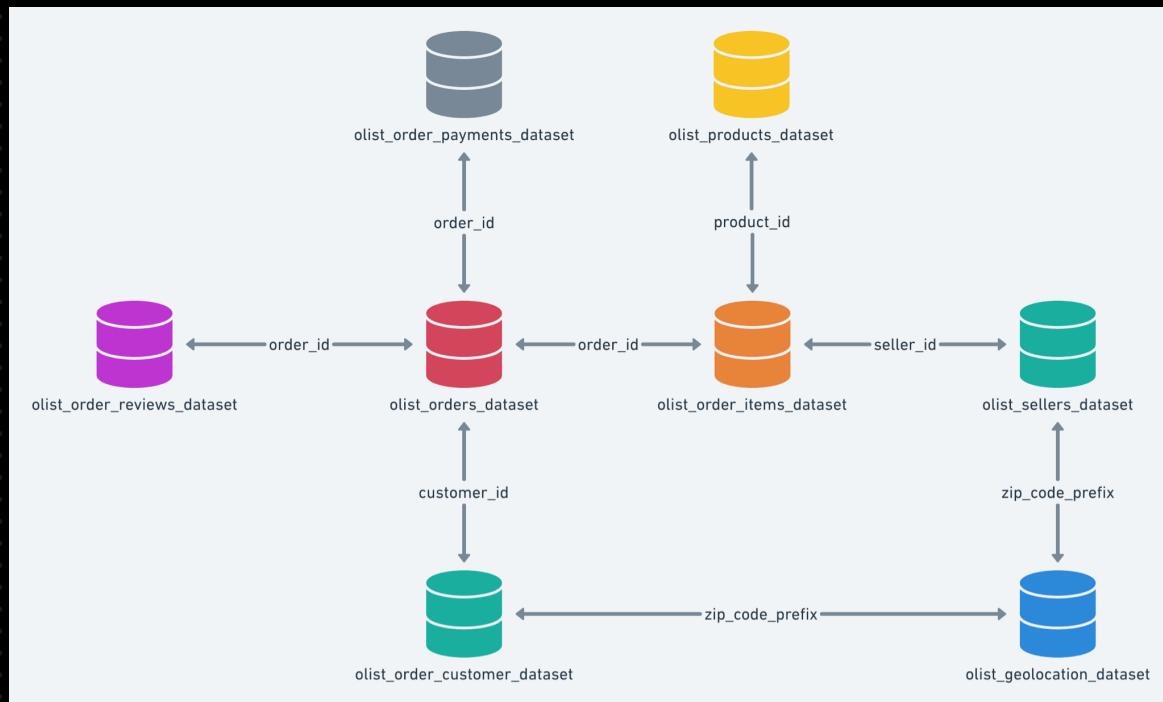


Table Name	Primary Key	Description
<code>olist_orders_dataset.csv</code>	<code>order_id</code>	Contains details of customer orders (timestamps, status, etc.)
<code>olist_order_items_dataset.csv</code>	<code>order_id , product_id , seller_id</code>	Links orders to products and sellers (price, shipping info)
<code>olist_customers_dataset.csv</code>	<code>customer_id</code>	Contains customer info like city, state, and unique ID
<code>olist_order_payments_dataset.csv</code>	<code>order_id</code>	Shows payment methods, installment plans, and values
<code>olist_order_reviews_dataset.csv</code>	<code>order_id , review_id</code>	Stores customer reviews and ratings
<code>olist_products_dataset.csv</code>	<code>product_id</code>	Includes product details like category, weight, and dimensions
<code>olist_sellers_dataset.csv</code>	<code>seller_id</code>	Seller information (location, ZIP code)
<code>olist_geolocation_dataset.csv</code>	<code>geolocation_zip_code_prefix</code>	Provides geographic data to map locations
<code>product_category_name_translation.csv</code>	<code>product_category_name</code>	Translates product categories from Portuguese to English

# Module 1: Data Ingestion and Exploration

We will be following a productionize industry approach.

- Storing data into hdfs/object storage
- Using spark dfl for scalable operation
- exploration will focus on practical business use case.
- we have to make sure to document everything.

# Module 2 - Data Cleaning & Transformation

## Steps in data cleaning and Transformation

1. Identify issues missing values, duplicate, invalid data
  2. Handle missing value drop or fill null values, impute
  3. Standardize formats Convert date time, normalize categorical fields
  4. Data type correction Correct data type if there is none for each column
  5. Deduplication Remove duplicate records
  6. Data Transformation Apply feature engineering
  7. Store cleaned data Data in HDFS, in parquet

# Module 3 - Data Integration & Aggregation

1. Join datasets efficiently
2. Optimizing Joins
3. Aggregations
4. Caching and optimizing queries for performance

\* Spend lots of time  
and understand  
data in depth

