# Batch and Streaming E-commerce Data Pipeline with PySpark

## Overview

This E-commerce data pipeline integrates **batch processing** for historical data and **real-time data streaming** using PySpark. It provides end-to-end ETL (Extract, Transform, Load) capabilities by ingesting data from multiple sources, cleaning and anonymizing sensitive information, enriching and transforming the data for meaningful analysis, and generating business insights through efficient querying. Additionally, the pipeline showcases real-time processing using Spark Streaming to handle incoming data on the fly.

The data sources processed include:

- **Customer Feedback** (CSV)
- **Transactions** (CSV)
- **Inventory** (JSON)

The pipeline ensures a robust, scalable, and optimized solution for handling both batch and streaming data, making it an ideal fit for e-commerce analytics.

## Prerequisites

Before running the pipeline, ensure the following dependencies are installed:

- **PySpark**: `pip install pyspark`
- **Matplotlib and Seaborn**: `pip install matplotlib seaborn` (for visualizations)

---

## Main Steps in the Pipeline:

1. **Data Ingestion**
2. **Data Cleaning and Anonymization**
3. **Data Transformation and Enrichment**
4. **Data Analysis and Reporting**
5. **Performance Optimization**
6. **Data Storage and Visualization**
7. **Real-Time Data Streaming**

# 1. Data Ingestion

The pipeline ingests batch data from CSV and JSON formats:

- **CSV**: For customer feedback and transactions.
- **JSON**: For inventory information.

## Key Operations:

- **Loading Datasets**: The CSV and JSON files are read into Spark DataFrames.
- **Schema Inference**: Automatically infers data types and ensures correct schema assignment to each field in the DataFrames.

By utilizing Spark's distributed capabilities, large datasets are loaded efficiently, allowing the pipeline to handle millions of records.

# 2. Data Cleaning and Anonymization

Raw data is prone to inconsistencies and missing values. This step ensures the data is clean and ready for further transformations.

## Key Operations:

- **Handling Missing Data**: Missing stock levels in the inventory dataset are filled with zeros, ensuring consistent analysis.
- **Anonymization**: To comply with privacy regulations, user IDs in the transactions dataset are anonymized using a secure hashing function, preventing exposure of sensitive customer information.

# 3. Data Transformation and Enrichment

Clean data often needs to be transformed and enriched to make it more suitable for business insights. This step applies necessary transformations and joins multiple datasets for a richer analysis.

## Key Operations:

- **Normalization**: Product names are standardized by converting them to lowercase and removing unnecessary white spaces, making them uniform across records.
- **Category Assignment**: A random category is assigned to each product to facilitate category-level analysis and reporting.

## Data Joining:

- **Join Operation**: A join operation is performed on the `product_id` between the transactions and inventory datasets. This enriches transaction records with details such as product names, categories, and stock levels, enabling deeper analysis across the business.

---

# 4. Data Analysis and Reporting

With the data ready, a variety of analysis tasks are performed using **Spark SQL**. Business intelligence metrics are extracted to provide insights into sales trends, inventory management, and customer behavior.

## Key Operations:

- **Top 10 Best-Selling Products**: Identifies the top-selling products based on total transaction amounts.
- **Monthly Revenue Trends**: Analyzes monthly sales figures to detect trends and seasonality in revenue.
- **Inventory Turnover Rates**: Measures how quickly products are sold out compared to the available stock, giving insight into the demand for various items.

## Additional Insights:

- **Maximum Transaction Date**: Extracts the most recent transaction date, allowing further time-based analysis (e.g., recent sales trends or last-month performance).

---

# 5. Performance Optimization

As the dataset grows, handling it efficiently becomes crucial for performance. To mitigate the overheads of large-scale data processing, several optimization techniques are applied.

## Key Operations:

- **DataFrame Caching**: Frequently accessed DataFrames are cached in memory to speed up repetitive tasks without having to re-read data from the disk.
- **Repartitioning**: DataFrames are partitioned by `product_id` to optimize join operations, ensuring that data is shuffled efficiently across the cluster.
- **Broadcast Join**: When the inventory dataset is small enough, a broadcast join is utilized. This method distributes the smaller DataFrame to all workers, significantly accelerating the join operation with larger datasets like transactions.

These optimizations ensure that the pipeline remains responsive, even when processing large datasets.

---

# 6. Data Storage and Visualization

Once the data analysis is complete, results are stored for future use and visualizations are generated to communicate insights.

## Data Storage:

- **Parquet Format**: The cleaned and enriched data is stored in the Parquet format, which provides efficient data compression and querying for large-scale datasets.

## Visualizations:

- **Top Products by Sales**: A bar chart is used to visually represent the top 10 best-selling products.
- **Monthly Revenue Trends**: A line chart displays monthly revenue trends, helping identify patterns and outliers.
- **Turnover Rate Heatmap**: A heatmap visualizes inventory turnover rates, enabling better inventory management by highlighting products with high or low turnover.

These visualizations make it easier for stakeholders to understand the business performance and take data-driven actions.

---

# 7. Real-Time Data Streaming with Spark Streaming

In addition to batch processing, this pipeline demonstrates real-time data ingestion and processing using **Spark Streaming**.

## Key Operations:

- **Socket Streaming**: A socket connection is created where data is streamed in real-time (using a Python script to simulate incoming JSON data). Spark Streaming ingests this data in real-time from a specified hostname and port.
- **Stream Processing**: As data is ingested, Spark processes each line in real-time, printing the results to the console and optionally storing or aggregating them for further analysis.

This real-time capability is crucial for businesses that need to process and react to data as it arrives, whether for monitoring customer interactions, updating stock levels, or tracking real-time sales.