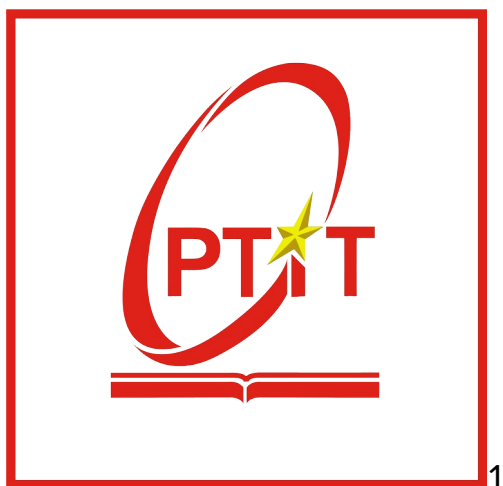


HỌC VIỆN CÔNG NGHỆ BƯU CHÍNH VIỄN THÔNG



Báo cáo hàng tuần

Môn học: Thực tập cơ sở

Giảng viên: Kim Ngọc Bách

Họ và tên: Nguyễn Hữu Phúc

Mã SV: B22DCAT224

Lớp: E22CQCN04-B

Báo cáo tuần 2: Quy trình Pipeline mới trong Dự án Glamira UserFlow Insights

Ngày: 22/03/2025

Dự án: Glamira UserFlow Insights

I. Tổng quan và Lý do thay đổi quy trình Pipeline

Trong dự án **Glamira UserFlow Insights**, mục tiêu là xử lý và phân tích **32 GB** dữ liệu log người dùng từ trang web Glamira để xây dựng các mô hình machine learning (ML) dự đoán hành vi khách hàng. Ban đầu, quy trình pipeline được thiết kế để đưa dữ liệu trực tiếp từ tệp BSON vào Kafka, sau đó thực hiện ETL qua Spark và lưu trữ vào MongoDB. Tuy nhiên, em nhận thấy quy trình này không tốt, nên đã điều chỉnh thành: **Data -> MongoDB -> Kafka -> ETL -> SQL Server -> Machine Learning**. Dưới đây là lý do và lợi ích của quy trình mới:

1. MongoDB làm điểm khởi đầu

- Linh hoạt:** MongoDB là cơ sở dữ liệu NoSQL không yêu cầu schema cố định, rất lý tưởng để nhập dữ liệu thô, không cấu trúc hoặc bán cấu trúc (như tệp BSON 32 GB). Điều này cho phép nhập dữ liệu mà không cần thiết kế schema trước.
- Khả năng mở rộng:** MongoDB xử lý tốt các tập dữ liệu lớn (32 GB không phải là nhỏ) và hỗ trợ mở rộng ngang (horizontal scaling) nếu cần.
- Tốc độ:** MongoDB có tốc độ ghi dữ liệu nhanh, phù hợp để nhập dữ liệu thô một cách nhanh chóng, đóng vai trò như một khu vực trung gian (data lake).
- Lý do bắt đầu ở đây:** MongoDB hoạt động như một hệ thống kiểu "data lake", cho phép thu thập toàn bộ dữ liệu thô trước khi quyết định cách xử lý. Điều này đặc biệt hữu ích khi dữ liệu log người dùng có thể chứa các trường không đồng nhất.

2. Kafka để truyền dữ liệu

- Tách biệt:** Kafka đóng vai trò trung gian, tách biệt giai đoạn nhập dữ liệu (MongoDB) khỏi các bước xử lý sau (ETL và SQL Server). Điều này tránh tắc nghẽn nếu hệ thống phía sau (SQL Server) chậm.
- Xử lý thời gian thực:** Mặc dù 32 GB là một lần tải, Kafka rất phù hợp nếu dữ liệu log người dùng được cập nhật liên tục (ví dụ: hàng ngày hoặc hàng giờ). Nó hỗ trợ xử lý dữ liệu theo thời gian thực hoặc gần thời gian thực.
- Độ tin cậy:** Kafka lưu trữ dữ liệu bền vững, đảm bảo không mất dữ liệu nếu ETL hoặc SQL Server gặp sự cố.
- Lý do thêm bước này:** Kafka tăng tính chịu lỗi và hỗ trợ kiến trúc dữ liệu hiện đại, kết hợp giữa xử lý hàng loạt (batch) và luồng (streaming), phù hợp với các hệ thống phân tích dữ liệu lớn.

3. ETL vào SQL Server

- **Chuyển đổi:** Dữ liệu thô từ MongoDB (dạng JSON) thường cần được làm sạch, chuẩn hóa, hoặc tổng hợp trước khi sử dụng cho ML hoặc phân tích. Quá trình ETL (Extract, Transform, Load) đảm nhiệm việc này, sử dụng Pyspark để thực hiện
- **Lưu trữ có cấu trúc:** SQL Server vượt trội trong việc xử lý dữ liệu quan hệ với schema rõ ràng, rất phù hợp để thực hiện các truy vấn có cấu trúc hoặc nối dữ liệu—những yêu cầu thường thấy trong giai đoạn chuẩn bị dữ liệu cho ML.
- **Lý do dùng SQL Server:** Nhiều quy trình ML yêu cầu dữ liệu dạng bảng (như CSV), và SQL Server tích hợp tốt với các công cụ như Python (qua ODBC) hoặc Power BI cho các bước phân tích sau.

4. Machine Learning là mục tiêu cuối

- **Dữ liệu sẵn sàng:** Các mô hình ML yêu cầu dữ liệu sạch, có cấu trúc, và được kỹ thuật hóa đặc trưng. Dữ liệu từ SQL Server đáp ứng nhu cầu này.
- **Tích hợp:** SQL Server có thể cung cấp dữ liệu trực tiếp cho các framework ML (như pandas trong Python hoặc Azure ML) hoặc lưu trữ tập dữ liệu huấn luyện.
- **Lý do có bước này:** Giá trị cuối cùng của 32 GB dữ liệu log nằm ở việc tạo ra các hiểu biết hoặc dự đoán (ví dụ: dự đoán khả năng mua hàng của khách hàng trên Glamira). ML là bước tự nhiên để đạt được mục tiêu này.

II. Phân biệt JSON và BSON

Dữ liệu log người dùng trong dự án được lưu trữ dưới dạng BSON (trong tệp summary.bson) và cần chuyển đổi sang JSON để đẩy lên Kafka. Dưới đây là sự khác biệt giữa JSON và BSON:

1. Cấu trúc dữ liệu

- **JSON:**
 - Là định dạng văn bản (text-based), cấu trúc dạng cặp khóa-giá trị (key-value pairs).
 - Được thiết kế để dễ đọc và hiểu đối với con người.
 - Hỗ trợ các kiểu dữ liệu cơ bản: đối tượng, mảng, chuỗi, số, giá trị boolean.
- **BSON:**
 - Là định dạng nhị phân (binary-based), có cấu trúc tương tự JSON nhưng hỗ trợ nhiều kiểu dữ liệu phức tạp hơn.
 - Hỗ trợ các kiểu dữ liệu đặc biệt như ObjectId, Date, BinData (dữ liệu nhị phân), ngoài các kiểu cơ bản.
 - Không dễ đọc bằng mắt người vì là dữ liệu nhị phân.

2. Kích thước và hiệu suất

- **JSON:**
 - Kích thước nhẹ vì là văn bản thuần, nhưng không hỗ trợ các kiểu dữ liệu đặc biệt như BSON.
 - Đọc/ghi dữ liệu dạng văn bản tốn nhiều tài nguyên tính toán hơn so với dữ liệu nhị phân.
- **BSON:**
 - Kích thước lớn hơn JSON một chút do lưu trữ thêm metadata, nhưng cho phép truy vấn và thao tác dữ liệu nhanh hơn trong MongoDB.
 - Đọc/ghi nhanh hơn vì không cần chuyển đổi qua lại giữa văn bản và nhị phân.

3. Dễ đọc và sử dụng

- **JSON:**
 - Dễ đọc, dễ hiểu đối với con người.
 - Dễ sử dụng, được hỗ trợ bởi hầu hết các ngôn ngữ lập trình.
- **BSON:**
 - Khó đọc vì là dữ liệu nhị phân.
 - Chủ yếu dùng trong các ứng dụng cần tốc độ cao hoặc lưu trữ dữ liệu phức tạp (như MongoDB).

4. Ứng dụng phổ biến

- **JSON:**
 - Dùng để trao đổi dữ liệu giữa các ứng dụng web và API (HTTP, RESTful APIs).
- **BSON:**
 - Dùng trong cơ sở dữ liệu NoSQL như MongoDB, phù hợp với dữ liệu phức tạp và truy vấn hiệu suất cao.

5. Tính tương thích

- **JSON:**
 - Tương thích với hầu hết các hệ thống và ngôn ngữ lập trình, dễ trao đổi giữa các hệ thống (ví dụ: REST APIs).
- **BSON:**
 - Hạn chế hơn, chủ yếu dùng trong MongoDB hoặc các ứng dụng tương thích với BSON.

Ứng dụng trong Glamira:

- Tập **summary.bson (32 GB)** được lưu trữ dưới dạng **BSON** vì nó chứa các kiểu dữ liệu đặc biệt của MongoDB (như ObjectId, Date).
 - Để đẩy lên Kafka, cần chuyển sang **JSON** vì Kafka thường xử lý dữ liệu dạng văn bản (JSON) để đảm bảo tính tương thích với các hệ thống sau (Spark, SQL Server).
-

III. Đưa dữ liệu BSON vào MongoDB

Mở terminal, sử dụng lệnh: mongorestore --db Glamira --collection summary
/home/nguyenphuc/Documents/GlamiraUserFlowInsightsProject/glamira_ubl_oct2019_nov2019/dump/countly/summary.bson

2025-03-18T20:39:18.938+0700	[.....]	Glamira.summary	491MB/31.2GB	(1.5%)
2025-03-18T20:39:21.939+0700	[.....]	Glamira.summary	984MB/31.2GB	(3.1%)
2025-03-18T20:39:24.938+0700	[#.....]	Glamira.summary	1.45GB/31.2GB	(4.6%)
2025-03-18T20:39:27.938+0700	[#.....]	Glamira.summary	1.92GB/31.2GB	(6.2%)
2025-03-18T20:39:30.938+0700	[#.....]	Glamira.summary	2.39GB/31.2GB	(7.6%)
2025-03-18T20:39:33.938+0700	[##.....]	Glamira.summary	2.87GB/31.2GB	(9.2%)
2025-03-18T20:39:36.938+0700	[##.....]	Glamira.summary	3.35GB/31.2GB	(10.7%)
2025-03-18T20:39:39.938+0700	[##.....]	Glamira.summary	3.82GB/31.2GB	(12.2%)
2025-03-18T20:39:42.938+0700	[###.....]	Glamira.summary	4.26GB/31.2GB	(13.6%)
2025-03-18T20:39:45.939+0700	[###.....]	Glamira.summary	4.72GB/31.2GB	(15.1%)
2025-03-18T20:39:48.938+0700	[###.....]	Glamira.summary	5.20GB/31.2GB	(16.6%)
2025-03-18T20:39:51.939+0700	[###.....]	Glamira.summary	5.66GB/31.2GB	(18.1%)
2025-03-18T20:39:54.938+0700	[####.....]	Glamira.summary	6.14GB/31.2GB	(19.7%)
2025-03-18T20:39:57.939+0700	[####.....]	Glamira.summary	6.61GB/31.2GB	(21.2%)
2025-03-18T20:40:00.938+0700	[####.....]	Glamira.summary	7.09GB/31.2GB	(22.7%)
2025-03-18T20:40:03.938+0700	[####.....]	Glamira.summary	7.55GB/31.2GB	(24.2%)
2025-03-18T20:40:06.938+0700	[####.....]	Glamira.summary	8.01GB/31.2GB	(25.7%)
2025-03-18T20:40:09.938+0700	[####.....]	Glamira.summary	8.48GB/31.2GB	(27.2%)
2025-03-18T20:40:12.938+0700	[####.....]	Glamira.summary	8.94GB/31.2GB	(28.6%)
2025-03-18T20:40:15.939+0700	[####.....]	Glamira.summary	9.38GB/31.2GB	(30.0%)
2025-03-18T20:40:18.938+0700	[####.....]	Glamira.summary	9.81GB/31.2GB	(31.4%)
2025-03-18T20:40:21.938+0700	[####.....]	Glamira.summary	10.2GB/31.2GB	(32.8%)
2025-03-18T20:40:24.939+0700	[####.....]	Glamira.summary	10.7GB/31.2GB	(34.1%)
2025-03-18T20:40:27.938+0700	[####.....]	Glamira.summary	11.1GB/31.2GB	(35.5%)
2025-03-18T20:40:30.938+0700	[####.....]	Glamira.summary	11.5GB/31.2GB	(36.9%)
2025-03-18T20:40:33.938+0700	[####.....]	Glamira.summary	12.0GB/31.2GB	(38.4%)
2025-03-18T20:40:36.938+0700	[####.....]	Glamira.summary	12.5GB/31.2GB	(39.9%)
2025-03-18T20:40:39.938+0700	[####.....]	Glamira.summary	12.9GB/31.2GB	(41.4%)
2025-03-18T20:40:42.939+0700	[####.....]	Glamira.summary	13.4GB/31.2GB	(42.8%)
2025-03-18T20:40:45.938+0700	[####.....]	Glamira.summary	13.8GB/31.2GB	(44.3%)
2025-03-18T20:40:48.939+0700	[####.....]	Glamira.summary	14.3GB/31.2GB	(45.8%)
2025-03-18T20:40:51.939+0700	[####.....]	Glamira.summary	14.8GB/31.2GB	(47.3%)
2025-03-18T20:40:54.938+0700	[####.....]	Glamira.summary	15.2GB/31.2GB	(48.8%)
2025-03-18T20:40:57.938+0700	[####.....]	Glamira.summary	15.7GB/31.2GB	(50.3%)
2025-03-18T20:41:00.938+0700	[####.....]	Glamira.summary	16.2GB/31.2GB	(51.7%)
2025-03-18T20:41:03.938+0700	[####.....]	Glamira.summary	16.6GB/31.2GB	(53.2%)
2025-03-18T20:41:06.939+0700	[####.....]	Glamira.summary	17.1GB/31.2GB	(54.7%)
2025-03-18T20:41:09.939+0700	[####.....]	Glamira.summary	17.6GB/31.2GB	(56.2%)
2025-03-18T20:41:12.939+0700	[####.....]	Glamira.summary	18.0GB/31.2GB	(57.6%)
2025-03-18T20:41:15.938+0700	[####.....]	Glamira.summary	18.4GB/31.2GB	(59.1%)
2025-03-18T20:41:18.939+0700	[####.....]	Glamira.summary	18.9GB/31.2GB	(60.6%)
2025-03-18T20:41:21.938+0700	[####.....]	Glamira.summary	19.4GB/31.2GB	(62.1%)
2025-03-18T20:41:24.938+0700	[####.....]	Glamira.summary	19.9GB/31.2GB	(63.6%)
2025-03-18T20:41:27.938+0700	[####.....]	Glamira.summary	20.3GB/31.2GB	(65.1%)
2025-03-18T20:41:30.938+0700	[####.....]	Glamira.summary	20.8GB/31.2GB	(66.6%)
2025-03-18T20:41:33.938+0700	[####.....]	Glamira.summary	21.3GB/31.2GB	(68.1%)
2025-03-18T20:41:36.938+0700	[####.....]	Glamira.summary	21.7GB/31.2GB	(69.6%)
2025-03-18T20:41:39.938+0700	[####.....]	Glamira.summary	22.2GB/31.2GB	(71.0%)
2025-03-18T20:41:42.938+0700	[####.....]	Glamira.summary	22.6GB/31.2GB	(72.5%)
2025-03-18T20:41:45.938+0700	[####.....]	Glamira.summary	23.1GB/31.2GB	(74.0%)
2025-03-18T20:41:48.938+0700	[####.....]	Glamira.summary	23.6GB/31.2GB	(75.5%)
2025-03-18T20:41:51.939+0700	[####.....]	Glamira.summary	24.0GB/31.2GB	(77.0%)
2025-03-18T20:41:54.939+0700	[####.....]	Glamira.summary	24.5GB/31.2GB	(78.4%)
2025-03-18T20:41:57.938+0700	[####.....]	Glamira.summary	24.9GB/31.2GB	(79.9%)
2025-03-18T20:42:00.938+0700	[####.....]	Glamira.summary	25.4GB/31.2GB	(81.3%)
2025-03-18T20:42:03.938+0700	[####.....]	Glamira.summary	25.8GB/31.2GB	(82.7%)
2025-03-18T20:42:06.938+0700	[####.....]	Glamira.summary	26.3GB/31.2GB	(84.2%)
2025-03-18T20:42:09.938+0700	[####.....]	Glamira.summary	26.8GB/31.2GB	(85.7%)
2025-03-18T20:42:12.938+0700	[####.....]	Glamira.summary	27.2GB/31.2GB	(87.1%)
2025-03-18T20:42:15.939+0700	[####.....]	Glamira.summary	27.7GB/31.2GB	(88.6%)
2025-03-18T20:42:18.938+0700	[####.....]	Glamira.summary	28.1GB/31.2GB	(90.0%)
2025-03-18T20:42:21.938+0700	[####.....]	Glamira.summary	28.6GB/31.2GB	(91.5%)
2025-03-18T20:42:24.939+0700	[####.....]	Glamira.summary	29.1GB/31.2GB	(93.0%)
2025-03-18T20:42:27.938+0700	[####.....]	Glamira.summary	29.4GB/31.2GB	(94.1%)
2025-03-18T20:42:30.938+0700	[####.....]	Glamira.summary	29.8GB/31.2GB	(95.4%)
2025-03-18T20:42:33.938+0700	[####.....]	Glamira.summary	30.2GB/31.2GB	(96.8%)
2025-03-18T20:42:36.938+0700	[####.....]	Glamira.summary	30.6GB/31.2GB	(98.1%)
2025-03-18T20:42:39.938+0700	[####.....]	Glamira.summary	31.1GB/31.2GB	(99.4%)
2025-03-18T20:42:41.287+0700	[####.....]	Glamira.summary	31.2GB/31.2GB	(100.0%)

```
2025-03-18T20:42:41.288+0700 finished restoring Glamira.summary (41432473 documents, 0 failures)
2025-03-18T20:42:41.289+0700 restoring indexes for collection Glamira.summary from metadata
2025-03-18T20:42:41.289+0700 index: &idx.IndexDocument{Options:primitive.M{"name":"time_stamp_1", "ns":"co
2025-03-18T20:42:41.289+0700 index: &idx.IndexDocument{Options:primitive.M{"name":"device_id_1", "ns":"cou
2025-03-18T20:44:54.320+0700 41432473 document(s) restored successfully. 0 document(s) failed to restore.
```