# Amazon Book Reviews Recommendation System

Big Data Pipeline on Google Cloud Platform - Final Project Presentation

# Introduction

## Dataset(from kaggle)

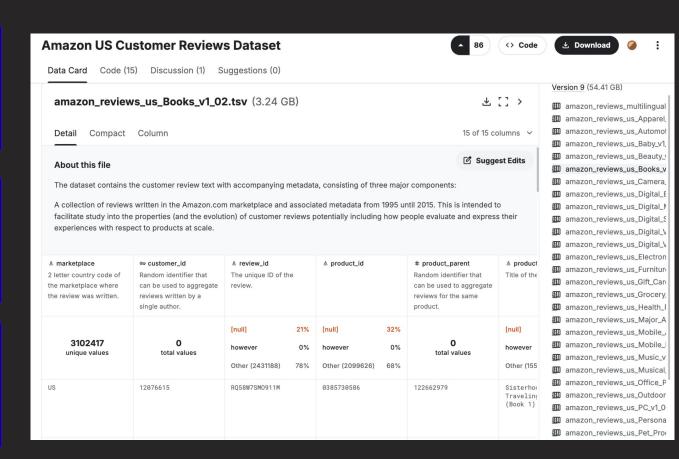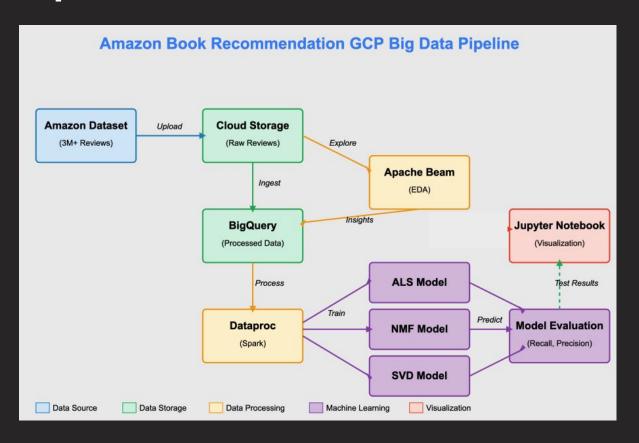3.1M reviews, 1.5M users, 779K products with ratings and metadata. 3 GB.

## Problem Statement

Build a scalable recommendation system for Amazon book reviews.

## Objective

Compare different recommendation algorithms on large-scale data using GCP.

# Pipeline Architecture & GCP Services



**Data Ingestion**

Google Cloud Storage

**Data Processing**

Dataflow (Apache Beam),
BigQuery

**Storage & Query**

BigQuery
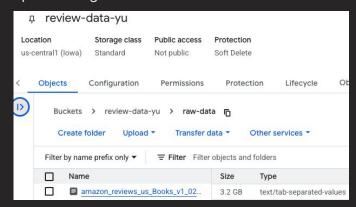
**Machine Learning**

PySpark ML, scikit-learn

# GCP Services Used

**Google Cloud Storage**

Stored raw TSV data and temporary storage for processing.



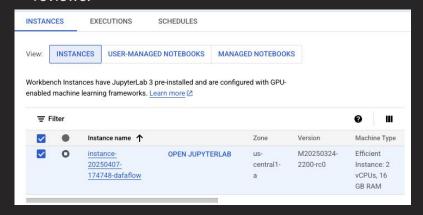**Dataflow**

Distributed Apache Beam pipelines for EDA on 3.1M reviews.



**BigQuery**

Structured storage, SQL queries, aggregations, and train/test splits.

**Dataproc Workbench**

Managed notebooks running PySpark and recommendation models.

# Pipeline Implementation - Part 1: EDA with Beam

**Data Loading**
Loaded TSV data from Cloud Storage into Dataflow.

**Distributed Processing**
Apache Beam pipeline analyzed ratings and user activity.

**Key Metrics**
3,105,520 records, 1,502,380 users, 779,733 products, avg rating 4.18.

# Part 2: Data Preprocessing

**Data Loading & BigQuery Setup:** Imports and organizes raw data.

**High-Frequency Filtering:** Focuses on active users and products.

**ID Mapping & Processing:** Creates numerical IDs for modeling.

**Train/Test Data Split:** Prepares data for model development.

**Validation:** Ensures data integrity throughout the process.

```
===== DATA PROCESSING VALIDATION =====
Stage | Count | Unique Users | Unique Products
--------------------------------------------------
raw | 3105520 | 1502380 | 779733
filtered | 143479 | 6238 | 11079
processed | 143479 | 6238 | 11079
train | 117283 | 6238 | 10957
test | 26196 | 5415 | 8302

===== INDEX MAPPING VALIDATION =====
User index range: 0 to 6237 (6238 unique)
Product index range: 0 to 11078 (11079 unique)
```

# Part 3: Recommendation Models

**ALS**

Alternating Least Squares using Spark MLlib.

Distributed matrix factorization algorithm.

**SVD**

Singular Value Decomposition using scikit-learn.

Classical matrix factorization technique.

**NMF**

Non-negative Matrix Factorization using scikit-learn.

Optimized for sparse matrix operations.

```
--- Model Comparison Summary
   Model   Recall@20   NDCG@20
0  ALS     0.023591    0.013878
1  SVD     0.070978    0.044329
2  NMF     0.042439    0.027186
```

# Challenges & Solutions

### Scale of the Dataset

**Problem:** 3.1M reviews too large for single-machine processing.

**Solution:** Used Dataflow for distributed EDA and BigQuery for filtering.

### Sparse User-Item Matrix

**Problem:** Very sparse interaction matrix.

**Solution:** Filtered for users/products with ≥30 interactions.

### Debugging

**Problem:** Hard to trace issues in complex pipeline logic.

**Solution:** Used Jupyter Notebook as the main environment to quickly test and debug.

# Future Work

- Submit Spark to Dataproc for distributed training

- Process reviews with NLP

- Hybrid recommendation models, fine tuning

- Deploy as web service with API access

THANK YOU