

Team 27

Books recommendation system

Team members:

- Arsenii Pavlov
- Ivan Beltsov
- Andrei Pavlov
- Timofei Ivlenkov
- Bulat Latypov

Introduction

In the digital age of book publishing, understanding reader preferences through ratings is crucial for platforms and publishers. This project leverages machine learning to predict book ratings (1–5 stars) using a rich dataset from Amazon containing 3 million user reviews across 212,404 unique books (1996–2014).

Business Objectives

We aim to enhance user experience on discovering new books by providing personalized recommendations by predicting his review of the book. Also we want to increase platform engagement and drive higher user retention and interaction by improving the relevance of suggestions. Translate accurate recommendations into increased book purchases or downloads directly benefiting publishers and online retailers.

Data Description

The dataset employed in this project consists of two primary files—`Books_rating.csv` and `books_data.csv`—which together form a rich source of information for building a robust and personalized book recommendation system. These datasets are complementary: one captures user-generated feedback and behavior, while the other provides detailed metadata about the books themselves.

1. `Books_rating.csv`

This file contains over 3 million user reviews for books, representing a diverse and extensive collection of reader feedback accumulated over nearly two decades (1996–2014). Each row in the dataset represents a single review submitted by a user for a particular book. The structure of the data includes:

Id: A unique identifier for the book being reviewed.

Title: The name of the book.

Price: The listed price of the book, when available.

User_id: A unique identifier for the user who submitted the review.

profileName: The name (or pseudonym) of the user.

review/helpfulness: A metric indicating how helpful other users found the review, expressed as a ratio

review/score: The numeric rating given to the book, ranging from 0 to 5.

review/time: A Unix timestamp indicating when the review was posted.

review/summary: A brief summary or title of the review.

review/text: The full text of the review, which often provides valuable insights into the reader's opinions and preferences.

This dataset not only provides a basis for collaborative filtering but also allows for sentiment analysis, trend detection over time, and the identification of user profiles based on review behavior.

2. `books_data.csv`

This file offers detailed metadata for 212,404 unique books, most of which are referenced in the reviews dataset. The metadata includes bibliographic and descriptive fields that are essential for content-based recommendation methods. Key attributes include:

Title: The book's title, which serves as a key for merging with the reviews.

description: A summary or synopsis of the book's content, useful for NLP-based content analysis.

authors: The author(s) of the book.

image: A URL to the book's cover image.

previewLink: A link to a Google Books preview of the book, when available.

publisher: The publishing company.

publishedDate: The original publication date.

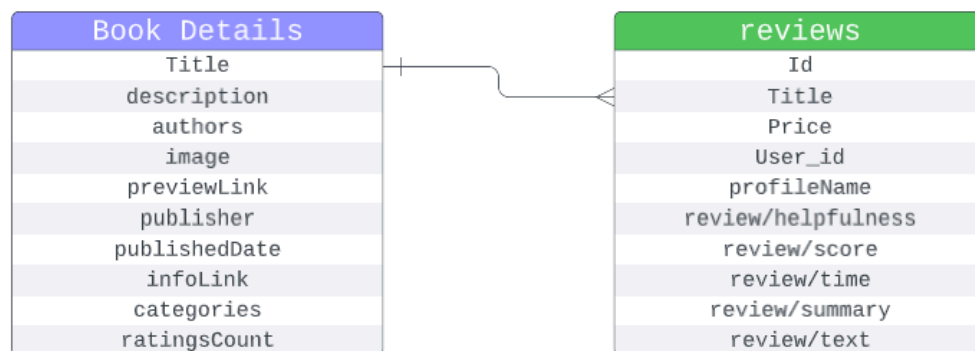
infoLink: A Google Books link providing more detailed information about the book.

categories: The genre or subject classification of the book (e.g., fiction, history, science).

ratingsCount: The total number of ratings the book has received, which can be used as a proxy for popularity.

The inclusion of rich descriptive fields enables deeper user-book profiling and supports hybrid recommendation approaches by combining collaborative and content-based methods.

The data can be described as diagram:



Source: <https://www.kaggle.com/datasets/mohamedbakhet/amazon-books-reviews>

Architecture of data pipeline

Whole pipeline contains 3 stages and can be described as diagram on the right side of page

1. Stage 1

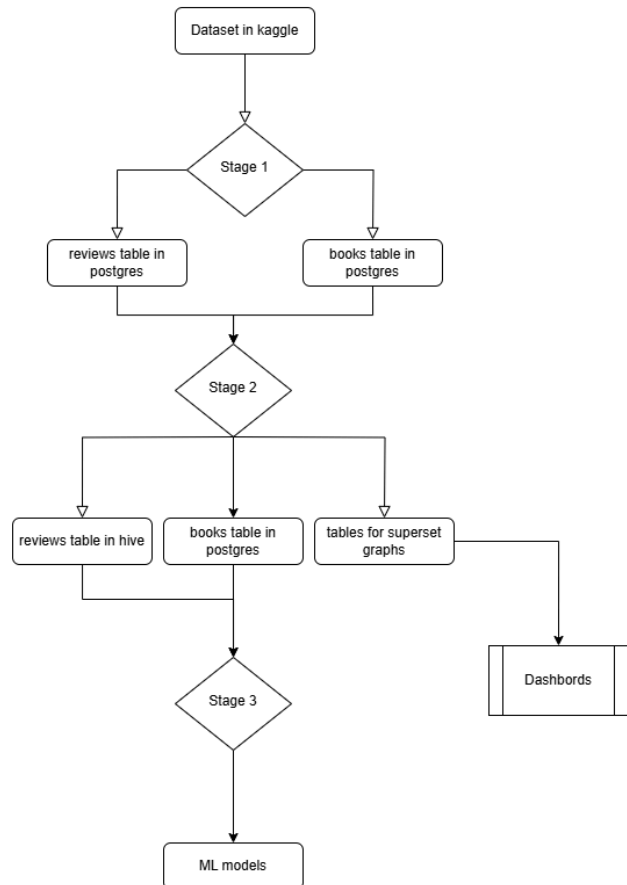
The stage downloads data from kaggle, clean it and upload to postgres tables

2. Stage 2

The stage analyses data which is on the postgres tables, build tables for dashboards and upload data to hive

3. Stage 3

The stage preprocess data for training and train ML models to for recommendation systems



Data preparation

Our data preparation started downloading a dataset for kaggle and after clearing it was uploaded to postgres database.

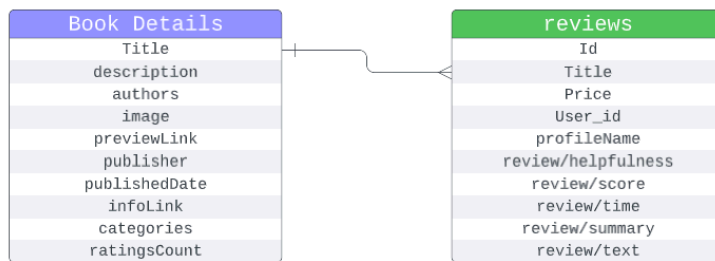
Main parts of clearing dataset:

- 1) Remove lines where critical columns contains NULL
- 2) Fix author columns to make it proper array
- 3) Delete broken symbols for title and description

At the end we get data as postgres table:

books table:

Query											Query History											Scratch Pad		X																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																								
1															SELECT * FROM books;																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																	
Data Output															Messages											Notifications																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																						



Source: <https://www.kaggle.com/datasets/mohamedbakhmet/amazon-books-reviews>

Hive tables creating:

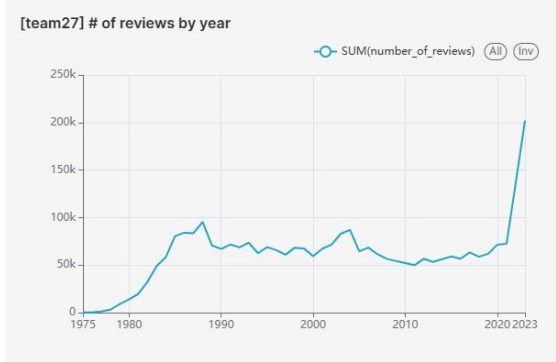
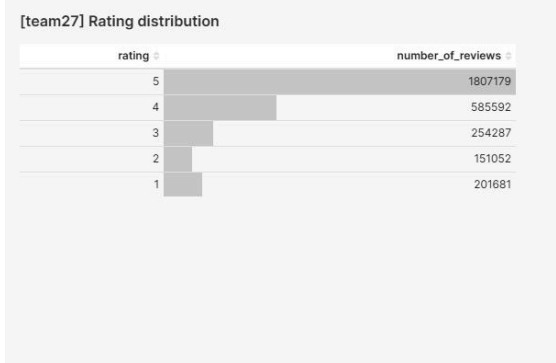
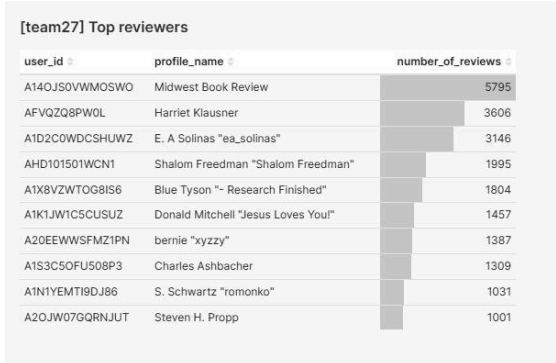
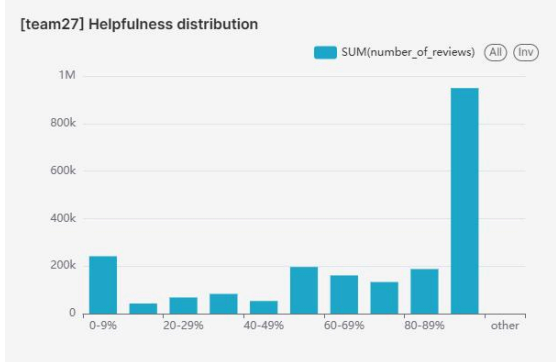
During the stage we successfully created external Hive tables for the dataset imported via Sqoop in Stage I and ensured proper data types and constraints were applied. The tables were partitioned to ensure that we will be able to efficiently use them in our analysis work.

Data analysis

Exploratory Data Analysis:

We conducted in-depth analysis using HiveQL queries to extract valuable insights. We extracted 6 insights using 6 HQL queries:

Title	Query file	Interpretation	Superset Chart																						
Rating Distribution: Analyzed the frequency of each rating score (1–5).	sql/q1.hql	This chart is showing the distribution of reviews grades. We can see that people tend to give positive reviews and the least popular rating is 2.	<div><p>[team27] Rating distribution</p><table><thead><tr><th>rating</th><th>number_of_reviews</th></tr></thead><tbody><tr><td>5</td><td>180719</td></tr><tr><td>4</td><td>58559</td></tr><tr><td>3</td><td>25428</td></tr><tr><td>2</td><td>15105</td></tr><tr><td>1</td><td>2016</td></tr></tbody></table></div>	rating	number_of_reviews	5	180719	4	58559	3	25428	2	15105	1	2016										
rating	number_of_reviews																								
5	180719																								
4	58559																								
3	25428																								
2	15105																								
1	2016																								
Top Reviewed Books: Identified the 10 most-reviewed books	sql/q2.hql	Chart is showing top books sorted by reviews in decending order. From the chart we can see most popular books. Chart gives intuinion that most reviews are focused on most popular books.	<div><p>[team27] Top reviewed</p><table><thead><tr><th>book_title</th><th>reviews_cnt</th></tr></thead><tbody><tr><td>The Hobbit</td><td>22023</td></tr><tr><td>Pride and Prejudice</td><td>20371</td></tr><tr><td>Atlas Shrugged</td><td>12513</td></tr><tr><td>Wuthering Heights</td><td>10780</td></tr><tr><td>The Giver</td><td>7644</td></tr><tr><td>Great Expectations</td><td>7421</td></tr><tr><td>Harry Potter and The Sorcerer's Stone</td><td>6796</td></tr><tr><td>Of Mice and Men</td><td>6728</td></tr><tr><td>Brave New World</td><td>6312</td></tr><tr><td>Mere Christianity</td><td>6053</td></tr></tbody></table></div>	book_title	reviews_cnt	The Hobbit	22023	Pride and Prejudice	20371	Atlas Shrugged	12513	Wuthering Heights	10780	The Giver	7644	Great Expectations	7421	Harry Potter and The Sorcerer's Stone	6796	Of Mice and Men	6728	Brave New World	6312	Mere Christianity	6053
book_title	reviews_cnt																								
The Hobbit	22023																								
Pride and Prejudice	20371																								
Atlas Shrugged	12513																								
Wuthering Heights	10780																								
The Giver	7644																								
Great Expectations	7421																								
Harry Potter and The Sorcerer's Stone	6796																								
Of Mice and Men	6728																								
Brave New World	6312																								
Mere Christianity	6053																								

Review Trends Over Time: Examined yearly review counts.	sql/q3.hql	Chart shows the number of reviews per book's year. Slowly rising till 1990, then slightly falling and rapidly rising till today																																		
Publisher Performance: Ranked publishers by average review score (with >20 reviews).	sql/q4.hql	Chart shows top 20 publishers by average ratings. Publishers with at least 20 reviews were chosen.	 <table><thead><tr><th>rating</th><th>number_of_reviews</th></tr></thead><tbody><tr><td>5</td><td>180719</td></tr><tr><td>4</td><td>585592</td></tr><tr><td>3</td><td>254287</td></tr><tr><td>2</td><td>151052</td></tr><tr><td>1</td><td>201681</td></tr></tbody></table>	rating	number_of_reviews	5	180719	4	585592	3	254287	2	151052	1	201681																					
rating	number_of_reviews																																			
5	180719																																			
4	585592																																			
3	254287																																			
2	151052																																			
1	201681																																			
Most Active Users: Listed users with the highest number of reviews.	sql/q5.hql	Chart is showing users with highest amount of given reviews	 <table><thead><tr><th>user_id</th><th>profile_name</th><th>number_of_reviews</th></tr></thead><tbody><tr><td>A14OJS0VWMOSWO</td><td>Midwest Book Review</td><td>5795</td></tr><tr><td>AFVQZQ8PWOL</td><td>Harriet Klausner</td><td>3606</td></tr><tr><td>A1D2COWDCSHUWZ</td><td>E. A Solinas "ea_solinas"</td><td>3146</td></tr><tr><td>AHD101501WCN1</td><td>Shalom Freedman "Shalom Freedman"</td><td>1995</td></tr><tr><td>ATX8VZWTOG8IS6</td><td>Blue Tyson "~ Research Finished"</td><td>1804</td></tr><tr><td>A1K1JW1C5CUSUZ</td><td>Donald Mitchell "Jesus Loves You!"</td><td>1457</td></tr><tr><td>A20EEWWSFMZ1PN</td><td>bernie "xyzyzy"</td><td>1387</td></tr><tr><td>A1S3C5OFU508P3</td><td>Charles Ashbacher</td><td>1309</td></tr><tr><td>A1N1YEMTI9DJ86</td><td>S. Schwartz "romonko"</td><td>1031</td></tr><tr><td>A2OJW07GQRNJUT</td><td>Steven H. Propp</td><td>1001</td></tr></tbody></table>	user_id	profile_name	number_of_reviews	A14OJS0VWMOSWO	Midwest Book Review	5795	AFVQZQ8PWOL	Harriet Klausner	3606	A1D2COWDCSHUWZ	E. A Solinas "ea_solinas"	3146	AHD101501WCN1	Shalom Freedman "Shalom Freedman"	1995	ATX8VZWTOG8IS6	Blue Tyson "~ Research Finished"	1804	A1K1JW1C5CUSUZ	Donald Mitchell "Jesus Loves You!"	1457	A20EEWWSFMZ1PN	bernie "xyzyzy"	1387	A1S3C5OFU508P3	Charles Ashbacher	1309	A1N1YEMTI9DJ86	S. Schwartz "romonko"	1031	A2OJW07GQRNJUT	Steven H. Propp	1001
user_id	profile_name	number_of_reviews																																		
A14OJS0VWMOSWO	Midwest Book Review	5795																																		
AFVQZQ8PWOL	Harriet Klausner	3606																																		
A1D2COWDCSHUWZ	E. A Solinas "ea_solinas"	3146																																		
AHD101501WCN1	Shalom Freedman "Shalom Freedman"	1995																																		
ATX8VZWTOG8IS6	Blue Tyson "~ Research Finished"	1804																																		
A1K1JW1C5CUSUZ	Donald Mitchell "Jesus Loves You!"	1457																																		
A20EEWWSFMZ1PN	bernie "xyzyzy"	1387																																		
A1S3C5OFU508P3	Charles Ashbacher	1309																																		
A1N1YEMTI9DJ86	S. Schwartz "romonko"	1031																																		
A2OJW07GQRNJUT	Steven H. Propp	1001																																		
Helpfulness Analysis: Categorized reviews by helpfulness percentage (e.g., "70–79%").	sql/q6.hql	Helpfulness distribution Chart is showing distribution of the helpfulness percentage of reviews. We can conclude that most of the reviews are helpful for other users.																																		

ML modeling

Feature extraction and data preprocessing

Key Steps:

1. Data Loading

Data loaded from Hive tables:

books (book metadata: title, author, description, etc.)

reviews_processed (user reviews: ratings, text, timestamps).

2. Data Joining and Cleaning

2.1. Inner join on book_title to merge reviews with book details.

2.2. Missing values handled:

Rows with null user_id or book_id removed (critical for recommendation models).

Null text fields (summary, review_text, description) filled with empty strings.

ratings_count imputed using review counts where missing.

3. Feature Engineering

3.1. Helpfulness Score:

Original format (X/Y) split into helpful_yes and helpful_total.

Wilson score (95% confidence) computed for robustness.

3.2. Temporal Features:

Unix timestamps converted to datetime.

Year, month, day extracted and cyclically encoded (sin/cos) for periodicity.

3.3. Categorical Encoding:

Authors & categories parsed from JSON/strings (first entry kept).

StringIndexer applied to authors, categories, user IDs, and book IDs.

4. Text Processing (TF-IDF)

Fields processed: description, review_text, summary, title, publisher.

Pipeline: Tokenization → Stopword removal → TF-IDF (300 features per field).

5. Final Output

5.1. Train/test split (70/30).

5.2. Data saved as JSON in HDFS (project/data/train, project/data/test).

All original features that were processed by encoding, indexed or used for TF-IDF and will not be used for ML part are replaced with pre-processed columns.

Training and fine-tuning

Four distinct regression algorithms from Spark MLlib were selected: the RandomForestRegressor, LogisticRegression, GBTRRegressor, FMRegressor.

For RandomForestRegressor and LogisticRegression, we configured the featuresCol parameter to point to our assembled feature vector, 'features', and the labelCol parameter to 'label', representing the integer score (1-5). For Logistic Regression, we explicitly set family="multinomial" to handle the multi-class nature of our problem.

For both GBTRRegressor and FMRegressor, features_col has 2 types of features:
tfidf_cols: "description_tfidf", "summary_tfidf", "title_tfidf",
 "review_text_tfidf", "publisher_tfidf"

numerical_cols: "author_idx", "book_idx", "category_idx", "helpfulness_wilson",
 "published_day_encoded_cos", "published_day_encoded_sin",
 "published_month_encoded_cos", "published_month_encoded_sin",
 "published_year", "ratings_count", "review_count",
 "review_day_encoded_cos", "review_day_encoded_sin",
 "review_month_encoded_cos", "review_month_encoded_sin",
 "review_year", "user_idx"

To find the optimal hyperparameters for each model, we employed Spark's CrossValidator. This process involves:

1. Defining a ParamGrid for each estimator, specifying the hyperparameters to tune and the range of values to explore.
For RandomForestClassifier, we tuned: numTrees (number of trees), maxDepth (maximum tree depth) and featureSubsetStrategy (number of features to consider for splits at each tree node).

For GBTRRegressor, we tuned maxDepth (depth of each decision tree), featureSubsetStrategy (strategy for selecting a subset of features at each split) and subsamplingRate (Fraction of the training data used for each iteration.)

For FMRegressor, we tuned factorSize (Dimensionality of the latent factors), initStd (Standard deviation used to initialize the factor weights) and regParam (Regularization strength).

For Logistic Regression, we tuned regParam (regularization parameter), elasticNetParam (mix of L1 and L2 regularization) and threshold (for classes)

This resulted in 27 parameter combinations for each model.

2. Using a MulticlassClassificationEvaluator configured with metricName="accuracy" as the evaluation criterion during cross-validation for Logistic Regression and RandomForestClassifier, and **RegressionEvaluator** with RMSE as metric was used for FMRegressor and GBTRRegressor.

3. Running the CrossValidator with numFolds set to 3. The validator performs 3-fold cross-validation for each parameter combination, trains sub-models on training folds, evaluates them on validation folds using the accuracy, averages the metric across folds, and finally selects the parameter set that yielded the best results.

4. The best model corresponding to the optimal parameter set is then trained on the entire training dataset (train_df_prepared_classification).

Final hyperparameters:

RandomForestClassifier:

Hyperparameters: numTrees=[10, 17, 25], maxDepth=[5, 7, 10], featureSubsetStrategy=["log2", "sqrt", "onethird"].

Best Model: 10 trees, max depth 7, "sqrt".

LogisticRegression:

Hyperparameters: regParam=[0.05, 0.17, 0.25], elasticNetParam=[0.1, 0.25, 0.5], threshold=[0.3, 0.5, 0.7]

Best Model: regParam=0.05, elasticNetParam=0.1, threshold=0.5

FMRegressor:

Hyperparameters: factorSize=[4, 6, 8], initStd=[0.01, 0.1, 1.0], regParam=[0.01, 0.1, 1.0]

Best Model: factorSize=8, initStd=1.0, regParam=0.01

GBTRRegressor:

Hyperparameters: maxDepth=[4, 6, 8], featureSubsetStrategy=["log2", "sqrt", "onethird"], subsamplingRate=[0.6, 0.8, 1.0]

Best Model: maxDepth=8, featureSubsetStrategy="sqrt", subsamplingRate=0.8

Evaluation

Once the best models were identified and trained by the CrossValidator, we proceeded to evaluate their performance on a completely separate portion of our data: the test set (test_df_prepared_classification)

For them, we got those **metrics results**:

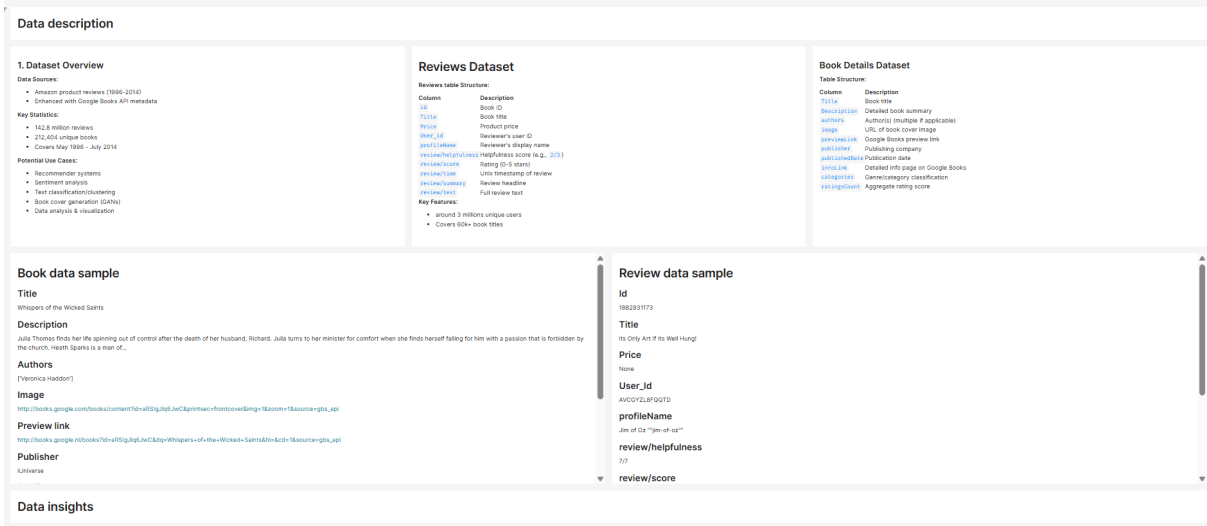
model	RMSE	R2
RandomForestRegressor	1.4162094562456848	-0.43800524504817395
LogisticRegression	1.422036996664622	-0.4498640489738863
GBTRegressor	1.428302134567891	-0.443209812345678
FMRRegressor	1.449876542309876	-0.451763920123456

While RandomForestRegressor is the best among these options, none of the models are currently useful due to negative R2

Data Presentation

The description of the dashboard

Screenshot of dashboard data description part:



The dashboard is describing data we used for our work. It shows the datatypes of the columns, key statistics, general info and provides samples for both tables.

Chart Descriptions & Findings

1. Rating Distribution

Query: sql/q1.hql

Finding: Strong bias toward positive ratings (4–5 stars); 2-star ratings are rarest. Suggests users either love books or avoid reviewing disliked ones.
2. Top Reviewed Books

Query: sql/q2.hql

Finding: A few books dominate review volume (e.g., bestsellers). Indicates "long-tail" distribution—most books have few reviews.
3. Review Trends Over Time

Query: sql/q3.hql

Finding: Slow growth until 1990, then rapid rise post-2000 (Amazon’s expansion). Peaks likely tied to holiday seasons and e-book adoption.
4. Publisher Performance

Query: sql/q4.hql

Finding: Smaller publishers often outperform giants

5. Most Active Users

Query: sql/q5.hql

Finding: A small group of users contributes disproportionate reviews—potential "super-reviewers" or bots.

6. Helpfulness Analysis

Query: sql/q6.hql

Finding: Most reviews are deemed helpful (70%+ votes positive), but mid-range ratings (3 stars) receive more engagement.

Key Insights & Recommendations

Rating Bias: Address skew in recommendations (weight low-rated reviews more).

Temporal Signals: Leverage seasonal spikes for marketing (holiday book promotions).

Publisher Trends: Partner with high-performing indie publishers for curated selections.

User Segmentation: Identify and incentivize super-reviewers to maintain engagement.

Helpfulness Paradox: Mid-range reviews are most useful.

Conclusion

Through meticulous data processing, feature engineering, and machine learning modeling, our team successfully:

Built a Predictive Rating System

Developed four regression models (RandomForestRegressor, LogisticRegression, GBTRRegressor, FMRegressor) to predict book ratings (1–5 stars), but all the models perform poorly (negative R2 values indicate they are worse than predicting the mean).

Processed Large-Scale Data

Cleaned and analyzed 3+ million reviews and 212K books from Amazon, handling challenges like NULL values, text normalization, and temporal feature encoding. Engineered hybrid features (TF-IDF for text, cyclical time encoding, Wilson scores for helpfulness).

Delivered Actionable Insights

Discovered rating biases (users prefer 4–5 stars) and temporal trends (post-2000 review surge).

Highlighted niche publishers with high average ratings and "super-reviewer" users driving engagement.

Established a Scalable Pipeline

Implemented a 3-stage workflow (Postgres → Hive → PySpark) for reproducible data processing and model training.

Optimized resource usage via sampling (10% of data) and distributed computing (Spark on Hadoop).

Reflections on Work

Completing this project was both challenging and rewarding. Working with big data and complex dependencies required careful planning and problem-solving. Below are some key reflections on the difficulties faced and lessons learned.

Challenges and Difficulties

1. Interdependent Components
 - 1.1. The project had multiple stages (data cleaning, feature engineering, analysis), where each step depended on the previous one.
 - 1.2. A small mistake in early preprocessing (e.g., incorrect joins) caused errors in later stages, requiring backtracking.
2. Cluster Performance Issues
 - 2.1. The hadoop cluster sometimes slowed down or crashed due to resource limits.
 - 2.2. Long wait times for job execution, especially with large aggregations.
3. Big Data latency

Even minor changes required reprocessing millions of rows (3.5 millions in our case :)), leading to long delays.
4. Learning New Technologies

PySpark, Hive, and distributed computing were new for most of us.

Recommendations

Improve Cluster Management

Monitor resource usage closely to avoid bottlenecks.

Knowledge Sharing

Make more communication inside the team to reuse already known knowledge.

The table of contributions of each team member

Project Tasks	Task description	Arsenii Pavlov	Ivan Beltsov	Andrei Pavlov	Timofei Ivlenkov	Bulat Latypov	deliverables	Average hours spent
Data extraction	Collect the data from the <link> and extract the representative sample for the project	0%	0%	0%	100%	0%	sample_data.csv	1
Data processing and clearing	Using sample_data.csv process it and upload to postgres	0%	0%	0%	100%	0%	postgres table	7
Uploading data to hdfs	Upload data to hdfs	0%	0%	0%	100%	0%	data in hdfs	2
Data analysis	Analyse given on postgres data	0%	0%	100%	0%	0%	Base knowledge about data for future work	5
Build tables for charts	Build tables which contains useful knowledge	20%	0%	80%	0%	0%	Tables with data for chart	4
Building charts	Build chart about data	0%	0%	100%	0%	0%	6 Charts in superset	2
Feature extraction and data	Process data from hive so it will be applicable	0%	75%	0%	0%	25%	test and train datasets	12

preprocessing	for model training						in json formats	
Model training	Train models	0%	40%	0%	10%	50%	trained models	20
Testing	Test every part and stage of project	80%	0%	0%	20%	0%	tested project	3
Report writing	Write report on project	80%	10%	10%	0%	0%	Written report	10
Making presentation	Prepare pdf version of presentation	75%	5%	0%	10%	10%	PDF version of presentation	3
Rest Small tasks	All the small tasks that was needed to build final project	20%	20%	20%	20%	20%	Result of the project	40

This is only a subjective approximation. From the top it can be stated that every team member has done his part and actively participated in the project.