

GameSphere: A Smart Gaming Analytics & Recommendation Application

Devdatta Gandole, Nakshatra Sachinbhai Desai, Rajarajeswari Premanand, Tejashree Kishor Badve, Vineeth Rayadurgam

San Jose State University

Course: Data 226 - Big Data Technologies and Applications

Abstract—GameSphere is a web-based platform that utilizes data to offer tailored video game suggestions and visual analytics for users on Steam. It employs PySpark, HDFS, and machine learning techniques to handle and examine an 8GB database of user reviews, facilitating actionable insights through sentiment analysis and collaborative filtering. The solution features an interactive frontend built with Streamlit to visualize trends and recommendations.

I. INTRODUCTION

The rapid expansion of user-generated content on gaming platforms offers both advantages and challenges for developers and researchers aiming to extract valuable insights at scale. Steam, one of the premier digital distribution platforms for PC games, receives millions of user reviews and gameplay logs on a daily basis. The analysis of this high-volume, fast-moving data requires a Big Data strategy capable of efficiently ingesting, processing, and modeling both streaming and batch data.

In this study, we introduce **GameSphere**, a web-based gaming analytics and recommendation system that utilizes a scalable pipeline incorporating Hadoop Distributed File System (HDFS) for storage and Apache Spark for distributed computing. The comprehensive workflow is illustrated in a detailed flow diagram (refer to FlowDetails.docx), which outlines the stages of raw data ingestion, preprocessing, exploratory data analysis (EDA), sentiment extraction, collaborative filtering, and result visualization: contentReference[oacite:0]index=0. The repository's organization—detailed in FolderStructure.rtf—ensures a clear distinction between raw and processed data, analysis notebooks, source code, and deployment scripts.

Our exploratory analysis phase ('EDA.ipynb') reveals significant patterns in review length, distribution of ratings, and temporal trends, which informs subsequent modeling choices. In 'sentiment analysis.ipynb', we derive sentiment-aware features using natural language processing techniques that transform free-text reviews into polarity scores, enhancing the recommendation engine. The fundamental recommendation logic, implemented through alternating least squares in 'game recommendation.ipynb', is compiled as a deployable Python module ('GameSphere.py'), along with a Streamlit frontend. All dependencies and environment settings are documented in 'requirements.txt' to ensure reproducibility.

This paper presents three primary contributions:

- A modular, fault-tolerant Big Data pipeline for analyzing large-scale gaming reviews.
- Integration of sentiment analysis with collaborative filtering to enhance recommendation relevance.
- An interactive dashboard that allows end users to explore both aggregate trends and personalized recommendations in real time.

II. DATASET AND INITIAL PREPROCESSING

Dataset: Steam Reviews 2021 from Kaggle

- **Size:** Approximately 8 GB
- **Entries:** Around 17 million reviews

Raw Schema Includes:

- Game identifiers: app_id, app_name
- User identifiers: author_steamid, author_num_games_owned, author_playtime_forever
- Review content: review, recommended, votes_helpful, timestamp_created, language

Preprocessing Pipeline:

1) Directory Setup in HDFS:

Structured directories such as /data/raw, /data/processed, /models, and /outputs were created to organize different stages of the data pipeline. This provided scalable and fault-tolerant storage for processing at various levels.

2) Schema Standardization:

Dot-separated column names (e.g., author.steamid) were converted to underscore format (e.g., author_steamid) to support SQL queries and PySpark compatibility. Entries missing essential fields were dropped to ensure data consistency.

3) Text Normalization:

Review text was cleaned by converting to lowercase, stripping HTML, emojis, and symbols, and removing entries with only numeric or junk content.

4) Data Type Conversion:

Fields were cast to types compatible with PySpark ML:

- Integers: app_id, author_num_games_owned

- **Floats:** `author_playtime_forever,`
`author_playtime_at_review`
- **Boolean:** `recommended`
- **Timestamps:** Converted from Unix format using `from_unixtime()` in Spark SQL

5) Language Filtering:

To ensure compatibility with the sentiment analysis model, only reviews where `language = 'english'` were selected. The filtered dataset was saved to `/processed/steam_review_english.parquet`.

Technology Role:

PySpark facilitated highly parallelized transformation using DataFrames, enabling millions of rows to be processed efficiently in memory. Its SQL support and functional chaining allowed seamless integration of schema transformations, filters, and column renaming steps, reducing ETL complexity and time to delivery.

III. TOOLS AND TECHNOLOGIES

To build GameSphere, a combination of big data tools, machine learning frameworks, and frontend visualization technologies were used. Each component played a specific role in the ETL pipeline, modeling, and dashboard interaction.

- **HDFS:** Distributed file system used to store raw input files, intermediate processed datasets, model outputs, and final recommendations. It ensures scalable and fault-tolerant storage across nodes.
- **PySpark:** Used for data preprocessing, text cleaning, feature extraction, and model training. PySpark enables parallel computation over large datasets and integrates well with HDFS and MLlib.
- **Pandas, Seaborn, Matplotlib:** These libraries were used for exploratory data analysis and generating static visualizations during the review behavior analysis stage.
- **HuggingFace Transformers (DistilBERT):** A lightweight transformer model used to perform binary sentiment classification (positive/negative) on user reviews.
- **Spark MLlib (ALS):** Spark's machine learning library was used to implement collaborative filtering via the Alternating Least Squares (ALS) algorithm for generating personalized game recommendations.
- **Streamlit:** Python-based open-source framework used to build the interactive dashboard for users to input Steam IDs, view recommendations, and explore game trends visually.

IV. METHODOLOGY

The GameSphere initiative employed a systematic approach from conception to implementation:

A. Problem Definition

Identified the issue of deriving valuable insights from millions of user-generated reviews on Steam. Objectives were set to offer game suggestions, analyze sentiment, and visualize user interactions.

B. Data Collection

The Steam Reviews 2021 dataset, containing approximately 17 million entries, was obtained from Kaggle. The dataset was uploaded and organized within HDFS to allow for scalable and concurrent data access.

C. Data Cleaning and Preprocessing

The schema was standardized, and irrelevant or noisy data entries were removed. Text content in reviews was normalized and refined for NLP compatibility. Only English reviews were retained, and the cleaned dataset was saved in Parquet format for efficient downstream use.

D. Exploratory Data Analysis (EDA)

User behavior, review trends, and language distribution were visualized. Data biases such as skewed playtime and review frequency were identified and corrected through adjusted filters and aggregations.

E. Sentiment Modeling

The DistilBERT model from HuggingFace was employed to classify reviews as positive or negative. The resulting sentiment scores were merged back into the original dataset to enrich downstream recommendation modeling.

F. Recommendation System

Collaborative filtering was implemented using the Alternating Least Squares (ALS) algorithm on the filtered review ratings. Hyperparameters such as rank and regularization were tuned, and model accuracy was evaluated using RMSE on a test split.

G. Dashboard Development

An interactive frontend was created using Streamlit. The interface accepts Steam IDs to generate personalized game recommendations and enables users to explore sentiment and popularity trends.

H. Testing and Results

The recommendations were cross-verified with popular game titles. Sentiment trends were analyzed in correlation with user playtime and review activity levels to ensure consistency and relevance.

I. Finalization

Models and outputs were stored in HDFS for accessibility. All backend processes were linked to the Streamlit dashboard to enable real-time insight delivery and a seamless user experience.

This organized methodology ensured that every stage—from raw data ingestion to insight visualization—was based on reproducible, scalable practices that leveraged contemporary big data and machine learning technologies.

V. EXPLORATORY DATA ANALYSIS (EDA)

The Exploratory Data Analysis (EDA) stage lays the groundwork for all future machine learning models and the development of dashboards. This phase includes statistical analysis, visual pattern identification, and data quality checks to uncover trends and anomalies.

A. Game Trends

The total number of reviews is calculated by `app_id` and `app_name` to highlight the top games with the greatest review count. The average recommendation score is utilized to rank the games that receive the most favorable feedback. Bar charts and horizontal graphs are employed to display the frequency of reviews and the proportion of positive sentiments for each game.

B. Temporal Patterns

The timestamp column (`timestamp_created`) is converted into user-friendly date formats and summarized by month and year. Review activities are illustrated as a time series to detect seasonal trends and spikes in user interaction, typically coinciding with major game launches or sales events. Temporal aggregation is further expanded to analyze the evolution of sentiment over time.

C. Language Distribution

The dataset contains reviews in various languages; a count of the `language` column indicates that over 80% of the entries are in English. This validation supports the decision to filter for English and guides the design of the sentiment model.

D. User Analysis

A histogram is created to evaluate the distribution of the number of games owned by each user (`author_num_games_owned`), helping to comprehend user diversity and spending habits. Violin and box plots for `author_playtime_forever` display skewed distributions, revealing that many users play for short periods while a few engage in extended play. Playtime information also assists in refining criteria for removing bot-like or anomalous behaviors.

E. Engagement Metrics

Additional metrics such as `votes_helpful` and review length are examined to assess the quality and impact of user reviews. Games with the highest median helpfulness in reviews and the longest review lengths are identified for potential promotion on the dashboard.

These insights are generated using Pandas for structured filtering and aggregation, while Seaborn and Matplotlib are utilized for clear and informative visualizations. The outcomes are stored in intermediate Parquet files for quick retrieval in the Streamlit dashboard. The cleaned datasets and visualizations are preserved for final presentation and analysis.

VI. SENTIMENT ANALYSIS PIPELINE

The sentiment analysis stage enriches user-generated reviews with emotional context, enabling a more personalized and insightful recommendation experience. This phase classifies reviews as **POSITIVE** or **NEGATIVE** using a pretrained transformer model and calculates confidence scores that indicate the model's certainty.

Model Used

`distilbert-base-uncased-finetuned-sst-2-english`

1. Step-by-Step Process

1) Loading Preprocessed Data:

English-language reviews, pre-filtered and stored as Parquet files, are loaded using PySpark DataFrames for distributed processing.

2) Review Classification:

Each review is passed through the HuggingFace DistilBERT pipeline, which tokenizes the text and returns a sentiment label (POSITIVE or NEGATIVE) along with a confidence score. This stage is optimized for batch processing and supports GPU acceleration.

3) Schema of the Output:

The enriched DataFrame includes:

- `review_id`
- `app_id`
- `app_name`
- `sentiment_label`
- `sentiment_score`

4) Game-Level Aggregation:

Using PySpark's `groupBy()` function, reviews are aggregated by `app_id` to calculate:

- Percentage of positive reviews
- Average sentiment score per game

5) Merging and Storing Output:

Aggregated sentiment scores are joined with the original cleaned dataset. The final output is written to HDFS:

`/processed/cleaned_steam_reviews_merged.parquet`

2. Contribution of Technology

- **HuggingFace DistilBERT:** A compact, fine-tuned transformer model optimized for sentiment classification. It balances high accuracy with low latency, suitable for large-scale textual analysis.
- **Spark DataFrames:** Enabled distributed and parallelized review processing. Aggregation and joins on massive datasets were handled efficiently using PySpark APIs.

This sentiment analysis pipeline bridges raw text with structured insights, improving the interpretability and personalization of game recommendations.

VII. GAME RECOMMENDATION SYSTEM

The game recommendation system employs collaborative filtering with the Alternating Least Squares (ALS) technique. The comprehensive process builds on each previous step to yield high-quality, individualized recommendations at scale.

1. Data Preparation

- The `recommended` column is transformed into a numeric rating: 5 signifies a positive recommendation, while 1 indicates otherwise. This standardization is essential for the ALS algorithm's input.
- Games that have received fewer than 200 reviews are eliminated to minimize noise and guarantee that the model learns from data that holds statistical significance.
- Users with abnormally low or high playtimes are excluded to filter out anomalies and potential bot activity, resulting in a more representative user base.

2. Feature Engineering

- User and game identifiers (e.g., `author_steamid`, `app_id`) are converted into numerical indices using `StringIndexer`. These indices are necessary for the ALS algorithm's matrix operations.

3. Training ALS

- The ALS model is set up utilizing the indexed fields: `userCol=author_index`, `itemCol=app_index`, and `ratingCol=rating`.
- Hyperparameters such as `rank` (number of latent features), `regParam` (regularization), and `maxIter` (number of iterations) are fine-tuned to enhance prediction accuracy.
- The dataset is divided into training (80%) and testing (20%) portions. This division ensures that the model is evaluated using data it has not previously encountered.

4. Evaluation

- The model's performance is measured using Root Mean Squared Error (RMSE). A lower RMSE reflects superior predictive ability.
- This evaluation is instrumental in detecting overfitting and verifying the selected hyperparameters.

5. Output

- Following the training phase, the model produces the top 5 game recommendations for every user.
- These recommendations are stored at the following HDFS location for future use:

`/outputs/user_recommendations.parquet`

6. Technological Advantage

- The ALS implementation in Spark MLlib takes advantage of distributed computing, enabling scalable matrix factorization across multiple nodes.
- This arrangement allows for efficient model training involving millions of users and games, greatly cutting down on computational time and memory requirements.

Each phase of this pipeline builds upon the previous one, refining data and optimizing model inputs to create a quick, precise, and scalable recommendation system that drives the GameSphere dashboard. This setup efficiently manages millions of user-game interactions, significantly speeding up both training and prediction processes across extensive clusters.

VIII. STREAMLIT DASHBOARD

Streamlit is a free, open-source Python library that facilitates the development of interactive web applications for data science initiatives with minimal coding effort. It empowers developers and analysts to convert data scripts into shareable web applications in just a few minutes by solely using Python.

Streamlit was selected for this initiative due to its user-friendliness, rapid development process, and seamless integration with Python-driven data workflows. Its reactive structure automatically refreshes the frontend in response to changes in the backend, making it ideal for real-time dashboards.

In **GameSphere**, Streamlit serves as the user interface that connects individuals to insights derived from extensive data analysis. It interfaces with cleaned datasets and model outputs to deliver a smooth experience through an easy-to-navigate layout.

Technical Implementation

- Streamlit scripts were written in Python and deployed either locally or via web hosting.
- Data was accessed from Parquet files generated by the PySpark processing pipeline.
- Caching techniques such as `@st.cache_data` were used to avoid unnecessary reloading of datasets.
- Charts were built using Streamlit's built-in tools and external libraries like Plotly and Matplotlib.
- The layout was structured with columns and tabs to separate visual analytics from personalized recommendations.

Tabs in the Dashboard

- **Game Recommendations**
 - **Input:** Steam ID
 - **Output:** Top 5 recommended games, predicted ratings, sentiment score overlays

- **Analytics Dashboard**

- **Filters:** Game selection and date range (monthly)
- **Visuals:** Pie charts, bar graphs, and time-series trendlines for reviews, sentiments, and recommendation frequency

Integration

- Accesses cleaned data from Parquet files
- Employs caching to enhance performance and minimize redundant I/O operations

Streamlit's ease of use enabled quick prototyping and iterative enhancement of user experience while maintaining full compatibility with the backend machine learning workflows.

IX. CONCLUSION AND FUTURE SCOPE

GameSphere showcases the effective integration of big data processing, sentiment analysis, and collaborative filtering to generate personalized, insightful gaming recommendations. The platform processes millions of user reviews using PySpark and HDFS, enriches them with sentiment via DistilBERT, and delivers real-time visualizations through Streamlit.

The structured pipeline—from raw data ingestion to front-end deployment—demonstrates a scalable, modular design that supports both analytical and predictive use cases.

Future Enhancements

- **Multilingual Support:** Expand sentiment analysis to support reviews in languages beyond English.
- **User-Level Insights:** Enable tracking of personalized sentiment and behavioral patterns.
- **Real-Time Processing:** Integrate Kafka and Spark Streaming for live review ingestion and dashboard updates.
- **Game Metadata Fusion:** Incorporate game genre, release dates, and pricing for more context-aware recommendations.
- **Mobile Optimization:** Adapt dashboard UI for mobile platforms to enhance accessibility.

Overall, GameSphere illustrates the powerful synergy of data engineering and machine learning in transforming large-scale unstructured data into actionable insights for the gaming ecosystem.

X. CONCLUSION

GameSphere successfully integrates big data analytics, machine learning, and web technologies to deliver an innovative platform for gaming insights and personalized recommendations, demonstrating practical use of modern big data technologies.

REFERENCES

- **Kaggle:** Steam Reviews 2021
- **HuggingFace Transformers**
- **Apache Spark MLlib**
- **Streamlit Documentation**
- **GitHub Repo:** https://github.com/tejashreebadve/DATA228_Group11