



Scalable Recommender System Pipeline

Design, Implementation, and Evaluation

Problem & Motivation



Large-scale Personalization

Modern digital platforms require highly personalized experiences, creating a growing demand for systems that understand user preferences and deliver context-aware recommendations at scale.



Massive Data Challenges

Recommender systems operate on vast datasets, including user interactions and item attributes. Managing this data demands robust infrastructure and efficient algorithms for scalability and responsiveness.



Data Sparsity & Cold-Start

New users or items face a cold-start problem due to insufficient interaction data. Data sparsity in preference matrices also makes accurate predictions challenging for conventional systems.



Single-Machine Limitations

Traditional, single-machine recommendation approaches often fail to keep up with the scale and complexity of modern data, leading to performance bottlenecks and high latency.

"Modern recommender systems must handle huge datasets with low latency and acceptable accuracy to meet today's personalized content demands."

Why Big Data Frameworks?

Apache Spark

A unified engine for large-scale data analytics, known for its **in-memory processing** and **ML capabilities**.

- In-memory Processing
- Machine Learning
- Data Science

Hadoop / HDFS

A robust framework for **distributed storage** and processing of massive datasets across clusters.

- Distributed Storage
- Fault-Tolerant
- Scalable Data Handling

Apache Parquet

An efficient **columnar storage format** designed for high-performance analytical queries and **I/O**.

- Columnar Format
- I/O Efficiency
- Optimized Queries



Spark + Parquet: Significantly reduce training time and I/O bottlenecks in large-scale pipelines.

Recommender Approaches

Collaborative Filtering (CF)



Groups users by similar behaviors to recommend new items. Finds users with similar preferences. **Powerful but suffers from sparsity.**

- Leverages user-item interactions
- Recommendations based on 'users like you'
- Core in systems like Netflix

Content-Based Filtering (CBF)



Recommendations based on item attributes and user's past preferences for those attributes. **Uses item information.**

- Focuses on item features (e.g., genre)
- Recommends based on 'items you liked'

Hybrid Systems



Combines CF and CBF to address individual limitations (like sparsity) and **enhance accuracy.**

- Mitigates sparsity and cold-start
- Enhances overall prediction accuracy
- Essential for real-world systems

Algorithms Compared



ALS

(Matrix Factorization)

Predicts user preferences in recommendation systems by factorizing user-item interaction matrices.



Linear Regression

(Baseline Model)

Foundational supervised model predicting continuous values, serving as a simple yet effective baseline.



Decision Tree

(Tree-based Model)

Non-linear supervised method that splits data based on features for classification or regression.



Random Forest

(Ensemble Method)

Builds multiple decision trees, combining outputs to improve accuracy and reduce overfitting.



GBT

(Ensemble Method)

Sequentially builds trees, correcting errors of previous ones to iteratively optimize performance.



Key Takeaway

*These models represent both **collaborative filtering** and **supervised learning** approaches.*

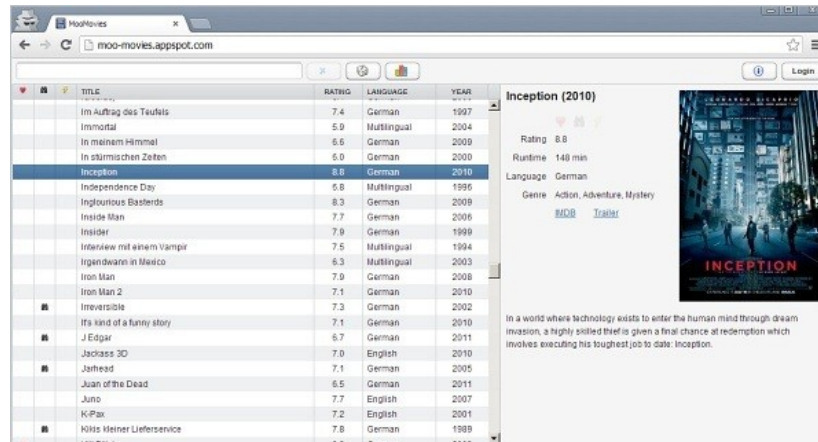
Dataset & Metrics

MovieLens Dataset Characteristics

The MovieLens dataset, collected by GroupLens Research, is a widely used benchmark for recommender systems.

It provides real-world feedback including user ratings, movie metadata, and tags, critical for research.

Millions of users, items, and ratings



	TITLE	RATING	LANGUAGE	YEAR
	Im Auftrag des Teufels	7.4	German	1997
	Immortal	5.9	Multilingual	2004
	In meinem Himmel	6.6	German	2009
	In stürmischen Zeiten	6.0	German	2000
	Inception	8.8	German	2010
	Independence Day	6.8	Multilingual	1996
	Ingolourous Bastards	8.3	German	2009
	Inside Man	7.7	German	2006
	Insider	7.9	German	1999
	Interview mit einem Vampir	7.5	Multilingual	1994
	Irren dürfen in Mexico	6.3	Multilingual	2003
	Iron Man	7.9	German	2008
	Iron Man 2	7.1	German	2010
	Inevitable	7.3	German	2002
	It's kind of a funny story	7.1	German	2010
	J Edgar	6.7	German	2011
	Jackass 3D	7.0	English	2010
	Jarhead	7.1	German	2005
	Juan of the Dead	6.5	German	2011
	Juno	7.7	English	2007
	K-Pax	7.2	English	2001
	Kills Kleiner Lieferservice	7.8	German	1989

Inception (2010)


Rating: 8.8

Runtime: 148 min

Language: German

Genre: Action, Adventure, Mystery

[Watch](#) [Trailer](#)



In a world where technology exists to enter the human mind through dream invasion, a highly skilled thief is given a final chance at redemption which involves executing his toughest job to date: Inception.

Key Evaluation Metrics

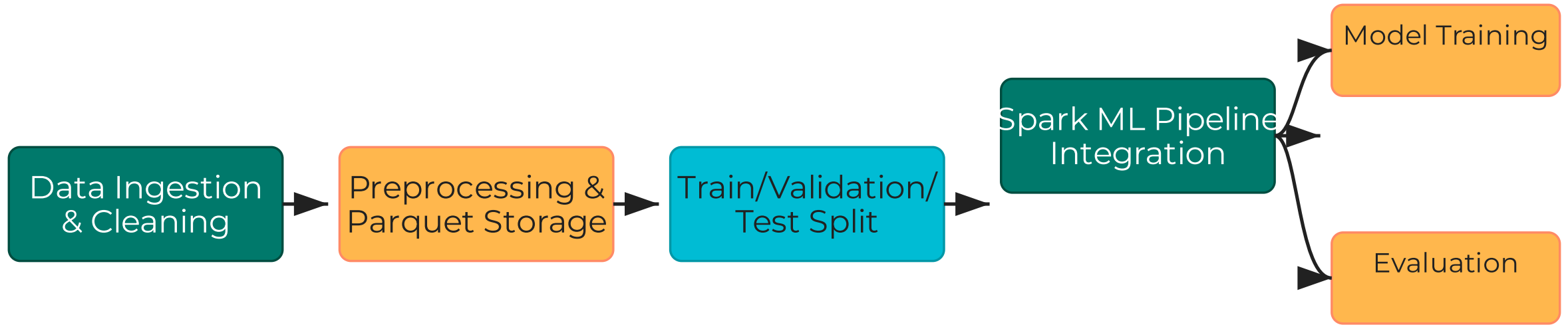
Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) are common accuracy metrics.

They measure how closely predicted user ratings align with actual ratings, quantifying prediction error.

RMSE is the standard metric for rating prediction, prominently used in the Netflix Prize and industry benchmarks.



System Pipeline



This end-to-end pipeline is fully distributed and optimized for scalability.

Experimental Results

RMSE Performance Comparison



Key Performance Insights

ALS: Exhibits superior RMSE performance for accurate recommendations. ✅

Random Forest & GBT: Demonstrate moderate performance, offering a balanced approach. ⚖️

Linear Regression: Registers the highest error rates, indicating less suitability. ⚠️

ALS outperforms supervised models by directly modeling latent user-item interactions, leading to more nuanced and effective recommendations.

Key Insights & Limitations

Matrix Factorization Effectiveness

Effectively captures latent features by decomposing user-item interaction matrices into lower-dimensionality components for collaborative filtering.

Latent Features

Collaborative Filtering

Cold-Start Problem

New users or items lack sufficient interaction data, leading to inaccurate or unavailable recommendations.

Absence of Ranking Metrics

Evaluations often overlook ranking-based metrics, which are critical for assessing recommendation order and presentation.

Batch-Only Processing

System primarily operates in batch mode, limiting real-time adaptability to changing user preferences or item dynamics.

Spark for Scalable Offline Training

Apache Spark enables scalable offline training through distributed processing and in-memory capabilities for efficient ML model development on big data.

Distributed Processing

In-Memory Processing

Accuracy alone is not sufficient; production recommender systems demand robust solutions addressing **cold-start**, **ranking**, and **real-time** challenges.

Conclusion & Future Work

Conclusion

- ▮ ALS as a **strong baseline**
- ▮ Scalable & reproducible **big-data pipeline**

Future Work

- 🔬 **Hybrid models** with metadata
- ★ **Implicit feedback** & ranking metrics
- ✳️ **Online/streaming** recommendations
- 🧠 **Deep learning** & graph models

Real-world recommender systems require hybrid, online, and multi-objective optimization.